

John K. Kruschke



Doing Bayesian Data Analysis



A Tutorial with R and BUGS



parameter values are negatively correlated, such that when we believe one parameter value is large, we believe that the other parameter value is small. negative correlation causes the conjoint distribution to straddle the line of equality.

In summary, the marginal distributions of single parameters do not indicate the relationship between the parameter values. The conjoint distribution of the two parameters might have positive or negative correlation (or even a non-linear dependency), and therefore the difference of the parameter values should be explicitly examined.

12.1.3 Region of Practical Equivalence (ROPE)

The estimation approach can be enhanced by including a region of practical equivalence (ROPE), which indicates a small range of values that are considered to be practically equivalent to the null value for purposes of the particular application. For example, if we wonder whether a coin is fair, for purposes of determining which team will kick off at a football game, then we want to know if the underlying bias in the coin is reasonably close to 0.50, and we don't really care if the true bias is 0.51 or 0.49, because those values are close enough for our application. As another example, if we are assessing the efficacy of a drug versus a placebo, we might only consider using the drug if it improves the probability of cure by at least three percentage points, and we would declare the drug to be counterproductive if it decreases the probability of cure at all, i.e., the lower boundary of the ROPE would be zero difference.

Once a ROPE is set, we make decisions according the following logic:

A parameter value is declared to be credible, or rejected, if its entire ROPE lies outside the 95% HDI of the posterior distribution of that parameter.

For example, suppose that we want to know whether a coin is fair, so we establish a ROPE that goes from .45 to .55. We flip the coin 500 times and observe 255 heads. If the prior is uniform, the posterior has a 95% HDI from .608 to .691, which is completely outside the ROPE. Therefore we declare that the null value of 0.5 is rejected for practical purposes.

Because the ROPE and HDI can overlap in different ways, there are different decisions that can be made. In particular, we can decide to “accept” or “reject” a value:

A parameter value is declared to be accepted for practical purposes if that value's ROPE completely contains the 95% HDI of the posterior of that parameter.

For example, suppose that we want to know whether a coin is fair, so we establish a ROPE that goes from .45 to .55. We flip the coin 1,000 times and observe 490 heads. If the prior is uniform, the posterior has a 95% HDI from .459 to .521, which is completely within the ROPE. Therefore we declare that the null value of 0.5 is accepted for practical purposes, because all of the credible values are practically equivalent to the null value.

This use of a ROPE around the null value also implies that if the null value really is true, we will eventually “accept” the null value as the sample size gets large enough. This is because, as the sample size gets larger, the HDI tends to be narrower and closer to the true value. When the sample size gets very large, the HDI is certain to be narrow enough, and close enough to the true value, to fall entirely within the ROPE. If we did not use a ROPE around the null value, and rejected the null every time that it falls outside the 95% HDI, then we would incorrectly reject the null 5% of experiments even if the null value is true. This 5% false alarm rate occurs because random samples of data,

when generated by the null value, will have coincidences ~~suffers~~, even for large N , that produce a 95% HDI that does not include the null value. Although 5% of those 95% HDIs exclude the null value, the HDIs do get narrower and closer to the null value as N gets large. Therefore the HDIs do eventually fall within the ROPE around the null, even though 5% of the 95% HDIs exclude the null value itself. As an example ~~of this~~ point: Suppose we establish a ROPE that goes from .45 to .55. We flip the coin 1000 times and observe 5200 heads. The 95% HDI goes from .510 to .530. Despite the fact ~~that~~ the HDI excludes the null value, we accept the null for practical purposes because the HDI falls entirely within the ROPE, which means that all of the credible values are practically equivalent to the null value. We would say that we believe the true bias in the coin ~~is very~~ nearly .52, but that's close enough to the null value of .50 for practical purposes.

How is the size of the ROPE determined? If the application domain in which costs and benefits of decisions can be determined, then the full theoretical machinery of expected utility can be brought to bear. In some domains ~~such as~~ medicine, expert clinicians can be interviewed, and their opinions can be translated into a reasonable consensus regarding how big of an effect is useful or important for the application. Otherwise, the ROPE might be established with somewhat arbitrary criteria, bearing in mind the key trade-off: When the ROPE is wider, there is a lower probability of ~~falsely~~ rejecting the null value (i.e., there is a smaller false alarm rate), but there is also a lower probability of rejecting the null value when it is false (i.e., there is a smaller hit). For further discussion of the ROPE, under somewhat different appellations of “range of equivalence” and “influence zone”, see e.g., Carlin and Louis (2009); Freedman, Lowe, Macaskill (1984); Hobbs and Carlin (2008); Spiegelhalter, Freedman, and Parmar (2019).

It is important to be clear that any discrete declaration ~~about~~ rejecting or accepting a null value does not exhaustively capture our beliefs about the parameter ~~so~~. Our beliefs about the parameter value are described by the full posterior distribution. When making a binary declaration, we have merely compressed all that detail into a single bit of information. The broader goal of Bayesian analysis is conveying ~~an informative summary of the posterior, and where the value of interest falls within that posterior.~~ Reporting the limits of an HDI region is more informative than reporting the declaration of a reject/accept decision. By reporting the HDI and other summary information about the posterior, different readers can apply different ROPEs to decide for themselves whether a parameter is practically equivalent to a null value.

12.2 The model-comparison (two-prior) approach

The previous section posed the question, is the null value ~~one of~~ the credible values, in terms of parameter estimation. We started with an informed prior distribution on the parameters, and then examined the posterior distribution ~~of the~~ parameters.

Some researchers prefer instead to pose the question ~~is the~~ model comparison. In this framing of the question, the focus is not on estimating the magnitude of the parameter. Instead, the focus is on deciding which of two hypothetical priors is least unbelievable. One hypothetical prior expresses the idea that the parameter is exactly the null value. The alternative hypothetical prior expresses the idea that the parameter could be anything (or anything but the null value). Notice that neither of these hypothetical priors is informed by prior knowledge. This lack of being informed is often taken as a desirable aspect of the approach, not a defect, because the method is therefore “automatic” insofar as it obviates

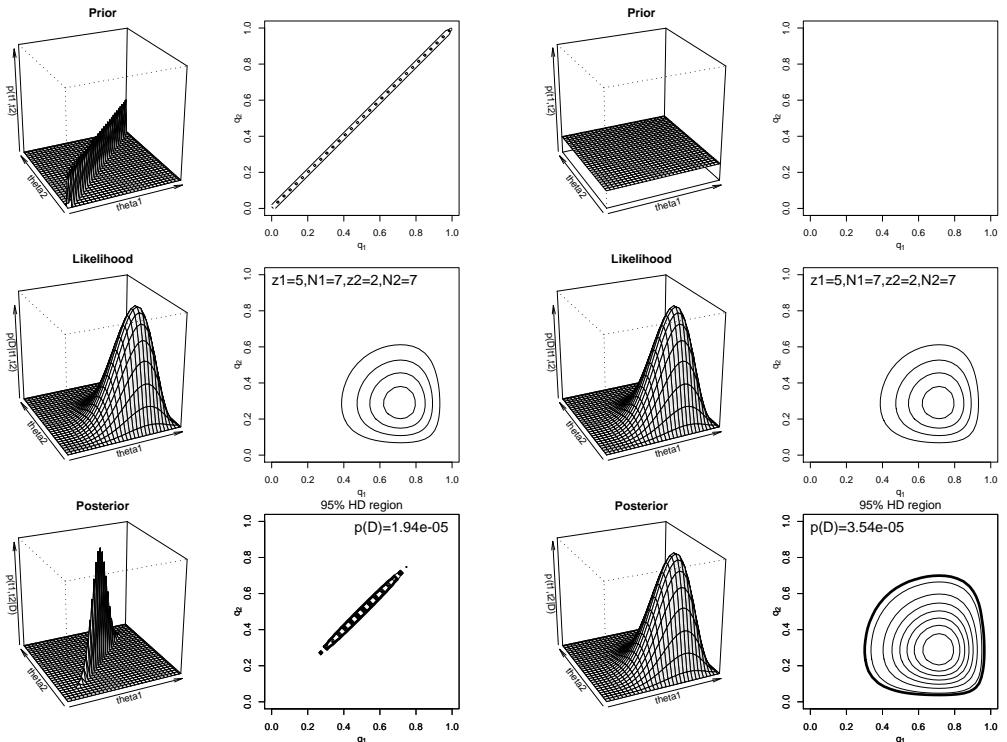


Figure 12.4: Two hypothetical priors for model comparison. The left columns show a ridge-shaped prior that allows only points for which $\theta_1 = \theta_2$. The right columns show a uniform prior, in which any combination of θ_1 , θ_2 values is entertained. The top-right contour plot is empty because the ~~priors~~ exactly horizontal, without any height changes to mark by a contour.

disputes about prior knowledge. We will see, however, that the model-comparison method can be extremely sensitive to the choice of “uninformed” prior for the alternative hypothesis (cf. Figure 10.7, p. 212), and that model comparison is necessarily meaningful unless both models are viable in the first place.

12.2.1 Are the biases of two coins equal or not?

Consider a situation in which we have two coins, and we want to infer whether their biases are equal or not. We pose the question as a model comparison such that one model expresses the “null” hypothesis that the two biases are equal, and the other model expresses the “alternative” hypothesis that the two biases could be different. The two models are distinguished by their priors (because the likelihood functions are the same in both models): The null model has zero prior probability on all combinations of biases except where they are equal, while the alternative model has uniform probability on all combinations of biases. The uniform alternative is chosen because it is a convenient expression of non-informed indifference.

Figure 12.4 shows an example of this approach. The “null” prior, $p(\theta_1, \theta_2 | \text{null})$, is shown in the left columns. Notice that this prior is shaped as a ridge along all values for which $\theta_1 = \theta_2$. The prior gives zero probability to any point at which, $\theta_1 \neq \theta_2$. Notice also that the height of the ridge is horizontal, meaning that it is not informed by any prior knowledge,

e.g., that the coins are probably fair. The “alternative” prior, $p(DjM_{alt})$, is shown in the right columns of Figure 12.4. This prior puts equal probability on all possible combinations of theta values. Notice that this hypothesis also is not informed by any prior knowledge, e.g., that the coins are probably fair.

The lower row of Figure 12.4 shows the posteriors resulting from two hypothetical priors, for the same data that have been used in previous figures such as Figure 8.1, p. 131. Notice that the posterior on the null model is a ridge in the prior, except that the extreme ends of the ridge have been attenuated because they are inconsistent with the data.

To do model comparison, we use the implication of Bayes rule:

$$\frac{p(M_{alt}|D)}{p(M_{null}|D)} = \frac{p(D|M_{alt})}{\left| \frac{p(D|M_{null})}{Z} \right|} \frac{p(M_{alt})}{p(M_{null})} \quad (12.1)$$

BF

where the ratio marked “BF” is the Bayes factor of the alternative model relative to the null model. The prior beliefs in each model are typically assumed equal, i.e. $p(M_{null}) = p(M_{alt}) = 0.5$, in the spirit of non-informed priors. The posterior plots in Figure 12.4 display the values of $p(D|M)$ for each model. The evidence for the null model $p(D|M_{null}) = 1:94 \times 10^{-5}$, while the evidence for the alternative model $p(D|M_{alt}) = 3:54 \times 10^{-5}$. The Bayes factor therefore slightly favors the alternative model but not by much. Because the ratio of posterior probabilities is not very extreme, we would conclude that either model remains reasonably credible, given the data. If the Bayes factor had turned out to be more extreme, we might decide to declare one or the other prior less unbelievable than the other prior.

12.2.1.1 Formal analytical solution

The Bayes factor for the null and alternative models can be computed analytically in this case. The alternative hypothesis has a uniform prior that is uniform because the alternative hypothesis in this approach is supposed to be a “conventional” prior that expresses a hypothesis complementary to the null hypothesis. A uniform prior on θ_1, θ_2 can be described as a product of beta distributions, namely $\text{beta}(\theta_1; 1, 1) \times \text{beta}(\theta_2; 1, 1)$. Therefore we can determine an exact value of $p(D|M_{alt})$ from Equation 8.5 (p. 131), which becomes

$$\begin{aligned} p(D|M_{alt}) &= \frac{B(z_1+1; N_1 - z_1+1) B(z_2+1; N_2 - z_2+1)}{B(1; 1)B(1; 1)} \\ &= B(z_1+1; N_1 - z_1+1) B(z_2+1; N_2 - z_2+1) \end{aligned} \quad (12.2)$$

because $B(1; 1) = 1$.

We can also derive a formal analytical expression for the evidence for the null hypothesis, $p(D|M_{null})$. First, notice this: Because the null hypothesis believes at any point $\theta_1 = \theta_2$ the ridge is zero, i.e. $p(\theta_1 = \theta_2) = 0$, the double integral for the evidence simplifies to a single integral over $\theta_1 = \theta_2$:

$$\begin{aligned} p(D|M_{null}) &= \int_Z d\theta_1 d\theta_2 p(D|\theta_1, \theta_2) p(\theta_1, \theta_2) \\ &= \int_Z d\theta_1 p(D|\theta_1) \end{aligned}$$

We now plug in the Bernoulli likelihood function to get

$$\begin{aligned}
 p(D_j M_{\text{null}}) &= \frac{d}{Z} p(D_j) \\
 &= \frac{d}{Z} (z_1)(1 -)^{(N_1 - z_1)} (z_2)(1 -)^{(N_2 - z_2)} \\
 &= \frac{d}{R} (z_1 + z_2)(1 -)^{(N_1 - z_1 + N_2 - z_2)} \\
 &= B(z_1 + z_2 + 1; N_1 - z_1 + N_2 - z_2 + 1)
 \end{aligned} \tag{12.3}$$

because, by definition $B(a; b) = \frac{d}{R} d^{(a-1)} (1 -)^{(b-1)}$.

The Bayes factor for the alternative hypothesis, relative to the null hypothesis, is then the ratio of Equations 12.2 and 12.3:

$$\frac{p(D_j M_{\text{alt}})}{p(D_j M_{\text{null}})} = \frac{B(z_1 + 1; N_1 - z_1 + 1) B(z_2 + 1; N_2 - z_2 + 1)}{B(z_1 + z_2 + 1; N_1 - z_1 + N_2 - z_2 + 1)} \tag{12.4}$$

As a concrete example, consider the case presented in Figure 12.4, where $z_1 = 5$, $N_1 = 7$, $z_2 = 2$, and $N_2 = 7$. Plugging those values into Equation 12.4 yields $p(D_j H_{\text{alt}}) = p(D_j H_{\text{null}}) = 1:82$. The grid approximation in Figure 12.4 displays $p(D_j H_{\text{alt}}) = 3:54 \times 10^{-5}$ and $p(D_j H_{\text{null}}) = 1:94 \times 10^{-5}$, yielding a ratio of 1.82, which is the same as the analytical solution.

12.2.1.2 Example application

Consider again the question of the “hot hand” in basketball. In Exercise 8.1, p. 153, which asked whether there are more successful shots after a preceding successful shot than after a preceding failed shot. For the case of one famous professional player, there were 251 successes after 285 initial successes, and 48 successes after 53 initial failures.

We can use model comparison to analyze this situation. Using Equation 12.4, we obtain a Bayes factor of 0.128, which, inverted, is a Bayes factor of 7.81 in favor of the null prior. If the priors on the models are 50, then the posterior probability of the null model is $p(M_{\text{null}}|D) = 88.7\%$, and the posterior probability of the alternative model $p(M_{\text{alt}}|D) = 11.3\%$.

The model comparison suggests that the null model is less believable than the alternative model. But should we conclude, therefore, that there is zero difference between the conditions? Maybe, but not necessarily. Prior belief says that it is extremely unlikely that there is literally zero correlation between the first and second shots of a pair of free throws in basketball. Therefore, even though the model comparison suggests that the null prior is less unbelievable than the alternative prior, what might really want is an estimate of the difference between success after success and success after failure. The model comparison does not provide such an estimate.

On the other hand, the direct estimation of the underlying effect did provide an estimate of the difference, as was shown in Figure 12.2, p. 242. There we saw that the measured difference of zero was among the credible differences, but we also saw that the measured difference was a bit larger than zero, and we also saw explicitly that there is still large uncertainty in the estimate of the difference. This conclusion from the direct estimate is a more informative inference from the data than the conclusion from model comparison, which states only that the ridge prior is less inconsistent with the data than the uniform prior. It should be remembered that the direct estimate of the difference began with a mildly informed prior, using

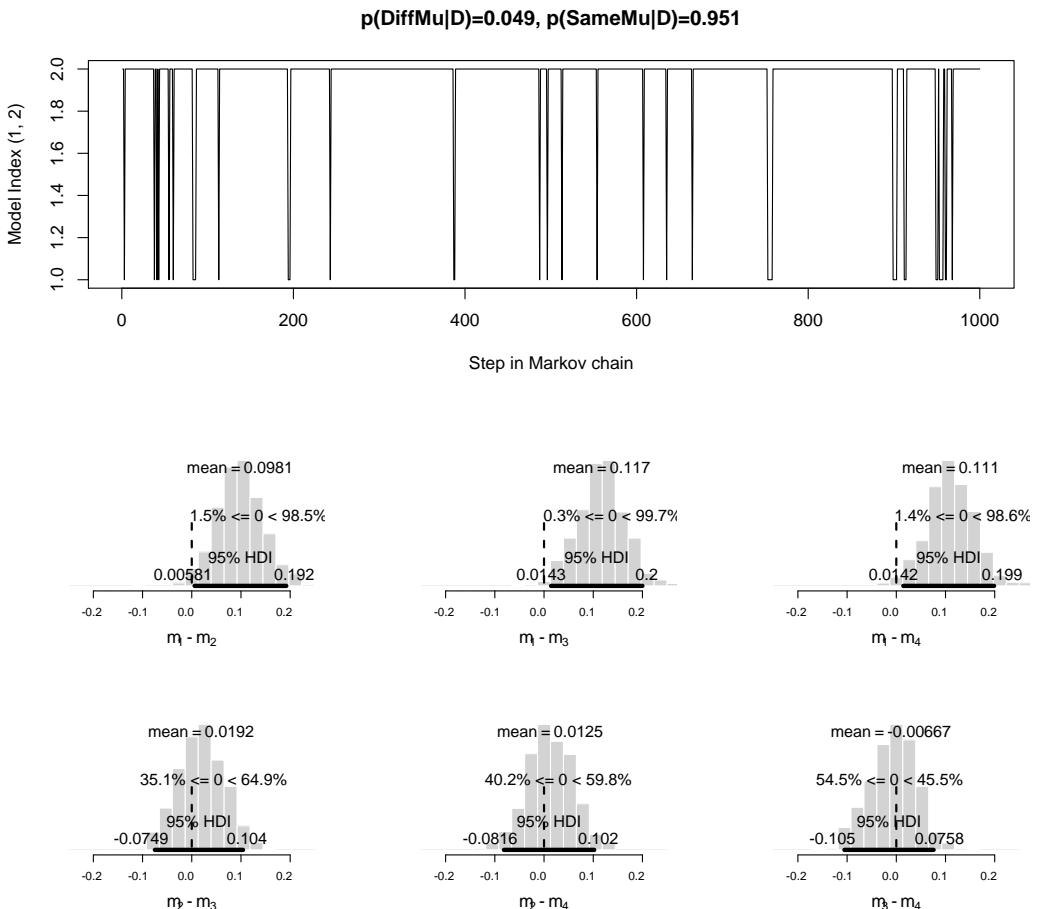


Figure 12.5: Upper Panel: Results of model comparison showing the model with a single mu parameter (“SameMu”) is preferred to a model with a separate mu parameter for each group (“DiffMu”). Histograms: Differences of posterior μ values for the four groups in the different-mu model. Notice that μ_1 is credibly different from μ_3 and μ_4 .

the knowledge that professional players tend to make about 50% of their free throws, unlike the alternative model’s uniform prior in the model comparison. Thus, the posterior from the alternative model is not the same as the posterior from the estimation.

12.2.2 Are different groups equal or not?

Suppose we conduct a study that has four conditions. Suppose every participant in every condition gets the same list of twenty words to try to memorize. The ability to recall a word is modeled as a Bernoulli distribution, with probability π_{ij} for the i^{th} person in the j^{th} condition. The individual recall propensity π_{ij} depends on hyperparameters α_j and β_j , that describe the overarching recall propensity in each condition because $\pi_{ij} \sim \text{beta}(\alpha_j + 1; (1 - \pi_{ij}) + 1)$. (Abstractly, this is a design just like the four category structures of the Iteration-condensation experiment of Section 9.3.1).

The only difference between the conditions is the type of music being played during learning and recall. For the four groups, the music *confounds*, respectively, the

death-metal band "Das Kruschke"¹Mozart, Bach, and Beethoven. For the four conditions, the mean number of words recalled is 11.85, 9.85, 9, and 9.60. The citations, randomly-generated data can be examined in detail by running the program in Section 12.4.1 (OneOddGroupModelComp).

We would like to know whether there is an effect of the type of music on ability to remember words. The most straight-forward way to find out is to estimate the parameters and then examine the posterior differences of the parameter estimates. The histograms in Figure 12.5 show the distributions of differences between the parameters. It can be seen that μ_1 is quite different than μ_3 and μ_4 ; a difference of zero falls well outside the 95% HDI intervals. From this we would conclude that Das Kruschke's memory better than Bach and Beethoven, but perhaps not better than Mozart.

A model-comparison approach addresses the issue in a different way. It compares the full model, which has distinct parameters for the four conditions, against a restricted model, which has a shared parameter to describe all the conditions simultaneously. If the two models have equal prior probabilities. The upper panel of Figure 12.5 shows the results of a model comparison, using transdimensional MCMC pseudopriors. The shared- μ_0 model is preferred, about 19 to 1 (0.951 to 0.049), over the full model. In other words, from this model comparison we would conclude that there is no difference in memory between the groups.

Which analysis should we believe? Is condition 1 different from conditions 3 and 4, or are all the conditions the same? Consider carefully what the model comparison actually says: Given the choice between one shared parameter and four group parameters, the one-parameter model is less unbelievable. But that does not mean that the one-parameter model is the best possible model. In fact, if a full model comparison is conducted, that compares the one-parameter model against a model that has one parameter for group 1 and a second parameter that is shared for groups 2 through 4, then the comparison favors the two-parameter model (see for yourself in Exercise 1).

In principle, we could consider all possible models formed by partitioning the four groups. For 4 groups, there are 15 distinct partitions. We can principle, put a prior belief on each of the 15 models, and then do a comparison of the models (Gopalan & Berry, 1998). From the posterior probabilities of the models we could ascertain which partition was most believable, and decide whether it is believable that other nearly-as-believable partitions. (Other approaches have been described by Berry & Hochberg, 1999; Mueller, Parmigiani, & Rice, 2007; Scott & Berger, 2006) Suppose that we conducted such a large-scale model comparison, and found that the most believable model partitioned groups 2–4 together, separate from group 1. Does this mean that we should truly believe that there is zero difference between groups 2, 3, and 4? Not necessarily. If the treatments are different, such as the four types of music in the present scenario, there is almost certainly at least some small difference between their outcomes. (In fact, the simulated data do come from groups with all different means.) We may still want to estimate the magnitude of those small differences. An explicit posterior estimate will reveal the range and uncertainty of those estimates. Thus, unless we have a reason to believe that different group parameters may be literally identical, an estimate of distinct group parameters will tell us what we want to know, without model comparison.

¹To find information regarding Das Kruschke, search www.archive.org/details/das_kruschke. The author has no relation to the band, other than, presumably, some unknown common ancestry many generations in the past. The author was, however, in a garage band as a teenager. That band didn't play death metal, although the music may have sounded that way to the critics.

12.3 Estimation or model comparison?

12.3.1 What is the probability that the null value is true?

There may be situations in which we want to know the probability that the null value is true. Proponents of the model-comparison approach may argue that this approach is especially applicable to these situations, because the model compares the posterior probability of the null hypothesis. Unfortunately, the posterior probability of the null is only a relative one, depending completely on the alternative hypothesis used in the model comparison. For example, as mentioned earlier for the fictitious memory experiment, a comparison of the null hypothesis (single-parameter) model against a two-parameter model yielded a strong preference in favor of the null hypothesis, but a comparison of the null hypothesis against a two-parameter model yielded a preference against the null hypothesis.

The posterior probability of the null hypothesis also depends crucially on the prior distribution assumed in the alternative model (e.g., Liu & Aitken 2008). Much of the effort in pursuing the model-comparison approach to null-hypothesis testing goes into justifying the form of the alternative-hypothesis prior (e.g., Edwards, Lindman, & Savage, 1963; Gallistel, 2009; Rouder, Speckman, Sun, Morey, & Iverson, 2009; Wagenmakers, 2007). Often the goal is to establish an “automatic” vague prior for the alternative model that has useful mathematical properties. The benefit of a conventionalized vague prior is that possible disputes regarding the best prior might be ignored by appeal to convention. Unfortunately, none of the automatic vague priors is usually the prior that is most appropriate for model comparison, which is an informed prior. A model comparison is useful if both models are genuinely viable, and viability is enhanced by incorporating prior information. For example, suppose a researcher is interested in investigating whether practitioners of therapeutic touch can detect the presence of another person's hand from a distance of a few centimeters (see Section 9.2.4, p. 177). The null hypothesis is that detection accuracy is at chance:

$\pi = 0.5$. Should the alternative hypothesis, which entertains chance values, be set at an “automatic” uniform distribution? Proponents of therapeutic touch would probably say no; instead, the prior distribution should incorporate knowledge about how well the practitioners might do in this unusual task. The practitioners might argue that accuracy would typically be only a little above chance, perhaps at 0.6, with fair uncertainty, as expressed perhaps in a $\text{beta}(6, 6)$ distribution.

The estimation approach, when combined with a ROPE, provides a direct measure of the probability of the null. This probability is simply the proportion of the posterior distribution that falls within the ROPE. This proportion of the HDI within the ROPE is not the probability that the null value is true; instead, this proportion is the probability that the parameter is practically equivalent to the null value. Obviously, this probability depends enormously on the limits of the ROPE, and it is meaningful only if the limits of the ROPE are meaningful landmarks for the particular application. If the proportion of the HDI within the ROPE is most interpretable when it is large, because the proportion can be small for the trivial reason that the posterior is broad and shallower to a small data set.

12.3.2 Recommendations

In general, Bayesian model comparison is only useful when both models are genuinely viable. If one or both models has little prior believability, then the Bayes factor and relative posterior beliefabilities are of little use. For example, suppose that we observe that there are

many nicely wrapped gifts strewn across the lawn beside the house. We might hypothesize two models for this observation: One hypothesis is that Santa Claus dropped them from the roof. Another hypothesis is that the Grinch dumped them from the roof. A Bayesian model comparison might prefer the Santa Claus hypothesis 50 to 1 over the Grinch hypothesis, but that does not mean we should believe in Santa Claus. Model comparison requires both models to have prior viability, not just one. To explain this on the lawn, one hypothesis is that Santa Claus lost his grip (i.e., a non-viable hypothesis) and the other hypothesis is that a philosophically impressionable teenager in the neighborhood scattered the gifts in a fanatical fit of anti-materialist divestment (i.e., a viable but also unlikely hypothesis). A Bayesian model comparison might prefer the hypothesis of anti-materialist teenager, 50 to 1, but only by virtue of the fact that the alternative hypothesis was not viable in the first place. (Cf. Section 10.4.2.)

The premise that Bayesian model comparison is most meaningful when both models are genuinely viable implies that “automatic” Bayesian hypothesis testing might not be very meaningful, unless it is carefully applied only to situations in which both the null hypothesis is theoretically viable and the viable alternative hypothesis uses a reasonably informed prior. It is up to the user of the method to decide whether these conditions are met, and to carefully interpret the results of the comparison.

In general, Bayesian model comparison is a very useful technique outside its specific application to null-hypothesis testing, as long as the models are viable, and as long as the models have priors that are equivalently informed. Even though models are theoretically viable, if one model has a prior distribution on its parameters that is already close to the posterior, but the other model has a prior distribution on its parameters that is far from the posterior, then the models are not starting with priors that are equivalently informed. These general points about Bayesian model comparison apply especially strongly to the specific case of Bayesian null-hypothesis testing.

There may be situations that specifically demand a dichotomous decision regarding a null value. In these situations, Bayesian null-hypothesis testing may be called for, but the alternative model and alternative prior should be carefully considered. Otherwise, the estimation approach can be both more informative and easier to implement.

12.4 R code

12.4.1 R code for Figure 12.5

This code is a straight-forward modification of the code from Section 10.2.2. Instead of the values being either distinct or shared across groups, the values are either distinct or shared across groups. In principle, the values could (or, indeed, should) simultaneously be either distinct or shared, but here the values are always distinct, merely to keep the code more readable. The only other difference from the code from Section 10.2.2 is that fictitious data are generated, whereby the groups have masses of .61, .50, .49, and .51, respectively.

(OneOddGroupModelComp.R)

```

5 # -----
6 # THE MODEL.
7
8 modelstring = "
9 # BUGS model specification begins here...

```

```

10 model {
11   for ( i in 1:nSubj ) {
12     # Likelihood:
13     nCorrOfSubj[i] ~ dbin( theta[i] , nTrlOfSubj[i] )
14     # Prior on theta (notice nested indexing):
15     theta[i] ~ dbeta( aBeta[ CondOfSubj[i] ] , bBeta[ CondOfSubj[i] ] )I(0.0001,0.9999)
16   }
17   # Re-parameterization of aBeta[j],bBeta[j] in terms of mu a      nd kappa:
18   for ( j in 1:nCond ) {
19     # Model 1: Distinct mu[j] each group. Model 2: Shared mu0 all g    rous.
20     aBeta[j] <-      ( mu[j]*(2-mdllIdx) + mu0*(mdllIdx-1) ) * kappa[j]
21     bBeta[j] <- ( 1 - ( mu[j]*(2-mdllIdx) + mu0*(mdllIdx-1) ) ) * kap      pa[j]
22   }
23   # Hyperpriors for mu and kappa:
24   for ( j in 1:nCond ) {
25     mu[j] ~ dbeta( a[j,mdllIdx] , b[j,mdllIdx] )
26   }
27   for ( j in 1:nCond ) {
28     kappa[j] ~ dgamma( shk , rak )
29   }
30   mu0 ~ dbeta( a0[mdllIdx] , b0[mdllIdx] )

31   # Constants for hyperprior:
32   # (There is no higher-level distribution of across-group re      lationships,
33   # merely to keep the focus here on model comparison.)
34   shk <- 1.0
35   rak <- 0.1
36   aP <- 1
37   bP <- 1
38
39   a0[1] <- .53*400      # pseudo
40   b0[1] <- (1-.53)*400 # pseudo
41
42   a0[2] <- aP # true
43   b0[2] <- bP # true
44
45   a[1,1] <- aP # true
46   a[2,1] <- aP # true
47   a[3,1] <- aP # true
48   a[4,1] <- aP # true
49   b[1,1] <- bP # true
50   b[2,1] <- bP # true
51   b[3,1] <- bP # true
52   b[4,1] <- bP # true
53
54   a[1,2] <- .61*100      # pseudo
55   a[2,2] <- .50*100      # pseudo
56   a[3,2] <- .49*100      # pseudo
57   a[4,2] <- .51*100      # pseudo
58   b[1,2] <- (1-.61)*100 # pseudo
59   b[2,2] <- (1-.50)*100 # pseudo
60   b[3,2] <- (1-.49)*100 # pseudo
61   b[4,2] <- (1-.51)*100 # pseudo
62
63   # Hyperprior on model index:
64   mdllIdx ~ dcat( modelProb[] )
65   modelProb[1] <- .5
66   modelProb[2] <- .5
67
68 }
```

```

69 # ... end BUGS model specification
70 " # close quote for modelstring
71 # Write model to a file:
72 writeLines( text=modelstring , con="model.txt" )
73 # Load model file into BRugs and check its syntax:
74 modelCheck( "model.txt" )
75
76 #-----
77 # THE DATA.
78
79 # For each subject, specify the condition s/he was in,
80 # the number of trials s/he experienced, and the number corre      ct.
81 # (Randomly generated fictitious data.)
82 npg = 20 # number of subjects per group
83 ntrl = 20 # number of trials per subject
84 CondOfSubj = c( rep(1,npg) , rep(2,npg) , rep(3,npg) , rep(4      ,npg) )
85 nTrlOfSubj = rep( ntrl , 4*npg )
86 set.seed(47401)
87 nCorrOfSubj = c( rbinom(npg,ntrl,.61) , rbinom(npg,ntrl,      .50) ,
88                 rbinom(npg,ntrl,.49) , rbinom(npg,ntrl,.51) )
89 nSubj = length(CondOfSubj)
90 nCond = length(unique(CondOfSubj))
91 # Display mean number correct in each group:
92 for ( condIdx in 1:nCond ) {
93     show( mean( nCorrOfSubj[ CondOfSubj==condIdx ] ) )
94 }
95
96 # Specify the data in a form that is compatible with BRugs mode      l, as a list:
97 dataList = list(
98     nCond = nCond ,
99     nSubj = nSubj ,
100    CondOfSubj = CondOfSubj ,
101    nTrlOfSubj = nTrlOfSubj ,
102    nCorrOfSubj = nCorrOfSubj
103 )
104
105 # Get the data into BRugs:
106 modelData( bugsData( dataList ) )
107
108 #-----
109 # INTIALIZE THE CHAINS.
110
111 nchain = 3
112 modelCompile( numChains=nchain )
113 modelGenInits()
114
115 #-----
116 # RUN THE CHAINS.
117
118 burninSteps = 5000
119 modelUpdate( burninSteps )
120 samplesSet( c("mu","kappa","mu0","theta","mdlIdx") )
121 nPerChain = 5000 ; nThin = 10
122 modelUpdate( nPerChain , thin=nThin )

```

12.5 Exercises

Exercise 12.1.[Purpose: Model comparison for different partitions of group means.]

The program in Section 12.4.1(OneOddGroupModelComp.R) does a model comparison in which M1 has different means for every group and M2 has the same mean for all groups. In this exercise, we consider a different model comparison, with a different partition of the conditions.

(A) Consider a comparison for which M1 has one mean for condition 1 and a second mean for conditions 2 through 4, and M2 has the same mean for all groups. A quick-and-dirty way to program this is by changing the model specification to (OneOddGroupModelCompEx12.1.R)

```
24   for ( j in 1:2 ) {
25     mu[j] ~ dbeta( a[j,mdlIdx] , b[j,mdlIdx] )
26   }
27   mu[3] <- mu[2]
28   mu[4] <- mu[2]
```

Explain in English what the modified code does.

Unfortunately, the modified model structure is not amenable to automatic initialization in BUGS. We can “manually” initialize it this way: (OneOddGroupModelCompEx12.1.R)

```
120 genInitList <- function() {
121   sqzData = .01+.98*datalist$nCorrOfSubj/datalist$nTrlO      fSubj
122   mu = aggregate( sqzData , list(datalist$CondOfSubj) , "mean" )[,"x"]
123   sd = aggregate( sqzData , list(datalist$CondOfSubj) , "sd"    )[,"x"]
124   kappa = mu*(1-mu)/sd^2 - 1
125   return(
126     list(
127       theta = sqzData ,
128       mu = c( mu[1] , mean( mu[2:4] ) , NA , NA ) ,
129       mu0 = mean(mu) ,
130       kappa = kappa ,
131       mdlIdx = 1
132     )
133   )
134 }
135 for ( chainIdx in 1 : nchain ) {
136   modellnits( bugsInits( genInitList ) )
137 }
```

The initialization uses the same technique introduced in Section 9.5.2 (ficonBrugs.R), with only one novel nuance: The last two components are specified as *not* stochastic because those components are not stochastic, but are instead fixed (in the model specification) to equal mu[2].

Run the modified program, and show the resulting graphs *against* to those in Figure 12.5. Why are the upper histograms all the same, and *whether* lower histograms (e.g., μ_2 – μ_3) so strange looking?

(B) What should we conclude from the model comparison of the *posterior* part? (Be careful to express your conclusion as a statement about beliefabilities.) Should we conclude that the means of conditions 2 through 4 are *approximately* equal?

Exercise 12.2.[Purpose: Estimating a difference, including a ROPE (and the hot hand example with a better prior).]

(A) Consider again the hot-hand example from Exercise 8.1, 3. In this part of the exercise we establish a better prior: We know from general basketball statistics that professional players tend to make about 75% of their free throws, some players almost as low as 50% and some players almost as high as 95%. A beta distribution nicely captures this prior knowledge. But what we do not know from basketball statistics is the difference between success after success and success after failure. To be as vague as possible about the difference, we'll put a uniform distribution on the difference. Notice that when the overall success rate is, say, 90%, the difference could be anything from +20% (i.e., 100% vs. 80%) to -20% (i.e., 80% vs. 100%), but when the overall success rate is, say, 70%, then the difference could be anything from -60% to +60%. The specification of the prior needs to accommodate this range that depends on overall success. Here is one way to specify this in the BUGS model:

```
theta1 <- mu + deflect
theta2 <- mu - deflect
mu    dbeta( 16 , 6 )
delta   dbeta( 1 , 1 )
deflect <- (delta-.5)*2 * min(mu,1-mu)
```

The variable `mu` is the overall success rate. The variable `deflect` is the deviation away from `mu` created by a previous success or a previous failure. The variable `delta` is just a linearly transformed value of the random variable `deflect`, which has a uniform prior. Incorporate this code into the model and run without the data, so you can see the posterior. Show your complete model specification and a graph of the MCMC sample, which should look much like the upper row of Figure 12.6.

(B) Suppose we establish a ROPE on the difference of success-after-success and success-after-failure. We will arbitrarily set it at 5%. Run the program with the data included, and display the posterior. The `Post.R` function allows specification of a ROPE, with output as shown in Figure 12.6. Include your output graph and say in English what the ROPE indicates.

Exercise 12.3. Purpose: Applying the approaches to a real-world example. The thematic apperception test (TAT) is a method for assessing personality and other aspects of interpersonal attitudes. The subject is shown pictures of people in ambiguous scenes, and the subject is asked to make up a story about what the pictured people are doing. In a study reported by Werner, Stabenau, and Pollin (1970, Table 4), 40 pictures showing parent-child interactions were shown to women who had children. The stories invented by the mothers were scored for whether or not they expressed "personally involved, dependent, exible interactions". Twenty mothers of normal children and twenty mothers of schizophrenic children were each shown 10 pictures. The number of stories, out of 10, expressed a flexible interaction, were as follows.

Mothers of normal children: 8, 4, 6, 3, 1, 4, 4, 6, 4, 2, 2, 1, 13, 4, 2, 6, 3, 4.

Mothers of schizophrenic children: 2, 1, 1, 3, 2, 7, 2, 1, 3, 2, 4, 2, 3, 3, 0, 1, 2, 2.

For purposes of this exercise, we will assume that the prior information from previous research, which indicates that the typical number of stories expressing flexible interactions for this picture set is around 3 or 4, and can be well described by a $\text{beta}(3.5, 6.5)$ distribution. This is the distribution of flexible-narrative tendencies across mothers.

(A) Estimate the difference between TAT scores of the two groups of mothers. Be sure

²The author found these data by perusing the collection by Daly, Lunn, McConway, and Ostrowski (1994).

to use the prior knowledge. Is zero among the credible ~~estimates~~? If you assume that the difference has a ROPE from -0.1 to $+0.1$, is the ROPE excluded from the 95% HDI? HINTS: Set up a hierarchical model in BUGS, using the `laptopcondensation` example as a guide. The prior applies to the distribution of individual `theta[i]` values, and the data are distributed as a binomial with bias `theta[i]`.

(B) Using model comparison, determine the posterior probability of a model that uses one hyperparameter to describe both groups, relative to a model that uses a different hyperparameter for each group. For this application, the prior on the hyperparameters should not be informed, but instead should be extremely vague and broad. Because the model comparison is supposed to be “automatic”. What is the posterior probability of the different-hyperparameter model? If it is the preferred model, can we trust the estimate of the difference between groups? (Answer: Not necessarily, because it ~~uses~~ the prior knowledge.)

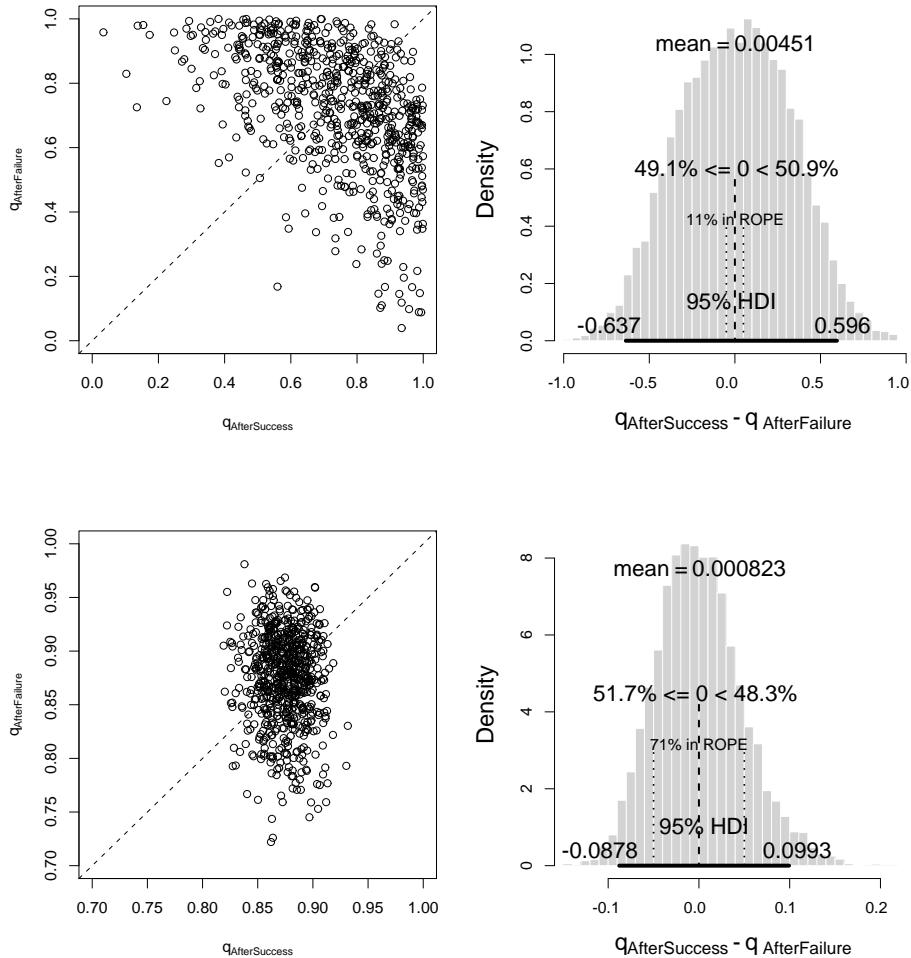


Figure 12.6: For Exercise 12.2 Upper row: The prior shows that we believe the overall success rate is between 50% and 95%, but the difference between the two 's could be anything. Lower row: The posterior. Notice the ROPE on the difference, from :05 to +:05, along with the display of what percentage of the posterior falls within the ROPE.

Chapter 13

Goals, Power, and Sample Size

Contents

13.1	The Will to Power	260
13.1.1	Goals and Obstacles	260
13.1.2	Power	261
13.1.3	Sample Size	262
13.1.4	Other Expressions of Goals	264
13.2	Sample size for a single coin	264
13.2.1	When the goal is to exclude a null value	265
13.2.2	When the goal is precision	266
13.3	Sample size for multiple mints	267
13.4	Power: prospective, retrospective, and replication	269
13.4.1	Power analysis requires verisimilitude of simulated data	270
13.5	The importance of planning	71
13.6	R code	272
13.6.1	Sample size for a single coin	227
13.6.2	Power and sample size for multiple mints	274
13.7	Exercises	281

Just how many times must I show her I care,
Until she believes that I'll always be there?
Well, while she denies that my value's enough,
I'll have to rely on the power of love.¹

Researchers collect data in order to achieve a goal. Sometimes the goal is to show that a suspected underlying state of the world is believable; sometimes the goal is to put a minimal degree of precision on whatever trends are observed. The goal may only be probabilistically achieved, as opposed to definitely achieved, because data are complete with random noise that can obscure the underlying state of the world. Statistical power is the probability of achieving the goal of an empirical study, when a suspected underlying state of the world is true. Scientists don't want to waste time and resources pursuing goals that have a small probability of being achieved. In other words, scientists不惜 power in their experiments.

¹The power of "Iu "? Sailors know that there's not much power in lug.

13.1 The Will to Power

In this section, a framework for research and data analysis is described, which leads to a more precise definition of power and how to compute it.

13.1.1 Goals and Obstacles

There are many possible goals for an experiment or observational study. For example, we might want to show that the rate of recovery for patients who take a drug is higher than the rate of recovery for patients who take a placebo. We might want to show that a coin is biased, i.e., that its tendency to show heads is not equal to tails. Goals such as those involve showing that a suspected difference or value really is tenable, or, complementarily, showing that a “null” value is not tenable. Other goals do not necessarily have a suspected value in mind. Instead, the goal is merely to put a desired degree of precision on whatever tendencies happen to be observed. For example, when polling a population for preferences of political candidates, the goal is not necessarily to show that an independent candidate is ahead, but to get a precise assessment of the population's preferences.

The various goals of research can be formally expressed in two ways. In this chapter I will focus on two goals, formalized in terms of the HDI.

Goal: Demonstrate that believable parameter values exclude a value.

- Formal expression: Show that the 95% HDI excludes the “null” value, or its ROPE.

Goal: Achieve precision in the estimate of believable values.

- Formal expression: Show that the 95% HDI has some specified maximal width.

In some research, the goal may be to demonstrate that a “value” is true, rather than false. This goal can also be formalized in terms of precision of the HDI: We want the entire 95% HDI to fall within the ROPE, as described in Section 13.1 (p. 244). There are other mathematical formalizations of the various goals, and they will be mentioned later. This chapter focuses on the HDI because of its natural interpretation for purposes of parameter estimation, as explained in the previous chapter.

If we knew the benefits of achieving our goal, and the costs of failing it, and if we knew the penalties for making a mistake while interpreting the data, then we could express the results of the research in terms of the long-run expected value. When we know the costs and benefits, we can conduct a full decision-theoretic treatment of the situation, and plan the research and data interpretation accordingly (e.g., Berger & Verdinelli, 1995; Lindley, 1997). In our applications we do not have access to those exact benefits, unfortunately. Therefore we rely on goals such as those outlined above.

The crucial obstacle to the goals of research is that a random sample is only a probabilistic representation of the population from which it came. Even if a coin is actually fair, a random sample of 100 flips will rarely show exactly 50% heads. Even if a coin is biased, it might come up heads 5 times in 10 flips. Drugs that actually work better than a placebo might happen to cure more patients in a particular random sample. And drugs that truly are

²Regarding the title of this section: Other than the fact that researchers desire statistical power, the notion of statistical power might have profound connections with concepts from Friedrich Nietzsche's work The Will to Power. See Exercise 13.1.

effective might happen to show little difference from a placebo in another particular random sample of patients. Thus, a random sample is a ~~clock~~ indicator of the true value in the underlying world. Whether the goal is showing that a ~~suspected~~ ~~state~~ is or isn't believable, or showing that a precise range of values is tenable, ~~randomization~~ is the researcher's bane. Noise is the nemesis.

13.1.2 Power

Because of random noise, the goal of a study can be achieved probabilistically. The probability of achieving the goal, given the (suspected) state of the world, is called the power of the experiment.

Scientists go to great lengths to try to increase the power of experiments or observational studies. There are three primary methods by which researchers can increase the chances of detecting an effect. First, we reduce measurement noise as much as possible. For example, if we are trying to determine the cure rate of a drug, to reduce other random influences on the patients, such as other drugs they might be taking or starting, changes in diet or rest, etc. Reduction of noise and control of influences is the primary reason for conducting experiments in the lab instead of in the ~~messes~~ of the real world. The second method, by which we can increase the chance of detecting an effect, is to amplify the underlying magnitude of the effect if we possibly can. For example, if we are trying to show that a drug helps cure a disease, we will want to administer as large a dose as possible (assuming there are no negative side effects). In non-experimental research, in which the researcher does not have the luxury of manipulating objects being studied, this second method is unfortunately unavailable. Sociologists, economists, and astronomers, for example, are often restricted to observing events that they cannot control or manipulate.

Once we have done everything we can to reduce noise in our measurements, and to amplify the effect we are trying to measure, the third way to increase power is to increase the sample size. The intuition behind this method is simple: With more and more measurements, random noises will tend to cancel themselves out, leaving on average a clear signature of the underlying effect. In general, as sample size increases, power increases. Increasing the sample size is an option in most experimental research, and in a lot of observational research (e.g., more survey respondents can help), but not in some domains where the population is finite, such as comparative studies of the states or provinces of a nation. In this latter situation, we cannot create a larger sample size, but Bayesian inference is still valid, and perhaps uniquely so (Western & Jackman, 2014).

What we accomplish in this chapter is precise calculations of power. We compute the probability of achieving a specific goal, given (i) a specified underlying distribution of biases or effects in the population being measured, and (ii) a specific generating process, such as collecting a fixed number of observations. These calculations are very useful for planning an experiment. To plan our research, we conduct "dress rehearsals" before the actual performance. We repeatedly simulate data we suspect we might get, and conduct Bayesian analyses on the simulated data sets. If the goal is achieved for most of the simulated data sets, then the planned experiment has power. If the goal is rarely achieved in the analyses of simulated data, then the planned experiment is likely to fail, and we must do something to increase its power.

In general, power can be approximated in the following manner:

1. Generate a random sample of data points, using the data-generating hypothesis. The sample should be generated according to how actual data would be gathered in the

eventual real experiment. For example, typically it is ~~assumed~~ that the number of data points is ~~fixed at N~~. It might instead be assumed that data will be collected for a fixed interval of time T , during which data points appear randomly at a known mean rate $\lambda = T$.

2. Compute the posterior estimate, using Bayesian analysis ~~with skeptical~~ priors. The data analysis must be convincing to the audience, ~~and the~~ the analysis must use priors that are agreeable to that skeptical audience.
3. From the posterior estimate, tally whether or not the goal is attained. The goal could be any of those outlined previously, such as having ~~95%~~ HDI exclude a ROPE around the null value, or having the 95% HDI be narrower than a desired width.
4. Repeat above steps many times, to approximate the power. Power is, by definition, the proportion of times that the goal is attained.

The result of this process is the power, i.e., the probability of achieving the goal when using the particular data-sampling procedure. Notice that if the data-sampling procedure uses a fixed sample of size N , then the process determines power as a function of N . If the data-sampling procedure uses a fixed sampling duration T , then the process determines power as a function of T .

Figure 13.1 illustrates the process of computing power. The top part of the figure shows the flow of information in an actual analysis. The real world generates a single sample of observed data. We use Bayes' rule, starting with a prior belief about the skeptical audience, to derive an actual posterior distribution. The lower part of the figure shows the flow of information in a power analysis. Instead of the real world, we have a hypothesis about the world. The hypothesis is expressed as a belief distribution over the parameters of the model. In many cases, this hypothesis is a posterior distribution derived from previous research. From the hypothesis, a sample of data is generated. This simulated sample of data is exemplary of data we anticipate if the hypothesis is true. We then conduct a Bayesian analysis on the anticipated data sample, thereby deriving an anticipated posterior distribution. For the anticipated posterior, we assess whether the goal was achieved. We then repeat the simulated data generation process. This process is indicated in Figure 13.1 by the layers of anticipated data samples and anticipated posteriors. Across repetitions, the proportion of times that the goal is achieved is the power. In some simple situations, power can be determined exactly, as is done, for example, in the program of Section 13.6.1 (`minNforHDIpower.R`). In more realistic situations, power is approximated through simulation (Wang & Gelfand, 2002), as is done, for example in the program of Section 13.6.2 (`FilconBrugsPower.R`).

13.1.3 Sample Size

Power increases as sample size increases (usually). Gathering data is costly, we would like to know the minimal sample size, or sample duration, required to achieve a desired power.

The goal of precision in estimation can always be attained with a large enough sample size. This is because the likelihood of the data, thought of logically as a function of the parameter, tends to get narrower and narrower as sample size increases. This narrowing

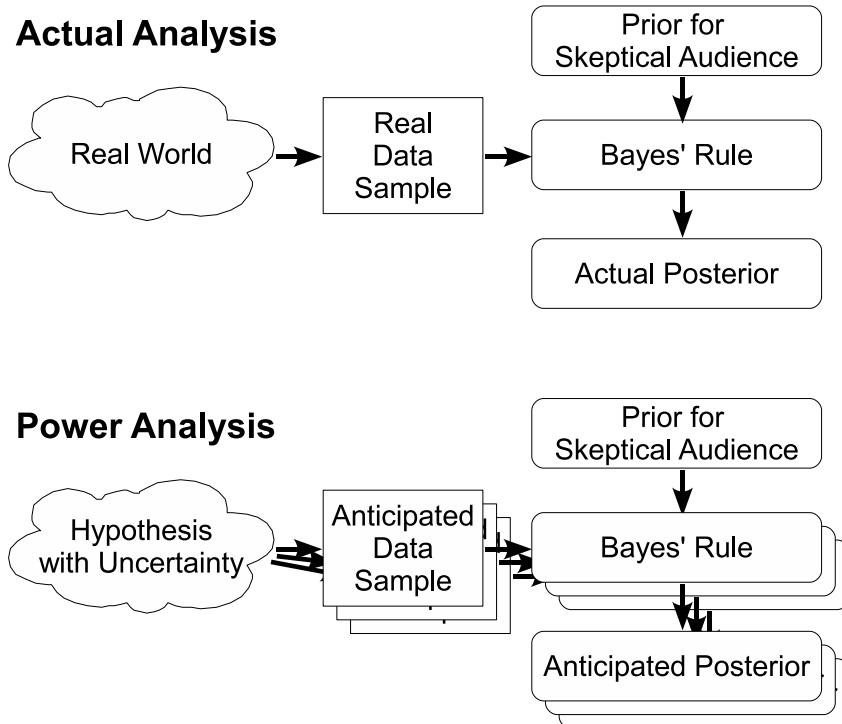


Figure 13.1: Upper diagram illustrates the flow of information in an actual Bayesian analysis. Lower diagram illustrates flow of information in a power analysis.

of the likelihood is also what causes the data to eventually overwhelm the prior beliefs as sample size increases. As we collect more and more data, the likelihood function gets narrower and narrower (on average), and therefore the posterior gets narrower and narrower. Thus, with a large enough sample, we can make the posterior distribution as narrow, i.e., precise, as we like.

The goal of showing that a parameter value is different from a null value might not be attainable with a high enough probability, however, no matter how big the sample size. Whether or not this goal is attainable with high probability depends on the data-generating distribution. The best that a large sample can do is exactly reflect the data-generating distribution. If the data-generating distribution has considerable mass around the null value, then the best we can do is get estimates that include mass around the null value. As a simple example, suppose that we think that a coin may be biased and the data-generating hypothesis is that $p(\theta = .5) = 25\%$, $p(\theta = .6) = 25\%$, $p(\theta = .7) = 25\%$, and $p(\theta = .8) = 25\%$. Because 25% of the simulated data come from a fair coin, the maximal probability of excluding $\theta = .5$, even with a huge sample, is 75%.

Therefore, when planning the sample size for an experiment, it is crucial to first decide what a realistic goal is. If there are good reasons to posit a very certain data-generating hypothesis, perhaps because of extensive previous results, a viable goal may be to exclude a null value. On the other hand, if the data-generating hypothesis is somewhat vague, then a more reasonable goal is to attain a desired degree of precision in the posterior.

13.1.4 Other Expressions of Goals

There are other ways to express mathematically the goal of precision in estimation. For example, another way of using HDIs was described by Joseph et al., and du Berger (1995a, 1995b). They considered an “average length criterion” which requires that the average HDI width, across repeated simulated data, does not exceed maximal value. There is no explicit mention of power, i.e., the probability of achieving the goal, because the sample size is chosen so that the goal is definitely achieved. The goal itself is probabilistic, however, because it regards an average: While some data sets have HDI width less than L, many other data sets will not have an HDI width less than L. Another goal considered by Joseph et al. (1995a) was the “average coverage criterion”. This goal starts with a specified width for the HDI, and requires its mass to exceed 95% (say, on average across simulated data). The sample size is chosen to be large enough to achieve this goal. Again, power is not explicitly mentioned, but the goal is probabilistic: Some data sets will have an HDI width mass greater than 95%, and other data sets will have an HDI width mass greater than 95%. Other goals regarding precision are reviewed by Adcock (1997) and by De Santis (2004, 2007). The methods emphasized in this chapter focus on limiting the worst precision, instead of the average precision.

A rather different mathematical expression of precision is entropy of a distribution. Entropy describes how spread out a distribution is, so small entropy connotes a more precise distribution. A distribution that consists of an infinite spike, that has an infinitesimally narrow width, has zero entropy. At the opposite extreme a uniform distribution has maximal entropy. The goal of high precision in the posterior distribution might be re-expressed as a goal of small entropy in the posterior distribution. For an overview of this approach, see, e.g., Chaloner and Verdinelli (1995). For an introduction to how minimization of expected entropy might be used spontaneously by people as they experiment with the world, see Kruschke (2008). Entropy may be a better measure of posterior precision than HDI width especially in cases of multimodal distributions for which HDI width is more challenging to determine. I will not further explicate the use of entropy because I think that HDI width is a more intuitive quantity than entropy by which to express precision, at least for most researchers in most contexts.

There are also other ways to express mathematically the goal of excluding a null value. In particular, the goal could be expressed as wanting a Bayes factor in a model comparison between the spike-null prior and the automatic alternative prior (e.g., Wang & Gelfand, 2002; Weiss, 1997). For example, we might set the desired Bayes factor at a ratio of 19 to 1 (i.e., .95 to .05). Kass and Raftery (1995) provide some guidelines for choice of criterial Bayes factor. I will not further address this approach, however, because the goal of a criterial BF for untenable caricatured priors has problems as discussed in the previous chapter. Instead, it will be assumed that the goal of the researcher is to obtain a certain number of parameter values, starting with a viable prior. The resulting posterior is then used to assess whether the goal was achieved.

13.2 Sample size for a single coin

As our first worked-out example, consider the simplest case: Data from a single coin. Perhaps we are polling a population and we want to know if there is a preference for candidate A or candidate B. Or perhaps we want to know if a drug has more than 50% cure rate. We will go through the steps listed in Section 13.1.2 to complete the exact sample size needed

Table 13.1: Minimal sample size required for 95% HDI to exclude = .5, when ipping a single coin.

Power	Generating Mean						
	.55	.60	.65	.70	.75	.80	.85
.7	642	148	64	33	18	13	7
.8	866	191	82	43	26	18	10
.9	1274	264	111	59	36	24	16

Note. The data-generating distribution is a beta density with mean = .65, as indicated by the column header, and with shape parameters of $\alpha = 1$ and $\beta = 2000$, where $N = 2000$. The audience-agreeable prior is uniform.

to achieve various degrees of power for different data-generating hypotheses.

13.2.1 When the goal is to exclude a null value

Suppose that our goal is to show that the coin is biased. In other words, we want to show that .5 is not among the credible values. More specifically, we want to show that the 95% HDI excludes = :50.

Because of previous research (or just a really strong hunch), we think that the true bias of the coin is very close to = :65. It is a strong belief, so we'll describe the data-generating hypothesis as a beta distribution with mean = :65 and certainty based on 2,000 flips of the coin. This puts 95% of the data-generating biases between :63 and :67.

Next, we generate simulated data from hypothetical coins that have those biases, and tally how often the HDI excludes = :50. One replication of the process goes like this: First, select a value for the “true” bias in the coin, from a beta-generating distribution that is narrowly centered on = :65. Suppose that the selected value is .638. Second, simulate ipping a coin with that bias N times. The simulated data have N heads and N tails. The proportion of heads, p , will tend to be around .638, but will be higher or lower because of randomness in the flips. Third, using the audience-agreeable prior for purposes of data analysis, determine the posterior beliefs regarding p if N heads in N flips were observed. Tally whether or not the 95% HDI excludes the null value of = :50. Notice that even though the data were generated by a coin with bias of .638, they might, by chance, show a low proportion of heads, and therefore the 95% HDI might exclude = :50. This process is repeated many times to estimate the power of the test, i.e., the probability that the null value will be excluded from the HDI. In fact, in this simple situation, the exact power can be computed analytically, as explained in Section 13.6.1 (minNforHDIpower.R), p. 272.

Table 13.1 shows the minimal sample size needed for the 95% HDI to exclude = :50 when ipping a single coin. As an example of how to read the table, suppose you have a data-generating hypothesis that the coin has a bias very near .65. This hypothesis is implemented, for purposes of Table 13.1, as a beta distribution with shape parameters of :65 and (1 - :65) = 2000. The value of 2000 is arbitrary; it's as if the generating mean of .65 was based on 2000 previous data containing 1002 heads. The table indicates that if we desire a 90% probability of obtaining a 95% HDI that excludes = :50, we need a sample size of $N = 111$, i.e., we need to flip the coin at least 111 times in order to have a 90% chance that = :5 falls outside the 95% HDI.

Notice, in Table 13.1, that as the generating mean increases, the required sample size decreases. This makes sense intuitively: When the generating mean is large, the sample proportion of heads will tend to be large, and so the HDI will fall toward the high end of the parameter domain. In other words, when the generating mean is large, it doesn't take a lot of data for the HDI to fall well above = :5. On the other hand, when the generating mean is only slightly above = :5, then it takes a large sample for the sample proportion of heads to be consistently above .5, and for the HDI to be consistently entirely above .5.

Notice also, in Table 13.1, that as the desired power increases, the required sample size increases quite dramatically. For example, if the data-generating mean is .6, then as the desired power rises from .7 to .9, the minimal sample size increases from 148 to 266.

13.2.2 When the goal is precision

Not all research has as its goal the exclusion of a particular value. Sometimes the goal is to establish a precise estimate of the parameter value. Sometimes, the research may have a null value of interest, but the data-generating hypothesis is too vague for the null value to be excluded most of the time, or, the goal may be to demonstrate that the HDI falls entirely within the ROPE. In these cases, the goal becomes precision of parameter estimation.

Here is an example in which there is a null value of interest, only a vague data-generating distribution. Suppose you are interested in the preferences of the general population regarding political candidates A and B. In particular, you would like to have high confidence in estimating whether the preference for candidate A exceeds = :5. A recently conducted poll by a reputable organization found of 10 randomly selected voters, 6 preferred candidate A and 4 preferred candidate B. We use a uniform pre-poll prior, our post-poll estimate of the population bias is a beta(6, 4) distribution. As this is our best information about the population so far, we can use a beta(7, 5) distribution as a data-generating distribution for planning the follow-up poll. Unfortunately, a beta(7, 5) distribution has a 95% HDI from = :318 to = :841, which means that = :5 is well within the data-generating distribution. How many more people do we need to poll so that 80% of the time we would get a 95% HDI that falls above = :5?

It turns out, in this case, that we can never have a sample size enough to achieve the goal of 80% of the HDIs falling above = :5. To see why, consider what happens when we sample a particular value from the data-generating distribution, e.g., = :4. We use that value to simulate a random sample of votes. Suppose the sample is huge, which implies that the HDI will be very narrow. What value of will the HDI focus on? Almost certainly it will focus on the value = :4 that was used to generate the data. To reiterate, when N is very large, the HDI essentially just reproduces the value that generated it. Now recall the data-generating distribution of our example: Beta(7; 5) has only about 72% of the values are above .5. Therefore, even with an extremely large sample size, we can get at most 72% of the HDIs to fall above .5.

There is a viable alternative goal, however. Instead of trying to reject a particular value of , we set as our goal a desired degree of precision in the posterior estimate. For example, our goal might be that the 95% HDI has width less than 0.2, at least 80% of the time. This goal implies that regardless of what values happen to be emphasized by the posterior distribution, the width of the posterior is usually narrow, that we have attained a suitably high precision in the estimate.

Table 13.2 shows the minimal sample size needed for the 95% HDI to have maximal

Table 13.2: Minimal sample size required for 95% HDI to have maximal width of 0.2, when flipping a single coin.

Power	Generating Mean						
	.55	.60	.65	.70	.75	.80	.85
.7	91	91	89	86	80	69	58
.8	92	92	91	90	86	79	67
.9	93	93	93	92	91	88	79

Note. The data-generating distribution is a beta density with mean mean , as indicated by the column header, and with shape parameters of $\alpha = 10$ and $\beta = 10$, where $\text{mean} = 10$. The audience-agreeable prior is uniform.

width of 0.2. As an example of how to read the table, suppose we have a data-generating hypothesis that the coin has a bias roughly around 0.6. This hypothesis is implemented, for purposes of Table 13.2, as a beta distribution with shape parameters of $\alpha = 10$ and $\beta = 10$. The value of 10 is arbitrary; it's as if the generating mean were based onitious previous data containing only 10 flips. The table indicates that if we desire a 90% probability of obtaining an HDI with maximal width of 0.2, we need a sample size of at least 93.

Notice in Table 13.2 that as the desired power increases, the required sample size increases only slightly. For example, if the data-generating mean is .6, then as the desired power rises from .7 to .9, the minimal sample size rises from 91 to 93. This is because the distribution of HDI widths, for a given sample size, has a very shrunken high tail, and therefore small changes in the mean can quickly pull the high tail across a threshold such as 0.2. On the other hand, as the desired HDI width decreases (not shown in the Table), the required sample size increases rapidly. For example, if the data-generating mean is .6 with $\alpha = 10$, and the desired HDI width is 0.1 (instead of 0.2), then the sample size needed for 80% power is 377 (instead of 92).

The R code in Section 13.6.1 (`NforHDIpower.R`), generated Tables 13.1 and 13.2. You can use the R code to specify your own data-generating distribution, goal, and desired power. You can also specify a ROPE. The R function will tell you the minimal sample size required. Exercise 13.2 has you work through some examples.

13.3 Sample size for multiple mints

Consider again the Iteration-condensation experiment described in Section 9.3.1, p. 178, in which people learned four different category structures. Two structures were “Iteration structures in which one stimulus dimension could be ignored” and “perfect accuracy could be achieved”. Two other structures were “condensation” structures in which both stimulus dimensions required processing to achieve perfect accuracy. Some theories of learning predict that Iteration should be easier to learn than condensation, and therefore one goal of the experiment was to show that the mean accuracy of people in the Iteration structure differed from the mean accuracy of people in the condensation structure. This goal was handily achieved: Using $N = 40$ participants per condition, the 95% HDI on the difference of parameters went from 0.173 to 0.259, with 100% of the MCMC differences falling well above 0.0; see the right-most panel of Figure 9.16, p. 18.

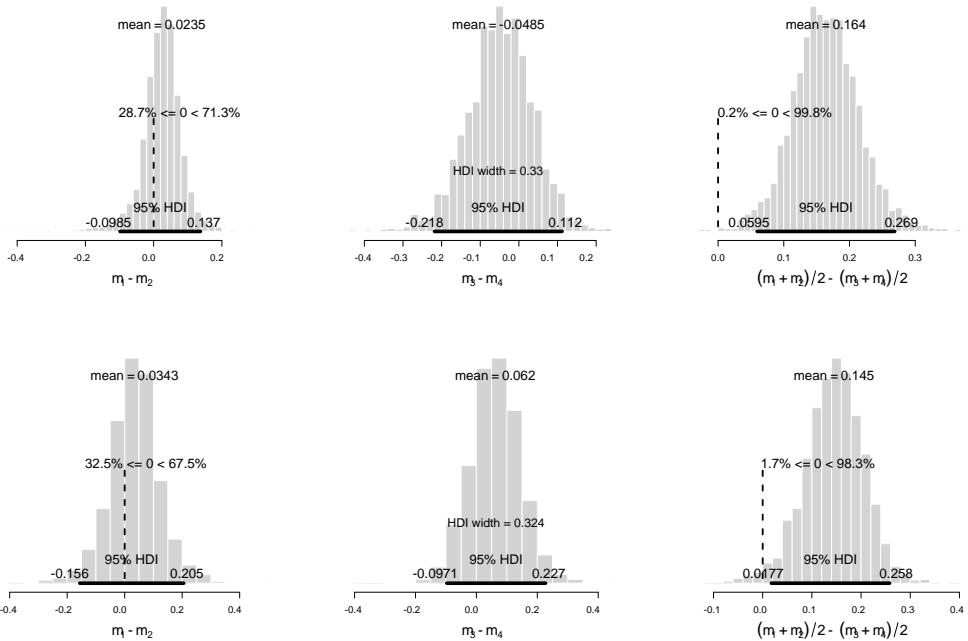


Figure 13.2: Posterior distributions from two runs of simulated data, using $N = 6$ per group. Upper row (right panel) shows a case in which the posterior for $(\mu_1 + \mu_2)/2 - (\mu_3 + \mu_4)/2$ excludes 0.0, as desired. Lower row (right panel) shows a case in which that goal is not achieved. Across many simulations, over 80% achieved the goal.

Is there a smaller sample size per group that could have been used to achieve high power? And what was the power of detecting the more subtle difference between the two iteration groups, shown in the left-most panel of Figure 9, p. 181? To answer these questions, we pursue the simulation procedure outlined at the beginning of this chapter (Section 13.1.2): Specify a data-generating hypothesis, simulate some data, do a Bayesian data analysis, determine whether the goal is achieved, repeat many times to estimate the probability of achieving the goal. R code for computing power for this experiment is listed in Section 13.6.2 (IconBrugsPower.R).

Suppose the goal of the experimenter is to show that the mean of the iteration groups exceeds the mean of the condensation groups. Specifically, suppose that the goal is achieved if the 95% HDI of $(\mu_1 + \mu_2)/2 - (\mu_3 + \mu_4)/2$ excludes 0.0. It turns out that using merely $N = 6$ in each of the groups achieves over 80% power for this goal. Figure 13.2 shows examples of the posteriors for two runs of simulated data with $N = 6$ per group. The upper row shows a case in which the posterior for $(\mu_1 + \mu_2)/2 - (\mu_3 + \mu_4)/2$ excludes 0.0, as desired. The lower row shows a case in which that goal was not achieved. Across many simulated runs using $N = 6$, 80% achieved the goal. In other words, instead of using $N = 40$, as in the actual experiment, the research could have been conducted with only $N = 6$, and had an 80% chance of inferring a believable difference between iteration and condensation.

This power analysis stems from simulated data that have been generated from moderately certain knowledge, based on results from 40 actual tests per condition. If we did

Table 13.3: Types of power analysis, including replication probability.

Analysis Type	Data Generator	Data Sample	Prior for Bayesian Analysis	Posterior
Actual	Real World	Observed Once	Skeptical Audience	Actual
Prospective Power	Hypothesis	Simulated Repeatedly	Skeptical Audience	Anticipated
Retrospective Power	Actual Posterior	Simulated Repeatedly	Skeptical Audience	Anticipated
Replication Power	Actual Posterior	Simulated Repeatedly	Actual Posterior	Anticipated

not have data-generating priors with such high certainty, our power predictions would not be so high either. For example, if we had data-generations on μ_1 with the same means as used here, but with wider (more uncertain) distributions, then the sample size would have to be much bigger than $N = 6$ to achieve the same power. This makes intuitive sense: If we are not very certain about what the data will look like, it takes a lot of data to reliably reveal the underlying trends.

Suppose that the goal of the experimenter was, instead to show that there was a difference between the two iteration groups. Specifically, we suppose that the goal is achieved if the 95% HDI of $\mu_1 - \mu_2$ excludes 0.0. Notice that this goal was not achieved in either simulated run shown in Figure 13.2. It turns out that using $N = 6$ in each of the groups achieves a power of only 4% for this goal. What is the power for this goal if $N = 40$, as used in the actual experiment? It turns out that the power is 43%. The fact that the actual data (Figure 9.16, p. 181) achieved the goal was merely a fluke by happenstance.

Finally, suppose that the goal of the experimenter was to reach a minimal precision on the difference between the two condensation groups. The motivation for this goal is as follows: Theory suggests that the two condensations should not differ much. Therefore the goal cannot be to show that the difference exceeds zero. Instead, the goal should be to achieve some desired degree of precision regarding the difference. Suppose we want the 95% HDI on $\mu_3 - \mu_4$ to have a width less than 0.20. To achieve this goal at least 80% of the time, i.e., with a power of 80%, we would need a sample size (per group) of approximately $N = 32$. You can verify this by running the program in Section 12.6 (FilconBrugsPower.R).

13.4 Power: prospective, retrospective, and replication

There are different types of power analysis, as shown in Table 13.3. For example, the top row of the table shows an actual analysis, in which the data generator is the real world, the data sample is observed once, the prior for the analysis is designed to be agreeable to a skeptical audience, and the actual posterior distribution is computed.

The second row of Table 13.3 shows prospective power analysis, for which the data generator is a research hypothesis based on theory or results from previous experiments that are not exactly the same as the present experiment. The research hypothesis is expressed as a distribution over model parameters, with both a tendency and uncertainty

that captures the researcher's beliefs about what would be appropriate parameter values if the data of the planned experiment were to appear as they would. The distribution over the model parameters could be set directly, theoretical considerations, by the researcher setting the shape constants of the prior distribution. On the other hand, the data-generating distribution could be set indirectly, first positing data consistent with the theory and then applying Bayesian updating to derive a posterior distribution over parameters that is consistent with the posited data. De Santis (2007) describes how to combine results of previous experiments to construct a data-generating distribution.

Contrast a prospective power analysis with the next row of Table 13.1, which shows a retrospective power analysis. This refers to a situation in which we have already collected data, and therefore we already have an actual posterior distribution. We want to determine the power of the experiment we ran. To do this, we use the laptop as the data generator. This is tantamount to using the data from a posterior predictive check. In other words, at a step in the chain, the parameter values are used to generate simulated data. The simulated data are then analyzed with the same Bayesian model as the actual data, and the posterior is examined for whether or not the goal is achieved. This process is readily programmed in BRugs, as explained in Section 13.6.2 (conBrugsPower.R). The example of the previous section, regarding the Itration-condon experiment, was a case of retrospective power analysis.

Finally, suppose that we have already collected some data and have an actual posterior distribution, and we want to know the probability that we would achieve our goal if we replicated the experiment. In other words, if we were to re-collect a new sample of data, what is the probability that we would achieve goal in the replicated study, also taking into account the first set of data? The bottom row of Table 13.3 indicates how we answer this question of replication power. We use the actual posterior derived from the first sample of data as the data generator. But for the analysis, we again use the actual posterior from the first sample of data, because that is the best-informed prior for the follow-up experiment. An easy way to execute this analysis in BUGS is as follows: Use the actual set of data with a skeptical-audience prior to generate a new sample from an actual posterior. Use that actual-posterior sample to generate simulated new data, as if in a posterior predictive check. Then concatenate the original data with the novel simulated data and update the original skeptical-audience prior with the new data set. This technique is tantamount to using the posterior of the original data set as the prior for the novel simulated data. Section 13.6.2 (conBrugsPower.R) shows BRugs code for this type of analysis.

Computation of replication power is natural in a Bayesian setting, but is difficult or impossible for traditional null-hypothesis significance testing (Miller, 2009). NHST has trouble when addressing replication probability because it has no good way to model a data generator: It has no access to the posterior distribution from the initial analysis.

13.4.1 Power analysis requires verisimilitude of simulated data

Power analysis is only useful when the simulated data imitate the actual data. We generate simulated data from a descriptive model that has uncertain parameter values, but we assume that the model is a reasonably good description of the actual data. If the model is instead a poor description of the actual data, then the simulated data do not imitate actual data, and inferences from the simulated data are not very useful. It is advisable, therefore, to perform a posterior predictive check that the simulated data accurately reflect the actual data.

When simulated data differ from actual data, strange results can arise in power analysis. Consider an analysis of replication probability in which simulated data are quite different than the actual data. The novel simulated data are combined with the original data to conduct the replication analysis. The combined data is a biased mixture of two different trends (i.e., the actual trend and the client simulated trend), and therefore the estimates of the parameters become more uncertain than for the original data alone. It is only when the simulated sample size becomes large, relative to the original sample size, that the simulated trend overwhelms the actual trend, and the replication uncertainty becomes smaller again. If you find in your power analyses that parameter uncertainty initially gets larger as the simulated sample size increases, then you may have a situation in which the model does not faithfully mimic the actual data.

13.5 The importance of planning

In the real conduct of research, scientists do not always specify a sample size and then stick to it. Nor should they have to. When experiments are set up such that the results from one participant do not affect the next participant, or when observational studies are conducted such that the results of one observation are independent of other observations —in other words when the next coin flip is not influenced by previous flips—then the researcher's intentions do not influence the data production or the data interpretation.

In some situations, data are collected for a particular period, during which the actual number of participants (coin flips) is a random value that depends on the random rate at which people happen to participate. For example, participants in a typical university psychology experiment might be recruited for a certain number of weeks during a semester. There may be a stable average rate of participation, but the number of participants in any given week is a random number. In this situation, our goal is to determine how many weeks we should plan on recruiting participants, so to have a high probability of achieving the goal of the experiment. We can do the power analysis in a manner analogous to the methods of the previous sections. As in those methods, we randomly generate parameter values from our data-generating hypothesis, but the sample size is a random value that depends on a specified rate of participation and amount of time during which participants are recruited. The sample size could be (but is not necessarily) modeled as a Poisson variable, as was done in Exercise 11.3. We simulate incrementally longer durations until we find the minimal duration that yields a sufficiently high power. Having thereby planned a duration for collecting data (and corresponding approximate number of subjects), we can decide whether or not it is feasible to actually conduct the research. If we do go ahead with the research, we are not stuck with the exact planned duration. We can stop whenever we want, assuming that the data are uninfluenced by our intentions to stop.

The point of the previous paragraphs is to say that sticking to a planned sample size or duration is not crucial to the data analysis, after the data have been collected. Instead, the planning is important for deciding whether or not the research has a reasonable chance of success, before the data are collected.

Conducting a power analysis in advance of collecting data is important and valuable. Often in real research, a fascinating theory and clever experimental manipulation imply a subtle effect. It can come as a shock to the researcher when power analysis reveals that detecting the subtle effect would take several hundred subjects! But the shock of power analysis is far less than the pain of running dozens of studies and finding highly uncertain

estimates of the sought-after effect.

Power analysis can reduce research pain in other ways. ~~In~~ In real research, an experiment or observational study is conducted merely objectively confirming what is anecdotally known to be a strong effect. A researcher may be tempted to conduct a study using the usual large sample size that is typical of related research. But a power analysis may reveal that the strong effect can be easily detected with a much smaller sample size.

Power analysis is also important when proposing research funding agencies. Proposals in basic research might have fascinating theories or clever experiments, but if the predicted effects are subtle, then reviewers of the proposal may be ~~skeptical~~ dubious, and want to be reassured by a power analysis. Proposals in applied research are even more reliant on power analysis, because the costs and benefits are ~~immediate~~ immediate and tangible. For example, in clinical research (e.g., medicine, pharmacology, psychiatry, counseling), it can be very costly to test patients, and therefore it is important to anticipate the probable sample size or duration.

While it is important to plan sample size in advance, it is also important, especially in clinical applications, to monitor data as they are collected and to stop the research as soon as possible. It behooves the researcher to discontinue the experiment as soon as the data clearly indicate a positive or negative outcome: It would be unethical to slavishly continue treating patients with an experimental treatment that is clearly detrimental, and it would be unethical to slavishly continue running patients in a placebo condition when the experimental treatment is clearly having positive effects. If the data are collected in such a way that they are uninfluenced by the experimenter's intent, then the data collection can be stopped at any time without influencing the interpretation of the data. The decision regarding when to stop collecting data is a topic of much ~~importance~~ significance, and goes under the name of Bayesian sequential design. It will not be discussed further here, but the interested reader is referred to books by J. O. Berger (1985) and DeGroot (2004), and various technical articles, for example those cited by Roy, Ghosal, and Rosenblub (2009, p. 427).

13.6 R code

13.6.1 Sample size for a single coin

The program described in this section was used to generate Tables 13.1 and 13.2. The program determines the minimal sample size needed to achieve a specified goal with a specified power, when flipping a single coin.

In principle, we could simulate random samples from the generating distribution, and then simulate random samples of coin flips based on the prior values of π . For each repetition, we would tally whether the goal was achieved. With a large number of repetitions, this procedure would provide a good approximation to the power.

In the present simple situation, however, we can determine ~~mathematically~~ the exact probabilities of each possible outcome. The simulated data generated by sampling a value according to our data-generating prior, and then ~~repeating~~ generating N flips of the coin according to the binomial distribution. Therefore the probability of getting z heads, across repeated sampling from the prior, is

$$p(z|N) = \sum_0^N p(z|N; \pi) p(\pi)$$

$$\begin{aligned}
 Z_1 &= d \text{ binomial}(z|N) \text{ Bernoulli}(j|a; b) \\
 &= \frac{Z^0}{0!} \frac{N!}{z} z(1 -)^{(N-z)} (a^{-1})(1 -)^{(b-1)} = B(a; b) \\
 &= \frac{N}{z} B(z+a; N-z+b) = B(a; b)
 \end{aligned} \tag{13.1}$$

(This probability of possible data is sometimes called the posterior marginal distribution of z ; cf. Eqn. 5 of Pham-Gia & Turkkan, 1992) For each possible outcome, z , we update the audience-agreeable prior to render an audience-agreeable posterior distribution, and then we decide whether the goal has been achieved for that outcome. Because the decision is determined by the outcome, the probability of the decisions is determined by the probability of the outcomes. Equation 13.1 is implemented in the program on lines 20–22, but in logarithmic form to prevent overflow errors.

The function `minNforHDIpower` has several arguments, most of which have default values. Two arguments without defaults are the mean `genPriorMean`, and effective sample size, `genPriorN`, of the data-generating hypothesis, which is assumed to be a beta density. For example, suppose that the experimenter believes that the bias of a coin is very nearly .65, and suppose that the belief is based on the equivalent 2000 of conscientious flips of the coin; then `genPriorMean = .65` and `genPriorN = 2000`.

The next two arguments, `HDImaxwid` and `nullVal`, specify the goal of the experiment. One and only one of the arguments must be `NULL`. For example, if the goal is for the HDI to exclude the value of 0.5, then `nullVal = .5` (and `HDImaxwid = NULL`). If the goal is for the HDI to have a maximal width of .2, then `HDImaxwid = .2` (and `nullVal = NULL`). The next argument specifies the limits of the ROPE, if relevant. The function in its present form checks only if the HDI excludes the null value, or ROPE if specified. The function does not check whether the ROPE fully contains the HDI.

Finally, the function also has arguments for the audience-agreed prior that is used in the Bayesian analysis. It defaults to a uniform prior, but an audience-agreed prior can be specified in terms of the mean `audPriorMean` and certainty `audPriorN` of a beta distribution.

Here is an example of how to use the program:

```
source("minNforHDIpower.R")
```

```
minNforHDIpower( genPriorMean=.65 , genPriorN=2000 , nullVal=.5 )
```

The program will compute power for increasing values of `sampleSize`, until stopping at $N=82$ (as shown in Table 13.1).

`(minNforHDIpower.R)`

```

1 minNforHDIpower = function( genPriorMean , genPriorN ,
2                               HDImaxwid=NULL , nullVal=NULL , ROPE=c(nullVal,nullVal) ,
3                               desiredPower=0.8 , audPriorMean=0.5 , audPriorN=2 ,
4                               HDImass=0.95 , initSampSize=20 , verbose=T ) {
5   if ( is.null(HDImaxwid) + is.null(nullVal) != 1 ) {
6     stop("One and only one of HDImaxwid and nullVal must be specified.")
7   }
8   # Convert prior mean and N to a,b parameter values of beta distribution.
9   genPriorA = genPriorMean * genPriorN
10  genPriorB = ( 1.0 - genPriorMean ) * genPriorN
11  audPriorA = audPriorMean * audPriorN
12  audPriorB = ( 1.0 - audPriorMean ) * audPriorN
13  # Initialize loop for incrementing sampleSize

```

```

14 sampleSize = initSampSize
15 notPowerfulEnough = TRUE
16 # Increment sampleSize until desired power is achieved.
17 while( notPowerfulEnough ) {
18     zvec = 0:sampleSize # All possible z values for N flips.
19     # Compute probability of each z value for data-generating pr      ior.
20     pzvec = exp( lchoose( sampleSize , zvec )
21                 + lbeta( zvec + genPriorA , sampleSize-zvec + genPriorB )
22                 - lbeta( genPriorA , genPriorB ) )
23     # For each z value, compute HDI. hdiMat is min, max of HDI for ea   ch z.
24     hdiMat = matrix( 0 , nrow=length(zvec) , ncol=2 )
25     for ( zIdx in 1:length(zvec) ) {
26         z = zvec[zIdx]
27         hdiMat[zIdx,] = HDIofICDF( qbeta , shape1 = z + audPriorA ,
28                                     shape2 = sampleSize - z + audPriorB )
29     }
30     hdiWid = hdiMat[,2] - hdiMat[,1]
31     if ( !is.null( HDImaxwid ) ) {
32         powerHDI = sum( pzvec[ hdiWid < HDImaxwid ] )
33     }
34     if ( !is.null( nullVal ) ) {
35         powerHDI = sum( pzvec[ hdiMat[,1] > ROPE[2] | hdiMat[,2] < ROPE[1] ] )
36     }
37     if ( verbose ) {
38         cat( " For sample size = ", sampleSize , ", power = " , powerHDI ,
39              "\n" , sep="" ) ; flush.console()
40     }
41     if ( powerHDI > desiredPower ) {
42         notPowerfulEnough = FALSE
43     } else {
44         sampleSize = sampleSize + 1
45     }
46 } # End while( notPowerfulEnough )
47 # Return the minimal sample size that achieved the desired po      wer.
48 return( sampleSize )
49 } # end of function

```

13.6.2 Power and sample size for multiple mints

This section explains how to use a BUGS model to estimate power given sample size. The first step is to create a working BRugs program `faingle` set of data, as if the data were from an actual experiment. Make sure that the `model` has appropriate burn-in and thinning so that the posterior sample is robust and useful. For a reminder of the issues of burn-in and thinning, see Section 23.2, p. 510. Once a working BRugs program is established, it is modified and wrapped in a function that's called repeatedly with different simulated sets of data.

The general scheme for estimating power with an BRugs `program` outlined in Figure 13.3. At the top of the code is the BUGS model wrapped in `reaction` called `GoalAchievedForSample`. The function takes simulated data as input, and returns a `logical` value indicating whether or not the goal was achieved. As noted in Figure 13.3, the BUGS model specification remains unchanged. The data `spotion` is modified, however, to accept the passed-in simulated data rather than `actual` data. The MCMC chains are then generated as in the original program, using the `burn` and `thin`ning that is already known to produce well-mixed chains. Finally, when the `chain`s are examined, new commands are included to test whether the goal has been achieved. For example, the code could

```

GoalAchievedForSample = function( simulatedData )      f
  library(BRugs) # needed inside function to re-initialize.
#
#-----#
# Model specification: Unchanged.
#
#-----#
# Data: NEW COMMANDS for using simulatedData.
#
#-----#
# Initialize and run chains: Unchanged.
#
#-----#
# Examine chains: NEW COMMANDS for checking whether goal is achieved.
# Denote result by true/false variable named goalAchieved.
  return( goalAchieved )
g

nSimulatedExperiments = 500 # An arbitrary large number.
nSuccess = 0 # Initialize counter.
for ( experimentIdx in 1:nSimulatedExperiments )      f
  simulatedData = ... # Create simulated data
  nSuccess = nSuccess + GoalAchievedForSample( simulatedData )
  estPower = nSuccess / experimentIdx
g

```

Figure 13.3: Outline of R code for computing power. This ~~outline~~ assumes that there is already a working BRugs script for a single set of ~~data~~. The working BRugs script is wrapped in a function at the beginning of ~~the~~ shown here.

check whether the 95% HDI excludes the ROPE. The result ~~is~~ ~~not~~ value denoted by the variable `goalAchieved`, which gets returned by the function.

The lower section of code in Figure 13.3 calls the function `GoalAchievedForSample` in the upper section of code. The idea is that we will call the function many times each time with different simulated data, and check whether the goal is achieved. The proportion of times that the goal is achieved is the estimated power. The number of times ~~all~~ the function is denoted `nSimulatedExperiments`. Then a `for` loop tallies the number of successes. As the number of simulated experiments increases, the estimate becomes ~~stable~~.

We can think of each time through the loop, with a new `simulatedData` set, as the flip of a coin that has a probability of heads equal to the true `prior` of the experiment. We are estimating the bias in that coin, i.e., the power of the experiment, by the observed number of successes. If our prior belief about the power is `uniform`, then the posterior belief is distributed as `beta(nSuccess, nSimExperiments-nSuccess)`. The HDI of the estimated power can be computed by calling the function `qbeta` (`qbeta(quantile(rBeta, 0.95), shape1 = nSuccess+1, shape2 = nSimExperiments-nSuccess+1)`).

A complete instantiation of the scheme is shown in the `program` below. It was used to generate Figure 13.2 and the other results mentioned in Section 13.3. The code is modified from the program used to analyze the real data, step 188. The program has many embellishments beyond the skeleton outlined in Figure 13.3.

The function, that takes a set of data as input and returns ~~whether~~ whether not the goal was achieved, spans lines 5–147 of the program. The function ~~has~~ takes extra arguments

that help with plotting graphs of the results, which we might want to see for at least a few sets of data to be sure that the simulations are behaving. One extra argument is `plotResults`, which is simply a `TRUE/FALSE` argument indicating whether or not to plot the results. Another extra argument is the `lename` to be used for saving plots.

The data section, lines 44–49, is simpler than in the original program, because the `completedataList` is provided from outside the function. The data section extends the `dataList` to BUGS.

In the results sections, beginning at line 89, there are three main changes. First, the chains are not checked for convergence and mixing, because it is assumed that the burn-in and thinning have already been checked for adequacy. (For a [further](#) of the issues of burn-in and thinning, see Section 23.2, p. 510.) Second, the producing code (lines 104–132) is wrapped in an `if` statement, which causes plots to be produced and saved if the `plotResults` argument is true.

The third main change is checking for whether the goal was reached, as computed in lines 134–144. The code checks three goals regarding the `c` values: (i) Does the HDI on iteration minus condensation exceed zero? (ii) Does the difference between the two iteration conditions exceed zero? (iii) Is the HDI difference of condensation condensations of width less than 0.2? The `TRUE/FALSE` values of these three goals are returned as a vector.

The function is called by code starting at line 150. The code includes the option for computing either retrospective or replication power, which is specified on line 152. The program then loads the original data, which are used for computing replication power, and the actual posterior sample, which is used for either type of power.

Various constants are specified in lines 170–176. In particular, the simulated number of subjects per condition is specified in line 170. Also note that the number of goals computed by the function is specified in line 176. The `numGoals` is used only to initialize the success counter in line 188.

Because we are simulating results from randomly generated data, and each data set may take a few minutes to run, the complete batch of `hds` data sets may take many hours to complete. Therefore we may want to run only a few simulated experiments at a time, or interrupt a run, and save the interim results for future continuation. Also, be warned that BUGS often decides to crash for no apparent reason on some simulated data sets and you will find yourself with interrupted runs whether you intended it or not. For these reasons, the program saves the interim results for every simulated data set, on lines 241–242. When the program is invoked, it checks whether there is a previously saved interim result, on lines 178–190, and restores those results if they exist.

The loop that repeatedly generates simulated data sets for the Bayesian analysis begins on line 193. The simulated data are generated in lines 205–209. It uses the `mu` and `kappa` values from the actual posterior to generate simulated data for each condition. The simulated data are assembled into `dataList` in lines 211–229. If retrospective power is being computed, then only the newly simulated data are added in `the dataList`. If replication power is being computed, then the original `actualData` are concatenated onto the simulated data. By including the actual data, the model effectively using the actual posterior as the prior for the new simulated data.

(`FilconBrugsPower.R`)

```

1 graphics.off()
2 rm(list=ls(all=TRUE))
3
```

```

4
5 GoalAchievedForSample = function( dataList , plotResults      =F ,
6                               fileNameRoot="DeleteMe" ) {
7
8 library(BRugs)          # Kruschke, J. K. (2010). Doing Bayesian dat a analysis:
9                      # A Tutorial with R and BUGS. Academic Press / Elsevier.
10 #-----
11 # THE MODEL.
12
13 modelstring = "
14 # BUGS model specification begins here...
15 model {
16   for ( subjIdx in 1:nSubj ) {
17     # Likelihood:
18     z[subjIdx] ~ dbin( theta[subjIdx] , N[subjIdx] )
19     # Prior on theta: Notice nested indexing.
20     theta[subjIdx] ~ dbeta( a[cond[subjIdx]] , b[cond[subjIdx]] ) I(0.001,0.999)
21   }
22   for ( condIdx in 1:nCond ) {
23     a[condIdx] <- mu[condIdx] * kappa[condIdx]
24     b[condIdx] <- (1-mu[condIdx]) * kappa[condIdx]
25     # Hyperprior on mu and kappa:
26     mu[condIdx] ~ dbeta( Amu , Bmu )
27     kappa[condIdx] ~ dgamma( Skappa , Rkappa )
28   }
29   # Constants for hyperprior:
30   Amu <- 1
31   Bmu <- 1
32   Skappa <- pow(meanGamma,2)/pow(sdGamma,2)
33   Rkappa <- meanGamma/pow(sdGamma,2)
34   meanGamma <- 10
35   sdGamma <- 10
36 }
37 # ... end BUGS model specification
38 " # close quote for modelstring
39 # Write model to a file:
40 writeLines(modelstring,con="model.txt")
41 # Load model file into BRugs and check its syntax:
42 modelCheck( "model.txt" )
43
44 #-----
45 # THE DATA.
46
47 # dataList supplied from outside the function.
48 # Get the data into BRugs:
49 modelData( bugsData( dataList ) )
50
51 #-----
52 # INTIALIZE THE CHAINS.
53
54 nChain = 3
55 modelCompile( numChains=nChain )
56
57 if ( F ) {
58   modelGenInits() # often won't work for diffuse prior
59 } else {
60   # initialization based on data
61   genInitList <- function() {
62     sqzData = .01+.98*dataList$z/dataList$N

```

```

63     mu = aggregate( sqzData , list(datalist$cond) , mean )[,"x"      ]
64     sd = aggregate( sqzData , list(datalist$cond) , sd )[,"x"]
65     kappa = mu*(1-mu)/sd^2 - 1
66     return(
67       list(
68         theta = sqzData ,
69         mu = mu ,
70         kappa = kappa
71       )
72     )
73   }
74   for ( chainIdx in 1 : nChain ) {
75     modellnits( bugslnits( genInitList ) )
76   }
77 }
78 #
79 #-----#
80 # RUN THE CHAINS.
81
82 burninSteps = 1000
83 modelUpdate( burninSteps )
84 cat("Got through burn in...");flush.console()
85 samplesSet( c("mu","kappa","theta","a","b") )
86 nPerChain = ceiling(3000/nChain)
87 modelUpdate( nPerChain , thin=5 )
88
89 #-----#
90 # EXAMINE THE RESULTS.
91
92 # Extract chain values from BUGS:
93 mu = NULL
94 kappa = NULL
95 for ( condIdx in 1:nCond ) {
96   nodeName = paste( "mu[" , condIdx , "]" , sep="" )
97   mu = rbind( mu , samplesSample( nodeName ) )
98   nodeName = paste( "kappa[" , condIdx , "]" , sep="" )
99   kappa = rbind( kappa , samplesSample( nodeName ) )
100 }
101 chainLength = NCOL(mu)
102
103 # Display results if desired:
104 if ( plotResults ) {
105   # Histograms of condition (i.e. group) mu differences:
106   windows(12,4)
107   layout( matrix(1:3,nrow=1) )
108   source("plotPost.R")
109   histInfo = plotPost( mu[1,]-mu[2,] , xlab=expression(mu[1]-mu[2]) ,
110                         compVal=0.0 , breaks=30 , main="" )
111   histInfo = plotPost( mu[3,]-mu[4,] , xlab=expression(mu[3]-mu[4]) ,
112                         breaks=30 , main="" )
113   HDIlim = HDIofMCMC( mu[3,]-mu[4,] )
114   text( mean(HDIlim) , .25*max(histInfo$density) , adj=c(.5,0) , cex=1.25 ,
115         bquote( "HDI width = " * .(signif(HDIlim[2]-HDIlim[1],3)) ) )
116   nSubjPerCond = round( datalist$nSubj / datalist$nCond )
117   histInfo = plotPost( (mu[1,]+mu[2,])/2 - (mu[3,]+mu[4,]) /2 , compVal=0.0 ,
118                         xlab=expression((mu[1]+mu[2])/2 - (mu[3]+mu[4])/2) ,
119                         breaks=30 , main="" )
120   dev.copy2eps( file = paste( fileNameRoot,N",nSubjPerCo      nd,"_,explIdx,".eps" ,
121                         sep="" ) )

```



```

181 # load previous expldx, nSuccess
182 load(paste(fileNameRoot,"N",nSubjPerCond,"Result.Rd      ata",sep=""))
183 prevExpldx = expldx
184 # Use just some of the MCMC steps, distributed among the whole chain:
185 chainIdxVec = round(seq(1,chainLength,len=(prevExpldx +nSimExperiments)))
186 } else { # ... otherwise, start a new record
187   prevExpldx = 0
188   nSuccess = rep(0,nGoal) # Initialize success counter.
189   chainIdxVec = round(seq(1,chainLength,len=nSimExperiments))
190 }
191
192 # Simulated experiment loop begins here:
193 for ( expldx in (1+prevExpldx):(nSimExperiments+prevExp     ldx) ) {
194
195   # Generate random data from posterior parameters
196   chainIdx=chainIdxVec[expldx]
197   CondOfSubj = sort( rep( 1:nCond , nSubjPerCond ) )
198   nTrlOfSubj = rep( nTrlPerSubj , nSubj )
199   nCorrOfSubj = rep( 0 , nSubj )
200   for ( condIdx in 1:nCond ) {
201     m = mu[condIdx,chainIdx]
202     k = kappa[condIdx,chainIdx]
203     a = m*k
204     b = (1-m)*k
205     # Generate random theta and z values for simulated subjects:
206     thetaVec = rbeta( nSubjPerCond , a , b )
207     zVec = rbinom( n=nSubjPerCond , size=nTrlPerSubj , prob=th     etaVec )
208     nCorrOfSubj[ CondOfSubj==condIdx ] = zVec
209   }
210
211   # Put data into list for BUGS program
212   if ( analysisType == "Retro" ) { # retrospective power
213     dataList = list(
214       nCond = nCond ,
215       nSubj = nSubj ,
216       cond = CondOfSubj ,
217       N = nTrlOfSubj ,
218       z = nCorrOfSubj
219     )
220   }
221   if ( analysisType == "Repli" ) { # replication probability
222     dataList = list(
223       nCond = nCond ,
224       nSubj = nSubj + nSubjOrig ,
225       cond = c( CondOfSubj , CondOfSubjOrig ) ,
226       N = c( nTrlOfSubj , nTrlOfSubjOrig ) ,
227       z = c( nCorrOfSubj , nCorrOfSubjOrig )
228     )
229   }
230
231   # Make plots for first 10 simulated experiments:
232   if ( expldx <= 10 ) { plotRes = T } else { plotRes = F }
233
234   # Call BUGS program and tally number of successes:
235   nSuccess = nSuccess + GoalAchievedForSample( dataList ,
236                                             plotRes , fileNameRoot )
237   estPower = nSuccess / expldx
238   cat( "\n*** nSubjPerCond:",nSubjPerCond, " , nSimExp:",e     xpldx,
239        " , nSuccess:",nSuccess, " , estPower:",round(estPower,2      ), "\n\n" )

```

```

240 flush.console()
241 save( nSuccess , expldx , estPower ,
242 file=paste(fileNameRoot,"N",nSubjPerCond,"Result.Rd      ata",sep="" ) )
243
244 } # end of simulated experiment loop.

```

13.7 Exercises

Exercise 13.1. [Purpose: Comic relief.] Read the complete oeuvre of Friedrich Nietzsche, with special attention to his posthumous work *The Will to Power* (Nietzsche, 1967). Provide a mathematical formalization of the Nietzschean concept of will and power, using Bayesian probability theory. Show that the notion of static power is a special case of formalized Nietzschean power, and vice versa. Write your answer in both English and German. Do not submit your answer to the instructor; do post your answer on your personal web blog.

Exercise 13.2. [Purpose: Understanding power for flipping a single coin, in Tables 13.1 and 13.2.] For this exercise, consider flipping a single coin and investigate its bias.

(A) Table 13.2 indicates that when the data-generating distribution is vague, with $\alpha = 10$ and $\beta = :80$, then 85 ips are needed for an 80% chance of getting the 95% width to be less than .2. What is the minimum N needed if the data-generating distribution is very certain, with $\alpha = 2000$? Show the command you used, and report the exact power for the smallest N that has power greater than .8.

(B) Regarding the previous part, why might a researcher pursue a goal of precision if the data-generating hypothesis is already very precise? (The audience prior may be different than the data-generating hypothesis. Discuss briefly perhaps with an example.)

(C) Table 13.1 indicates that when the data-generating distribution is highly certain, with $\alpha = 2000$ and $\beta = :80$, then 18 ips are needed for an 80% chance of getting the 95% HDI to exclude $\mu = :5$. What is the minimum N needed if the data-generating distribution is vague, with $\alpha = 2$? Show the command you used, and report the exact power for the smallest N that has power greater than .8.

(D) For the previous part, the goal was for the HDI to exclude the value (i.e., 0.5). Notice that the goal can be satisfied if the HDI is above the value or if the HDI is below the null value. (i) When the data-generating prior is a beta distribution with $\alpha = 0:8$ and $\beta = 2$, as in the previous part, what proportion of the data-generating biases are greater than the null value? (ii) If the goal is for the HDI to fall entirely above the null value, what sample size is needed to achieve a power of .8? Hint: Use `hdiPower.R` with the argument `ROPE=c(0,.5)`. Watch the sample size increase and increase and increase, as the power creeps toward an asymptote. Why does the power exceed the proportion you computed for (i)?

Exercise 13.3. [Purpose: Power determination for groups of coins, when goals precision.] Consider the iteration-condensation experiment summarized in Sect 13.3. Suppose we want the 95% HDI on the difference, $\mu_3 - \mu_4$, to have a width less than 0.20. What sample size (N per group) is needed to achieve this goal at least 80% of the time? Determine the answer by running the program in Section 13.6 (`binBrugsPower.R`) with various values

for nPerGroup. Hint: N = 17; your job is to find the minimaN and discuss how you did it.

Exercise 13.4. [Purpose: This is a capstone exercise that uses real data to review many techniques of the previous chapters, including generating priors in BUGS, checking credibility of null values, estimating retrospective power, and conducting a posterior predictive check.]

This exercise examines a learning experiment that investigated how easy it is for people to learn new category structures after having previously learned an initial structure (Kruschke, 1996). Some new structures had relevant stimulus dimensions that were also previously relevant in the initial structure, while other new structures had relevant stimulus dimensions that were previously irrelevant in the initial structure. The initial structure, and the subsequently learned structures, are outlined in Table 13.4.

Table 13.4: Design of relevance shift experiment reported by Kruschke (1996). Cells indicate category assignment (X or O) of each stimulus. Learners were trained in the Initial Phase, and then seamlessly continued one of the Reversal, Relevant, Irrelevant or Compound Phases.

Phase	Stimulus							
	j	j	j	j	j	j	j	j
Initial	X	O	O	X	X	O	O	X
Reversal	O	X	X	O	O	X	X	O
Relevant	X	X	O	O	X	X	O	O
Irrelevant	X	X	X	X	O	O	O	O
Compound	X	O	X	O	O	X	O	X

Note. The actual stimuli had somewhat different proportions than the ones displayed here. The assignment of physical stimuli to abstract structural items was randomly permuted for each subject.

The stimuli were simple pictures of freight-train box cars that had three dimensions: position of “door” (left or right), height of box (short or tall) and color of “wheels” (blank or wheeled). The initial phase involved a structure in which all dimensions are relevant. For example, in Table 13.4, the initial phase can be described as “it’s short and left-doored, or tall and right-doored, then it’s an X”. Notice that wheelcolor is irrelevant in the initial phase. After learning that classification accurately, people would be trained on one of the four shifts listed in Table 13.4. The Reversal shift reverses the labels. The Relevant shift uses just one of the previously relevant dimensions. The example in Table 13.4 can be described as “If it’s short, then it’s an X”. The Irrelevant shift uses the one previously irrelevant dimension; the example in Table 13.4 can be described as “If it’s blank-wheeled, then it’s an X”. Finally, the Compound shift requires attention to two dimensions, one of which was previously relevant and one of which was previously irrelevant. The example in Table 13.4 can be described as “If it’s blank-wheeled and left-doored, or tall-wheeled and right-doored, then it’s an X”.

Various theories of learning predict different degrees of difficulty in learning the new structures after the old structures (see p. 230 of Kruschke 1996). For example, if learn-

ers simply memorize the eight stimuli and their assignments. ~~the Reversal requires all eight associations to be changed, whereas Relevant, and Compound shifts all require only four assignments to be changed. Therefore, reversal should be most difficult to learn, and the other three changes should be equally easy to learn.~~ One of the novel contributions of the experiment was the shift design ~~that allowed direct comparison of “reversal” shift with “intradimensional” and “extradimensional” shifts, which had not been done previously. Another novel contribution was that the ~~phase~~ had dimensions that were relevant to the outcome without being individually related with the outcome. Variations of the design have subsequently been used by other researchers (e.g., D. N. George & Pearce, 1999; Oswald et al., 2001).~~

For each shift, the number of correct responses in the ~~trial~~ was recorded for each learner. These data are shown in the left column of Figure 13.286, and are included in the `Kruschke1996CSdatsum.Rdata`. It appears that Reversal shift is easiest (i.e., has highest accuracy), followed by Relevant, Irrelevant, and Compound shifts. The Bayesian analysis will tell us the credibility of those apparent differences between conditions.

(A) [Purpose: Create a BUGS model that has an estimated hyperprior on α and on κ .] In the hierarchical diagram on the right side of Figure 9.17, p., ~~188~~ α parameters are drawn from a gamma distribution that has its parameters estimated rather than fixed. For our new model, we want to do the analogous estimation for κ parameters also. Instead of setting A and B as constants, they will be estimated hyperparameters, $a_{\text{prior}} = m - k + 1$ and $b_{\text{prior}} = (1 - m) - k + 1$, with $m \sim \text{dunif}(0, 1)$ and $k \sim \text{dgamma}(1:1)$. Here is one way to specify the model: (`Kruschke1996CSbugs.R`)

```

11 model {
12   for ( i in 1:nSubj ) {
13     nCorrOfSubj[i] ~ dbin( theta[i] , nTriOfSubj[i] )
14     theta[i] ~ dbeta( a[ CondOfSubj[i] ] , b[ CondOfSubj[i] ] )( 0.0001,0.9999 )
15   }
16   for ( cond in 1:nCond ) {
17     a[cond] <- mu[cond] * kappa[cond]
18     b[cond] <- (1-mu[cond]) * kappa[cond]
19     mu[cond] ~ dbeta( aMu , bMu )
20     kappa[cond] ~ dgamma( sGamma , rGamma )
21   }
22   aMu <- max( .01 , mMu * kMu )
23   bMu <- max( .01 , (1-mMu) * kMu )
24   mMu ~ dunif(0,1)
25   kMu ~ dgamma(1,.1)
26   sGamma <- max( .005 , pow(muGamma,2)/pow(sigmaGamma,2) )
27   rGamma <- max( .005 , muGamma/pow(sigmaGamma,2) )
28   muGamma ~ dgamma(1,.1)
29   sigmaGamma ~ dgamma(1,.1)
30 }
```

Discuss why we would want to estimate higher-level distributions across the α and κ . You may want to mention commonalities across conditions, ~~such as~~ the learners being the same species, and all learners experiencing the same ~~stimuli~~ on the same computer display. It is also important to discuss shrinkage of estimates. Finally, discuss why the particular higher-level distribution might not be appropriate.

(B) [Purpose: Check for convergence, mixing, and autocorrelation.] With the data included, run the BUGS model and check for convergence, mixing, and autocorrelation. The model can be initialized automatically, using `modelGenInits()`, but unfortunately takes forever for

the chains to converge, because some are initialized too far from the mode of the posterior. Instead, manually initialize the chains at reasonable values, as indicated by the data. The same method as was used in Section 9.5.2 (Kruschke1996CSbugs.R) is used and extended here: (Kruschke1996CSbugs.R)

```

77 genInitList <- function() {
78   sqzData = .01+.98*datalist$nCorrOfSubj/datalist$nTrlO      fSubj
79   mu = aggregate( sqzData , list(datalist$CondOfSubj) , "mean" )[, "x"]
80   sd = aggregate( sqzData , list(datalist$CondOfSubj) , "sd"      )[, "x"]
81   kappa = mu*(1-mu)/sd^2 - 1
82   mMu = mean( mu )
83   kMu = mMu * (1-mMu) / sd(mu)^2
84   muGamma = mean( kappa )
85   sigmaGamma = sd( kappa )
86   return( list( theta = sqzData ,
87                 mu = mu ,
88                 kappa = kappa ,
89                 mMu = mMu ,
90                 kMu = kMu ,
91                 muGamma = muGamma ,
92                 sigmaGamma = sigmaGamma ) )
93 }
94 for ( chainIdx in 1 : nchain ) {
95   modellnits( bugslnits( genInitList ) )
96 }
```

Run the model and determine reasonable burn-in and thinning. Show the autocorrelation and chain plots for μ_c and σ_c . You may find it useful to refer to Section 23.2, p. 510.

(C) [Purpose: Examine and interpret the posterior distribution of differences.] Which groups are different from each other, on which parameters? Hint: See Figure 23.14

(D) [Purpose: Check the robustness of the posterior when the prior is changed in reasonable ways.] Do the posterior differences change much if the prior changes? (i) Specifically, try this alternative vague prior: Wherever the top-level specification `dgamma(1,.1)`, change it to `dgamma(0.1,.1)`. (ii) Also, try this prior that forces all μ_c to be close to 15:

```

muGamma dgamma(22500,1500)
sigmaGamma dgamma(25,250)
```

Show your results. Which prior is more reasonable, and why?

(E) [Purpose: Conduct a posterior predictive check.] For this part, use the plausible `dgamma(1,.1)` prior, not the less-plausible other priors explored in the previous part. Conduct a posterior predictive check by generating simulated data from the sampled posterior parameter values. This can be done in R as follows: (Kruschke1996CSbugs.R)

```

203 nSimExps = min(500,chainLength) # number of simulated experiments
204 nSubjPerCond = sum( CondOfSubj == 1 ) # number of subjects per condition
205 nTrlPerSubj = nTrlOfSubj[1] # number of trials per subject
206 nCorrOfSubjPredMat = matrix( 0 , nrow=nSimExps , ncol=nSubjPerCond*nCond )
207 CondOfSubjPredMat = matrix( 0 , nrow=nSimExps , ncol=nSubjPerCond*nCond )
208 nTrlOfSubjPredMat = matrix( 0 , nrow=nSimExps , ncol=nSubjPerCond*nCond )
209 for ( stepIdx in 1:nSimExps ) {
210   for ( condIdx in 1:nCond ) {
211     m = mu[condIdx,stepIdx]
212     k = kappa[condIdx,stepIdx]
213     a = m*k + 1
214     b = (1-m)*k + 1
215     for ( subjIdx in 1:nSubjPerCond ) {
```

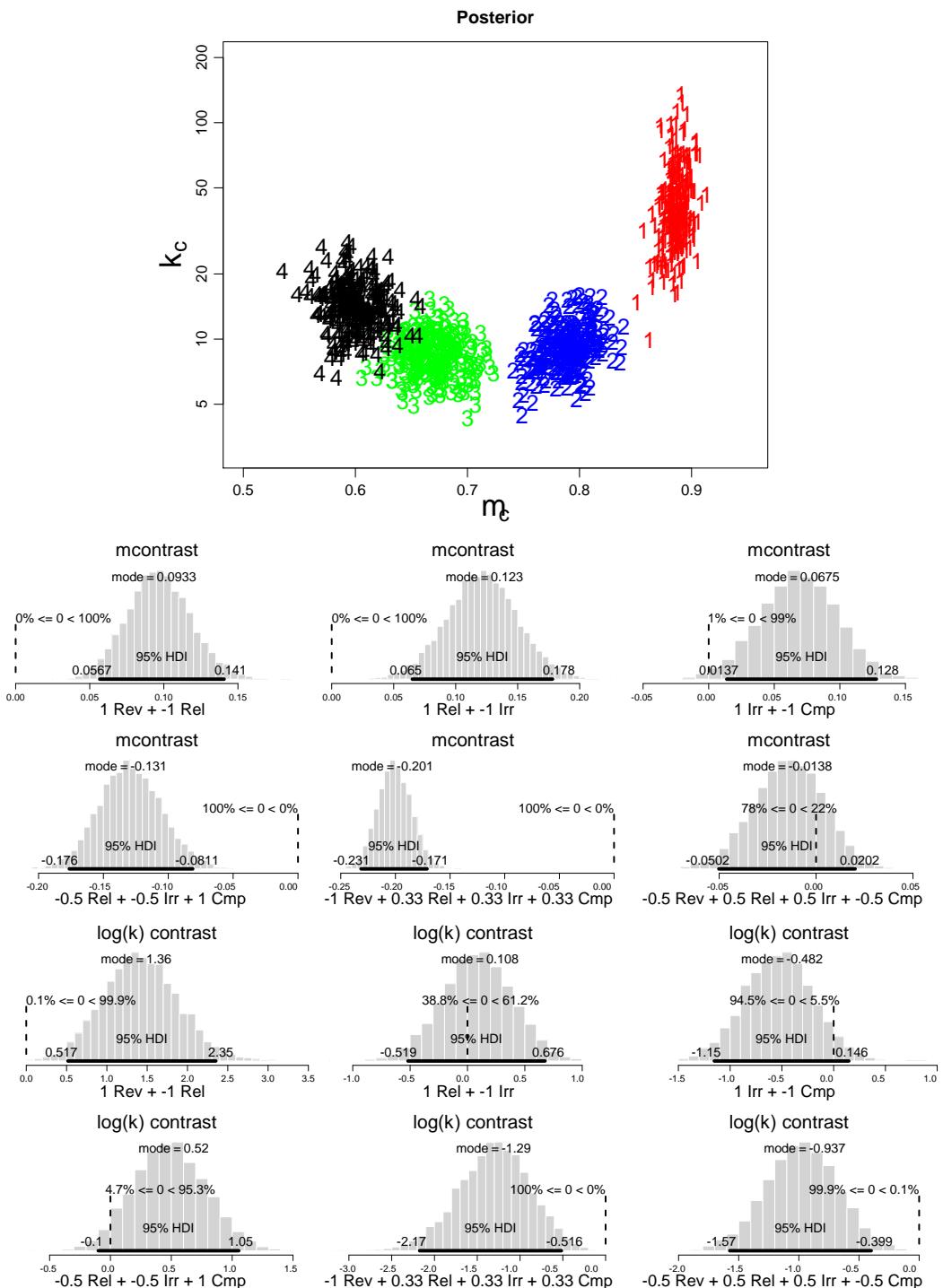


Figure 13.4: For Exercise 13.4, Part C. Posterior m_c and k_c values, and their differences. In the scatterplot, numerals represent a step in the chain. mcontrast condition.

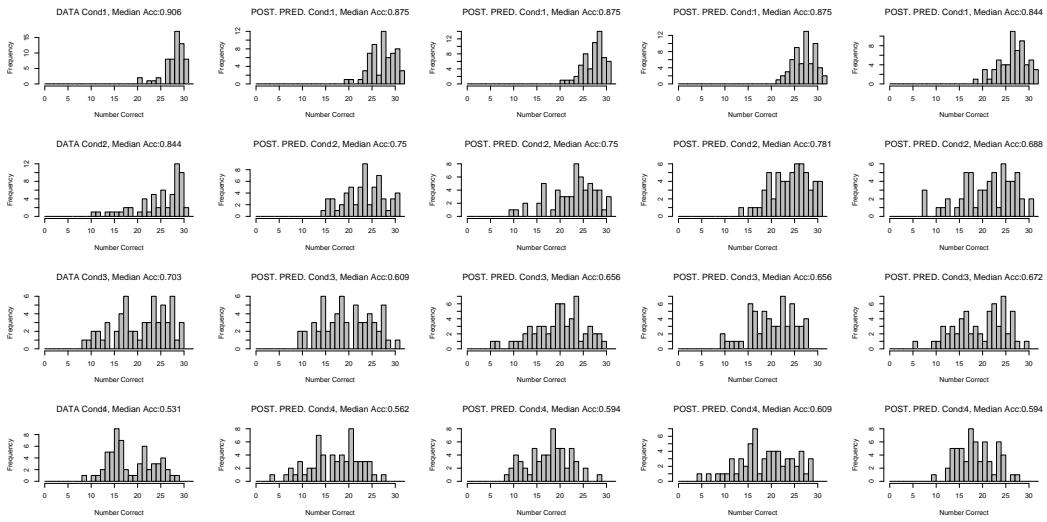


Figure 13.5: For Exercise 13.4, Part E. Actual data in `leftcon`, with examples of posterior-predicted data in other columns. Top row is `Rev`, second row is `Rel`, third row is `Irr`, and fourth row is `Compd`. There is a suggestion of bimodality in the data from the Compound Shift (`to-left histogram`) that is not often generated by simulated posterior-predicted data. Follow-up modeling and empirical work might want to investigate that bimodality if it is deemed theoretically interesting.

```

216 colIdx = (condIdx-1)*nSubjPerCond + subjIdx
217 theta = rbeta( 1 , a , b )
218 nCorrOfSubjPredMat[stepIdx,colIdx] = rbinom( 1 ,
219 size=nTriPerSubj , prob=theta )
220 CondOfSubjPredMat[stepIdx,colIdx] = condIdx
221 nTriOfSubjPredMat[stepIdx,colIdx] = nTriPerSubj
222 }
223 }
224 }
```

This code, above, puts each randomly generated sample `atFind` at a row of the matrix `nCorrOfSubjPredMat`. Figure 13.5 shows histograms of simulated data. The above code relied on previously extracting the parameter chain from fBugs, as follows: (Kruschke1996CSbugs.R)

```

127 mu = NULL
128 kappa = NULL
129 for ( condIdx in 1:nCond ) {
130   nodeName = paste( "mu[" , condIdx , "]" , sep="" )
131   mu = rbind( mu , samplesSample( nodeName ) )
132   nodeName = paste( "kappa[" , condIdx , "]" , sep="" )
133   kappa = rbind( kappa , samplesSample( nodeName ) )
134 }
```

(F) [Purpose: Do a retrospective power analysis.] Using the method of Figure 13.3, which is exemplified by the code explained in Section 13.6 (`conBrugsPower.R`), conduct a retrospective power analysis, and estimate the power of `foals`: $\text{Rev} - \text{Rel} > 0$, $\text{Rel} - \text{Irr} > 0$, $\text{Irr} - \text{Com} > 0$, and $\text{Rev} - \text{Rel} > 0$ (condition 1 is Reversal, condition 2 is Relevant, condition 3 is Irrelevant, and condition 4 is Compound). To do this, first analyze the actual

data, and get a posterior sample of c values. Then, wrap the BUGS program into a function, into which you can pass data and from which you check whether each goal is achieved. Then make a loop that goes stepwise through the chain and generates simulated data from those credible parameter values. At each step, pass the simulated data into the BUGS analysis function. Run at least 100 simulated steps, and report your tally for how many times the goals were achieved. Cautions and hints: If each BUGS analysis takes two minutes to run, a batch of a few hundred simulated steps will take hours. Plan accordingly. Sometimes BUGS will run fine for dozens of simulated data sets, and then inexplicably crash on the next random data set. Therefore, be sure that your results regarding power estimation are saved at the end of each data set. The goal of showing

$\text{Rev} - \text{Rel} > 0$ has a retrospective power of about 100%, but at least one of the other goals involving c has retrospective power under 50%, which indicates that with 60 subjects per condition, this experiment might have been underpowered. Bayesian analysis does not care whether there are equal numbers of subjects in each condition, so follow-up work could use more subjects in some conditions and fewer subjects in other conditions. The goal of showing $\text{Rev} - \text{Rel} > 0$ has a retrospective power of only 20%, which is caused by large shrinkage of c estimates.

(G) [Purpose: Do a replication power analysis.] Suppose you repeat the experiment, using the posterior of the first experiment as the prior for the second experiment. Estimate the probability of achieving these four goals: $\text{Rev} - \text{Rel} > 0$, $\text{Rel} - \text{Irr} > 0$, $\text{Irr} - \text{Com} > 0$, and $\text{Rev} - \text{Rel} > 0$ (condition 1 is Reversal, condition 2 is Relevant, condition 3 is Irrelevant, and condition 4 is Compound). In other words, do an analysis of replication probability, regarding the same goals as the previous part. This task is easy to do if you successfully accomplished the previous part. In each step of the chain, instead of using only the simulated data, concatenate the simulated data to the actual data. Cautions and hints: The runs will take even longer than in the previous part because of the larger data sets.

Part III

The Generalized Linear Model

Chapter 14

Overview of the Generalized Linear Model

Contents

14.1	The generalized linear model (GLM)	292
14.1.1	Predictor and predicted variables	292
14.1.2	Scale types: metric, ordinal, nominal	293
14.1.3	Linear function of a single metric predictor	294
14.1.3.1	Reparameterization x threshold form	296
14.1.4	Additive combination of metric predictors	296
14.1.4.1	Reparameterization x threshold form	298
14.1.5	Nonadditive interaction of metric predictors	298
14.1.6	Nominal predictors	300
14.1.6.1	Linear model for a single nominal predictor	303
14.1.6.2	Additive combination of nominal predictors	302
14.1.6.3	Nonadditive interaction of nominal predictors	303
14.1.7	Linking combined predictors to the predicted	304
14.1.7.1	The sigmoid (a.k.a. logistic) function	305
14.1.7.2	The cumulative normal (a.k.a. Phi) function	307
14.1.8	Probabilistic prediction	308
14.1.9	Formal expression of the GLM	308
14.2	Cases of the GLM	311
14.2.1	Two or more nominal variables predicting frequency	313
14.3	Exercises	315

Straight and proportionate, deep in your core

All is orthogonal, ceiling to oor.

But on the outside the vines creep and twist

'round all the parapets shrouded in mist.

The previous part of the book explored all the basic concepts of Bayesian analysis applied to a simple likelihood function, namely the Bernoulli distribution. The focus on a simple likelihood function allowed the complex concepts of Bayesian analysis, such as MCMC methods and hierarchical priors, to be developed ~~without~~ inference from additional complications of elaborate likelihood functions with multiple parameters.

Kruschke, J. K. (2010) Doing Bayesian Data Analysis: A Tutorial with R and BUGS. Academic Press Elsevier. Copyright 2010 by John K. Kruschke. Draft of May 11, 2010. Please do not distribute this preliminary draft. If you report Bayesian analyses based on this book, please do cite it! “

In this new part of the book, we apply all the concepts to a more complex but versatile model known as the generalized linear model (GLM; McCullagh & Nelder, 1989; Nelder & Wedderburn, 1972). This model comprises the traditional “off the shelf” analyses such as t-tests, analysis of variance (ANOVA), multiple linear regression, logistic regression, etc. Because we now know from previous chapters the concepts and mechanisms of Bayesian analysis, we can focus on applications of this versatile tool. The present chapter is important for understanding subsequent chapters because it outlines the framework for all the models in the remainder of the book.

14.1 The generalized linear model (GLM)

To understand the generalized linear model and its many species, we must build up a variety of component concepts regarding relationships between variables and how variables are measured in the first place.

14.1.1 Predictor and predicted variables

Suppose we want to predict someone's weight from their height. In this case, weight is the predicted variable and height is the predictor. Or, suppose we want to predict high school grade point average (GPA) from Scholastic Aptitude Test (SAT) score and family income. In this case, GPA is the predicted variable, while SAT and income are predictor variables. Or, suppose we want to predict the blood pressure of patients who either take drug A, or take drug B, or take a placebo, or merely wait. In this case, predicted variable is blood pressure, and treatment-group membership is the predictor.

The key mathematical difference between predictor and predicted variables is that the likelihood function expresses the probability of values of the predicted variable as a function of values of the predictor variable. The likelihood function does not describe the probabilities of values of the predictor variable. The value of the predictor variable comes from “outside” the system being modeled, whereas the value of the predicted variable depends on the value of the predictor variable.

Because the predicted variable depends on the predictor, at least mathematically in the likelihood function if not causally in the real world, the predicted variable can also be called the “dependent” variable. The predictor variables are sometimes called “independent” variables. The key conceptual difference between independent and dependent variable is that the value of the dependent variable depends on the value of the independent variable. The term “independent” can be confusing because it can be used strictly or loosely. In experimental settings, the variables that are really manipulated and set by the experimenter are the independent variables. In this context of experimental manipulation, the values of the independent variables truly are (in principle at least) independent of the values of other variables, because the experimenter has control to arbitrarily set the values of the independent variables. But sometimes a non-manipulated variable is also referred to as “independent”, merely as a way to indicate that it is treated as a predictor variable.

Among non-manipulated variables, the roles of predictor and predictor are arbitrary, determined only by the interpretation of the analysis. Consider for example, people's weights and heights. We could be interested in predicting a person's weight from his/her height, or we could be interested in predicting a person's height from his/her weight. Prediction is merely a mathematical dependency, not necessarily a picture of underlying causal

relationship. Although height and weight tend to co-vary as people, the two variables are not directly causally related. When a person slouches by getting shorter, he does not lose weight. And when a person drinks a glass of water, by weighing more, he does not get taller.

Just as “prediction” does not imply causation, “prediction” also does not imply any temporal relation between the variables. For example, we want to predict a person’s sex, male or female, from his/her height. Because males tend to be taller than females, this prediction can be made with better than chance accuracy. Peterson’s sex is not caused by his/her height, nor does a person’s sex occur only after the height is measured. Thus, we can “predict” a person’s sex from his/her height, but this does not mean that the person’s sex occurred later in time than his/her height.

In summary: All manipulated independent variables are predicted variables, not predicted. Some dependent variables can take on the role of predicted variables, if desired. All predicted variables are dependent variables. The likelihood function specifies the probability of values of the predicted variables as a function of values of the predictor variables.

Why we care. We care about these distinctions between predicted and predicted variables because the likelihood function is a mathematical expression of the dependency of the predicted variable on the predictor variable. The thing we have to do in statistical inference is identify what variables we are interested in predicting, on the basis of what predictors.

14.1.2 Scale types: metric, ordinal, nominal

Items can be measured on different scales. For example, the participants in a foot race can be measured either by the time they took to run the race, by placing in the race (1st, 2nd, 3rd, etc.), or by the name of the team they represent. These measurements are examples of metric, ordinal, and nominal scales, respectively (Stevens, 1946).

Examples of metric scales include response time (i.e., latency or duration), temperature, height, and weight. Those are actually cases of a specific type of metric scale, called ratio scale, because they have a natural zero point on the scale. A zero point on the scale corresponds to there being a complete absence of the thing measured. For example, when the duration is zero, there has been no time elapsed; when the weight is zero, there is no downward force. Because these scales have a natural zero point, it is meaningful to talk about ratios of amounts being measured, and that is why they are called ratio scales. For example, it is meaningful to say that taking 2 minutes to solve a problem is twice as long as taking 1 minute to solve the problem. On the other hand, a scale of historical time has no known absolute zero. We cannot say, for example, there is twice as much time in January 2nd as there is time in January 1st. We cannot define duration since some arbitrary reference point, but we cannot talk about the absolute amount of time in any given moment. Scales that have no natural zero are called interval scales because all we know about them is the amount of time in an interval on the scale, not the amount of time at a point on the scale. Despite the conceptual difference between ratio and interval scales, we will lump them together into the category of metric scales.

A special case of metric-scaled data is count data, also called frequency data. For example, the number of cars that pass through an intersection in an hour is a count. The number of poll respondents who say they belong to a particular political party is a count. Count data can only have values that are non-negative integers. Distances between counts have meaning, and therefore the data are metric, but because the data cannot be negative

and are not continuous, they are treated with different mathematical forms than continuous, real-valued metric data.

Examples of ordinal scales include placing in a race, or rating of degree of agreement. When we are told that, in a race, Jane came in first, Jill came second, and Jasmine came in third, we only know the order. We do not know whether Jane beat Jill by a nose or by a mile. There is no distance or metric information in an ordinal scale. As another example, many polls have ordinal response scales: Indicate how much agree with this statement: "Bayesian statistical inference is better than null hypothesis significance testing", with 5 = strongly agree, 4 = mildly agree, 3 = neither agree nor disagree, 2 = mildly disagree, and 1 = strongly disagree. Notice that there is no metric information in the response scale, because we cannot say the difference between ratings of 5 and 4 is the same amount of difference as between ratings of 4 and 3.

Examples of nominal a.k.a. categorical scales include political party affiliation, the face of a rolled die, and the result of a flipped coin. For nominal scales, there is neither distance between categories nor order between categories. For example, suppose we measure the political party affiliation of a person. The categories of the scale might be Green, Democrat, Republican, Libertarian, and Other. While some political theories might infer that the parties fall on some underlying liberal-conservative scale, there is no such scale in the actual categorical values themselves. In the actual categorical labels there is neither distance nor ordering.

In summary, if two items have different nominal values, all we know is that the two items are different (and what categories they are in). On the other hand, if items have different ordinal values, we know that the two items are different and we know which one is "larger" than the other, but not how much larger. If two items have different metric values, then we know that they are different, which one is larger, and how much larger.

Why we care. We care about the scale type because the likelihood function must specify a probability distribution on the appropriate scale. If the scale has two nominal values, then a Bernoulli likelihood function may be appropriate. If the scale is metric, then a normal distribution may be appropriate as a likelihood function. Whenever we are choosing a model for data, we must answer the question, What kind of scale are we dealing with?

In the following sections, we will first consider the case of metric predicted variable with metric predictors. In that context of all metric variables, we will develop the concepts of linear functions and interactions. Once those concepts are established for metric predictors, the notions will be extended to nominal predictors.

14.1.3 Linear function of a single metric predictor

Suppose we have identified one variable to be predicted, which we'll call y , and one variable to be the predictor, which we'll call x . Suppose we have determined that both variables are metric. The next issue we need to address is how to model the relationship between x and y . There are many possible dependencies of y on x , and the particular form of the dependency is determined by the specific meanings and nature of the variables. But in general, across all possible domains, what is the most basic or simplistic dependency on x that we might consider? The usual answer to this question is, a linear relationship. A linear function is the generic, "vanilla", off-the-shelf dependency that is used in statistical models. Methods can be generalized to other models when needed.

Linear functions preserve proportionality. If you double the input, then you double the output. If cost of a book is a linear function of the number of pages, then when the number

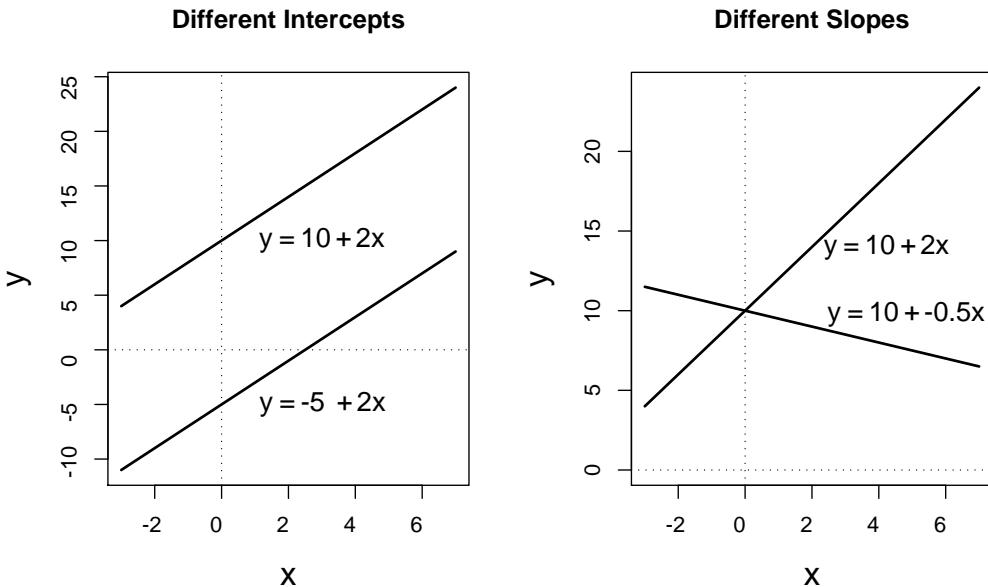


Figure 14.1: Examples of linear functions of a single variable. The left panel shows examples of two lines with the same slope but different intercepts. The right panel shows examples of two lines with the same intercept but different slopes.

of pages is reduced 10%, the cost should be reduced 10%. ~~of~~ ~~the~~ ~~the~~ speed is a linear function of gas delivery to the engine, then when you press the pedal 20% further, the car should go 20% faster. Non-linear functions do not preserve proportionality. For example, in actuality, car speed is not a linear function of gas delivery. At higher and higher speeds, it takes proportionally more and more gas to make the car ~~stop~~. Despite the fact that many real-world dependencies are non-linear, most are ~~not~~ approximately linear over moderate ranges of the variables. For example, if you have ~~to~~ paint the wall area, it takes approximately twice the amount of paint. It is also the case that linear relationships are intuitively prominent (Brehmer, 1974; Honan, Earle, & Slovic, 1981; Kalish, Griths, & Lewandowsky, 2007). Linear relationships are the easiest to think about: Turn the steering wheel twice as far, and the car should turn twice as sharp; turn the volume knob 50% higher, the loudness should increase 50%.

The general mathematical form for a linear function of a single variable is

$$y = \beta_0 + \beta_1 x \quad (14.1)$$

When values of x and y that satisfy Equation 14.1 are plotted, they form a line. ~~Examples~~ are shown in Figure 14.1. The value of parameter β_0 is called the y -intercept because it is the value where the line intersects the y -axis when $x = 0$. The left panel of Figure 14.1 shows two lines with different y -intercepts. The value of parameter β_1 is called the slope because it indicates how much y increases when x increases by 1. The right panel of Figure 14.1 shows two lines with the same y -intercept but different slopes.

In strict mathematical terminology, the type of transformation in Equation 14.1 is called a linear function. When $\beta_0 = 0$, the transformation does not preserve proportionality. Example, consider $y = 10 + 2x$. When x is doubled from $x = 1$ to $x = 2$, y increases from $y = 12$ to

$y = 14$, which is not doubling. Nevertheless, the rate of increase is the same for all values of x : Whenever x increases by 1, y increases by 2. Equation 14.1 can be algebraically re-arranged so that it does preserve proportionality, as shown next.

14.1.3.1 Reparameterization to threshold form

Equation 14.1 can be algebraically re-arranged as follows:

$$\begin{aligned} y &= \beta_0 + \beta_1 x \\ &= \beta_1 x + |\beta_0| \end{aligned} \quad (14.2)$$

This form of the equation is useful because it explicitly shows the value of the x -intercept, a.k.a. threshold, denoted (Do not confuse this use of the symbol with different uses in previous chapters.) The threshold is the value of y when x is zero. This is sometimes also called the x intercept.

The x -threshold form preserves proportionality for x . As an example, consider again the case of $y = 10 + 2x$. When changed to threshold form, it becomes $y = 2(x + 5)$. When x changes from 1 to 2, $x + 5$ changes from 6 to 7, which is an increase of $(7 - 6) = 1 = 1$. The resulting change in y is from 12 to 14, which is an increase of $(14 - 12) = 2 = 1$. Thus, a 1% increase in x results in a 1% increase in y .

The threshold (i.e. x intercept) is often more meaningful than the x -intercept. For example, suppose we are piloting a tugboat upstream on the Mississippi river, and we want to predict how much headway we will gain against the current for a given setting of the throttle x . Suppose it is the case that $y = -2 + 4x$. This form of the equation indicates that when we apply zero engine power, that is when $x = 0$, then we lose 2 miles an hour, i.e., $y = -2$. In other words, the x intercept tells us the baseline speed of the river current that we are trying to overcome. What may be more useful to know, however, is the amount of engine power we need to apply in order to overcome the current. How big must x be so that we are just matching the downstream pressure? That is, this question is the threshold, i.e., the value of x that makes $y = 0$. In our example, where $y = -2 + 4x$, the threshold is $x = (-2 - 0) / 4 = 0.5$. In other words, when the throttle is set above the threshold of 0.5, then we make progress upstream because $y > 0$, but when the throttle is set below the threshold of 0.5, we drift downstream because $y < 0$. Thus, the more intuitive form of the “headway” equation is the intercept form $y = 4(x - 0.5)$, because it shows explicitly that our headway is proportional to how much the throttle exceeds 0.5.

Summary of why we care. The likelihood function specifies the form of the dependency of y on x . When y and x are metric variables, the simplest form of dependency, both mathematically and intuitively, is one that preserves proportionality. The mathematical expression of this relation is a so-called linear function. The usual mathematical expression of a line is the intercept form, but often a more intuitive expression is the threshold form. Linear functions form the core of most statistical models. It is important to become facile with their algebraic forms and graphical representations.

14.1.4 Additive combination of metric predictors

If we have more than one predictor variable, what function do we use to combine the influences of all the predictor variables? If we want the combination to be linear in each

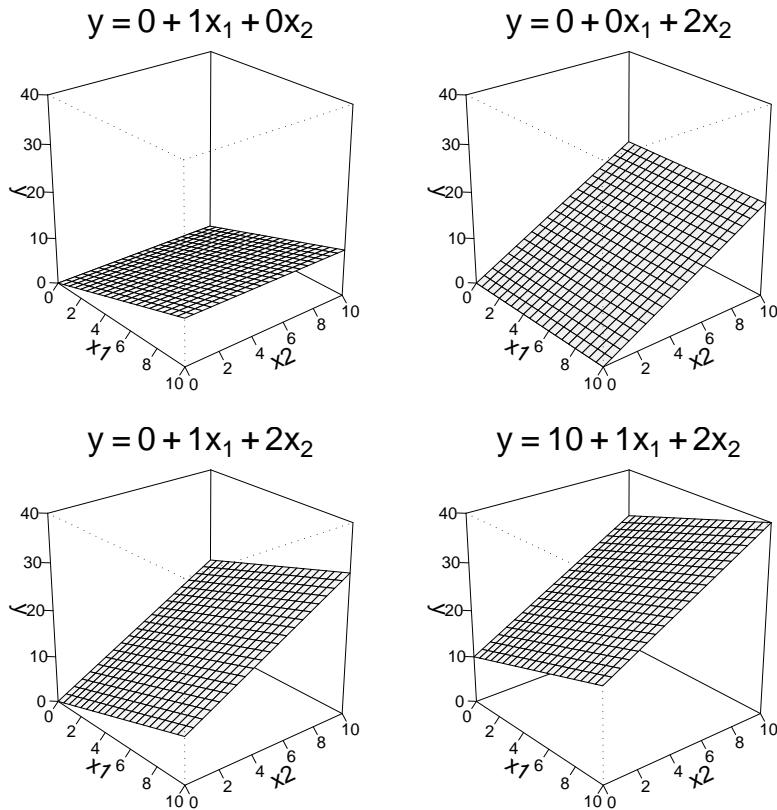


Figure 14.2: Examples of linear functions of two variables x_1 and x_2 . Upper left: Only x_1 has an influence on y . Upper right: Only x_2 has an influence on y . Lower left: x_1 and x_2 have an additive influence on y , compare with upper panels. Lower right: Non-zero intercept is added; compare with lower-left panel.

of the predictor variables, then there is just one answer to the question. In other words, if we want an increase in one predictor variable to predict the same proportional increase in the predicted variable for any value of the other predictor variables, then the predictions of the individual predictor variables must be added.

In general, a linear combination of predictor variables has the form

$$\begin{aligned} y &= \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k \\ &= \beta_0 + \sum_{k=1}^K \beta_k x_k \end{aligned} \tag{14.3}$$

Figure 14.2 shows examples of linear functions of two variables, x_1 and x_2 . The graphs show y plotted only over a domain with $0 \leq x_1 \leq 10$ and $0 \leq x_2 \leq 10$. It is important to realize that the plane extends from minus to plus infinity. The graphs only show a small region. Notice in the upper-left panel, where $y = 0 + 1x_1 + 0x_2$, that when $x_1 = 10$, then $y = 10$, regardless of the value of x_2 . The plane tilts upward in the x_1 direction, but the plane is horizontal in the x_2 direction. The opposite is true in the upper-right panel. The plane tilts upward in the x_2 direction, but the plane is horizontal in the x_1 direction, because there $y = 0 + 0x_1 + 2x_2$. The lower-left panel shows the two influences added: $y = 0 + 1x_1 + 2x_2$. Notice that the slope in the x_2 direction is steeper than in the x_1 direction. Most importantly,

notice that the slope in the direction is the same at any specific value of x_2 . For example, when $x_1 = 0$, y rises from $y = 0$ to $y = 20$, i.e. an increase of 20, when x_2 goes from $x_2 = 0$ to $x_2 = 10$. And when $x_1 = 10$, y rises from $y = 10$ to $y = 30$, again an increase of 20, when x_2 goes from $x_2 = 0$ to $x_2 = 10$.

14.1.4.1 Reparameterization to threshold form

For notational convenience, define the length of a vector $\mathbf{h} = h_1, \dots, h_k$ to be $P_k^{h_k^2} = 1^2$. This may look complicated, but it's merely the everyday ~~formula~~ for distance from the Pythagorean theorem. Carpenters all memorize ~~the~~ a special case of this relationship, known as the "3-4-5 rule": When the lengths of the edges of a right-angled triangle are 3 and 4, then the length of the long edge is ~~exact~~ because $5 = (3^2 + 4^2)^{1/2}$. (Carpenters actually use the 3-4-5 rule to infer angle ~~lengths~~, rather than infer a third length from two lengths and an angle. Carpenters know if a triangle has edge lengths of 3, 4, and 5, then the angle between the short edges is exactly 90 degrees.) With this new notation for length, Equation 14.3 can be algebraically-expressed as

$$\begin{aligned} y &= \underset{k=1}{\overset{K}{+}} h_k x_k \\ &= \underset{k=1}{\overset{K}{+}} \sqrt{h_k^2} x_k \\ &= \underset{k=1}{\overset{K}{+}} \sqrt{z_k} x_k \quad | \quad \underset{0}{\overset{1}{\{z\}}} \end{aligned} \tag{14.4}$$

Notice that when there is only a single predictor variable, when $K = 1$, then $h_1 = j_1$ and Equation 14.4 reduces to Equation 14.2.

In Equation 14.4, the value of y is the (Euclidean) length of when $y = 0$ and when x is in the direction of vector h_1, \dots, h_k . In other words, when $x = \underset{0}{\overset{1}{\{z\}}}$, then $y = 0$. When the length of x (in that direction) exceeds the threshold $y > 0$. The threshold form in Equation 14.4 becomes especially useful when we ~~den~~ logistic regression in future chapters.

Summary of section: When the influence of every individual predictor is unchanged by changing the values of other predictors, then the influence is additive. The combined influence of two or more predictors can be additive even if the individual influences are nonlinear. But if the individual influences are linear, and the combined influence is additive, then the overall combined influence is also linear. The form of Equation 14.3, or its reparameterization in Equation 14.4, is known as the ~~linear~~ model. It forms the core of many statistical models.

14.1.5 Nonadditive interaction of metric predictors

The combined influence of two predictors does not have to be additive. Consider, for example, a person's self-rating of happiness, predicted ~~by~~ ~~from~~ ~~other~~ overall health and annual

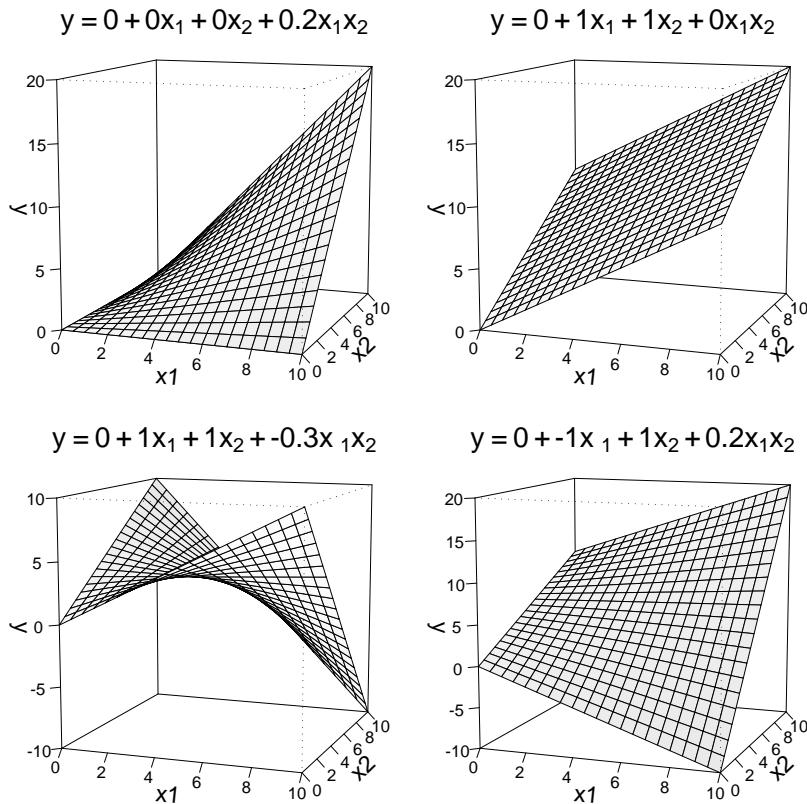


Figure 14.3: Multiplicative interaction of two variables, x_1 and x_2 . Upper right panel shows zero interaction, for comparison. Figure 14.3, provides additional perspective and insight.

income. It's likely that if a person's health is very poor~~the~~ the person is not happy, regardless of ~~his~~ income. And if the person has zero income, then the ~~person~~ probably not happy, regardless of ~~his~~ health. But if the person is both healthy and rich, then the person is probably happy (despite celebrated counter-examples in the popular media).

A graph of this sort of non-additive interaction between ~~predictors~~ appears in the upper left panel of Figure 14.3. The vertical axis, labeled y , represents happiness. The horizontal axes, x_1 and x_2 , are health and income. Notice that if either $x_1 = 0$ or $x_2 = 0$, then $y = 0$. But if both $x_1 > 0$ and $x_2 > 0$, then $y > 0$. The specific form of interaction plotted here is multiplicative $y = 0 + 0x_1 + 0x_2 + 0.2x_1x_2$. For comparison, the upper-right panel of Figure 14.3 shows a non-interactive (i.e., additive) combination of x_1 and x_2 . Notice that the graph of the interaction has a twist in it, but the graph of the additive combination is ~~flat~~.

The lower-left panel of Figure 14.3 shows a multiplicative interaction in which the individual predictors increase the outcome, but the combined variables decrease the outcome. A real-world example of this occurs with some drugs: Individually, each of two drugs might improve symptoms, but when taken together, the two drugs interact and cause a decline in health. As another example, consider lighter-than-air travel, i.e., ballooning. The buoyancy of a balloon is increased by ~~re~~, as in hot air balloons. And the buoyancy of a balloon is increased by hydrogen, as in many early-20th century blimps and dirigibles. But the buoyancy of a balloon is dramatically decreased by the combination of ~~methane~~ and hydrogen.

The lower-right panel of Figure 14.3 shows a multiplicative interaction in which the direction of influence of one variable depends on the magnitude of the other variable. Notice that when $x_2 = 0$, then an increase in the variable leads to a decline in y . But when $x_2 = 10$, then an increase in the variable leads to an increase in y . Again, the graph of the interaction shows a twist and is not straight.

A non-additive interaction of predictors does not have to be multiplicative. Other types of interaction are possible. The type of interaction is motivated by idiosyncratic theories in different variables in different application domains. Consider, for example, predicting the magnitude of gravitational force between two objects, given three predictor variables: mass of object one, mass of object two, and the distance between the objects. The force is proportional to the product (i.e., multiplication) of the two masses. But the force is proportional to the masses divided by the squared distance between them.

14.1.6 Nominal predictors

14.1.6.1 Linear model for a single nominal predictor

The previous sections assumed that the predictor was metric. What if the predictor is nominal, such as political party affiliation or gender? What is the simplest generic model for a metric variable predicted from a nominal variable? We can consider the overall average height across both sexes as the baseline height. When an individual has the value “male”, that adds an upward deviation to the predicted height. When an individual has the value “female”, that adds a downward deviation to the predicted height.

Expressing that idea in mathematical notation can get tricky. First consider the nominal predictor. We can't represent it appropriately as a single scalar value, such as 1 through 5, because that would mean that level 1 is closer to 2 than it is to level 5, which is not true of nominal values. Therefore, instead of representing the value of the nominal predictor by a single scalar value, we will represent the nominal predictor by a vector $x = [x_1; \dots; x_J]$, where J is the number of categories that the predictor has. When an individual has level j of the nominal predictor, this is represented by setting $x_{i,j} = 1$ and $x_{i,j} = 0$ for all other i . For example, suppose sex, with level 1 being “male” and level 2 being female (so $J = 2$). Then “male” is represented as $x = [1; 0]$ and “female” is represented as $x = [0; 1]$. As another example, suppose that the predictor is political affiliation, with Green as level 1, Democrat as level 2, Republican as level 3, Libertarian as level 4, and Other as level 5. Then Democrat is represented as $x = [0; 1; 0; 0; 0]$, and Libertarian is represented as $x = [0; 0; 0; 1; 0]$. Political party affiliation is being treated here as a categorical label only, with no ordering along a liberal-conservative scale.

Now that we have a formal representation for the nominal predictor variable, we can create a formal representation for the generic model of how the predictor influences the predicted variable. As mentioned above, the idea is that there is a baseline level of the predicted variable, and each category of the predictor adds a deviation above or below that baseline level. We will denote the baseline value of the prediction as β_0 . The deviation

¹Division by squared distance can be thought of as multiplying by the reciprocal of squared distance, but that amounts to re-parameterizing the model.

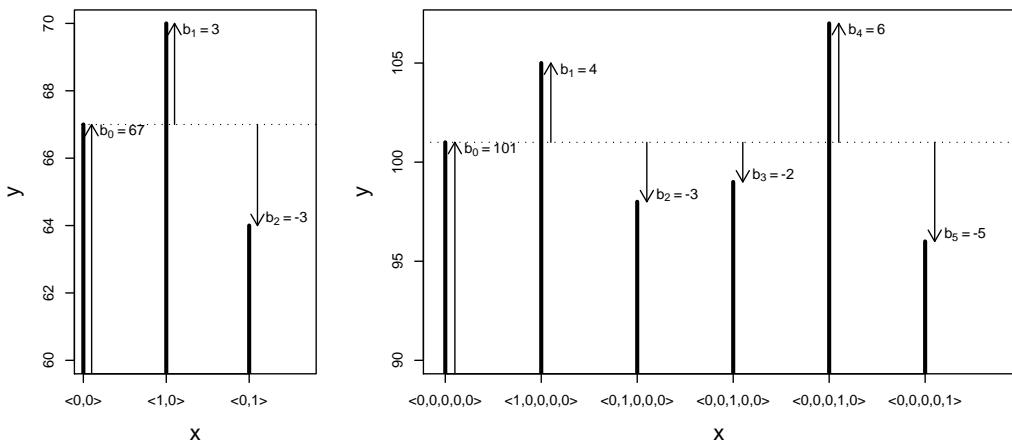


Figure 14.4: Examples of a nominal predictor (Equations 14.5 and 14.6). Left panel shows a case with $k=2$, right panel shows a case with $k=5$. In each panel, the baseline value y_0 is on the far left, when all the components are zero. Notice that the deviations from baseline sum to zero.

for the j^{th} level of the predictor is denoted. Then the predicted value is

$$\begin{aligned} y &= \beta_0 + \beta_1 x_1 + \dots + \beta_j x_j \\ &= \beta_0 + \beta^T x \end{aligned} \quad (14.5)$$

where the notation $\beta^T x$ is sometimes called the “dot product” of the vectors.

Notice that Equation 14.5 has a form very similar to the linear form of Equation 14.1. The conceptual analogy is this: In Equation 14.1 for a metric predictor, the slope β_1 indicates how much y changes when x changes from 0 to 1. In Equation 14.5 for a nominal predictor, the coefficient β_j indicates how much y changes when x changes from neutral to category j .

There is one more consideration when expressing the in ~~the~~ nominal predictor as in Equation 14.5: How should the baseline value be set? For example, predicting height from sex. We could set the baseline height to be zero. The deviation from baseline for male might be 5'10" (say), and the deviation from baseline for female might be 5'4" (say). On the other hand, we could set the baseline height to be 5'7". Then the deviation from baseline for male would be 3", and the deviation from baseline for female would be -3". The second way of setting the baseline is the typical way done in generic statistical modeling. In other words, the baseline is ~~chosen~~ so that the deviations sum to zero across the categories:

$$\sum_{j=1}^k \beta_j = 0 \quad (14.6)$$

The expression of the model in Equation 14.5 is not complete without the constraint in 14.6.

Figure 14.4 shows examples of a nominal predictor, expressed in terms of Equations 14.5 and 14.6. The left panel shows a case for which $k=2$, and the right panel

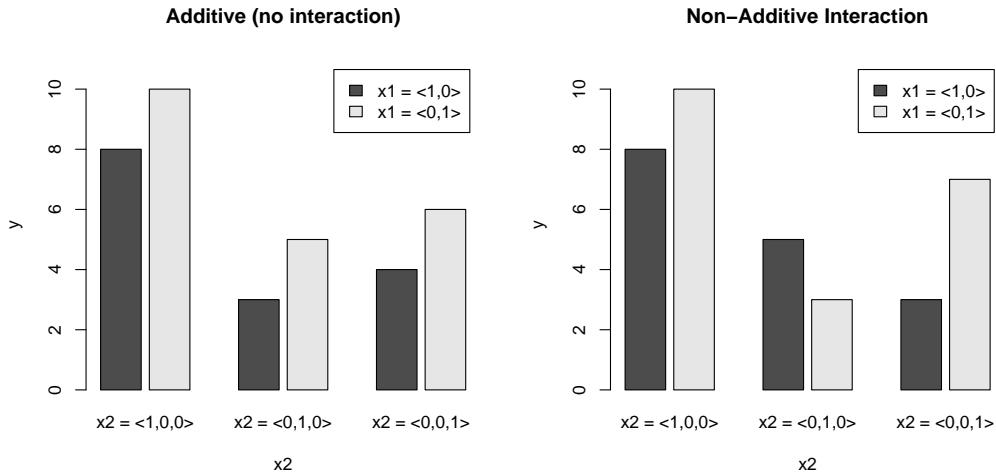


Figure 14.5: Combinations of two nominal variables. Left: Additive combination. Notice that the difference between adjacent dark-grey and light-grey bars is the same for every level of x_2 . Right: Non-additive interaction. Notice that the difference between adjacent dark-grey and light-grey bars is not the same for every level of x_2 . Figure 19.1, p. 423, provides additional perspective on this.

shows a case in which $\beta = 5$. Notice that the deviations from baseline sum to zero, as demanded by the constraint in Equation 14.6.

14.1.6.2 Additive combination of nominal predictors

Suppose we have two (or more) nominal predictors of a metric. For example, we might be interested in predicting income as a function of political party affiliation and gender. Figure 14.4 showed examples of each of those predictors individually. What we do now is consider the joint influence of those predictors. If the influences are merely additive, then the model from Equation 14.5 becomes

$$\begin{aligned} y &= \beta_0 + \beta_1 x_{1j} + \beta_2 x_{2j} \\ &= \beta_0 + \sum_{j=1}^J \beta_{1j} x_{1j} + \sum_{j=1}^J \beta_{2j} x_{2j} \end{aligned} \quad (14.7)$$

with the constraints

$$\sum_{j=1}^J \beta_{1j} = 0 \quad \text{and} \quad \sum_{j=1}^J \beta_{2j} = 0 \quad (14.8)$$

The left panel of Figure 14.5 shows an example of two nominal predictors that have additive effects on the predicted variable. In this case, the overall bias is 6. When $x_1 = h1; 0i$, there is a deviation of -1, and when $x_1 = h0; 1i$, there is a deviation of +1. This deviation by x_1 is the same at every level of x_2 . The deviations for the three levels of x_2 are +3, -2, and -1. These deviations are the same at all levels of x_2 . Formally, the left panel of Figure 14.5 is expressed mathematically as an additive combination:

$$y = 6 + \beta_1 x_{11} + \beta_2 x_{21}$$

14.1.6.3 Nonadditive interaction of nominal predictors

When the predictor variables are non-metric, it does not make sense to talk about a multiplicative interaction, because there are no numbers to multiply. For example, consider predicting annual income from political party affiliation and gender. Both predictors are nominal, so it makes no sense to “multiply” them. But it does make sense to consider non-additive combination of their influences.

For example, the overall influence of gender is that men, on average, have a higher income than women. The overall influence of political party affiliation is that Republicans, on average, have higher income than Democrats. But it may be that influences combine non-additively: Perhaps people who are both Republican and male have a higher average income than would be predicted by merely adding the average boosts for being Republican and for being male. (This interaction is noted to be true; it is being used only as a hypothetical example.)

We need new notation to formalize the non-additive influence combination of nominal values. Just as x_1 refers to the value of predictor 1, and x_2 refers to the value of predictor 2, the notation $x_{1,2}$ will refer to a particular combination of values of predictors 1 and 2. If there are J_1 levels of predictor 1 and J_2 levels of predictor 2, then there are $J_1 \times J_2$ combinations of the two predictors.

A non-additive interaction of predictors is formally represented by including a term for the influence of combinations of predictors, beyond the additive influences, as follows:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_{1,2} x_{1,2}$$

Whenever the interaction coefficient $\beta_{1,2}$ is non-zero, the predicted value of y is not a mere addition of the separate influences of the predictors.

The right panel of Figure 14.5 shows a graphical example of nominal predictors that have interactive (i.e., non-additive) effects on the predicted variable. Notice, in the left pair of bars ($x_2 = h1; 0; 0$), that a change from $x_1 = h1; 0$ to $x_1 = h0; 1$ produces an increase of 2 in y , from $y = 8$ to $y = 10$. But for the middle pair of bars ($x_2 = h0; 1; 0$), a change from $x_1 = h1; 0$ to $x_1 = h0; 1$ produces an increase of 2 in y , from $y = 5$ to $y = 3$. Thus, the influence of x_1 is not the same at all levels of x_2 .

An interesting aspect of the pattern in the right panel of Figure 14.5 is that the average influences of x_1 and x_2 are the same as in the left panel. Overall, on average, going from $x_1 = h1; 0$ to $x_1 = h0; 1$ produces a change of 2 in y , in both the left and right panels. And overall, on average, for both panels it is the case that $x_1 = h1; 0$ is +3 above baseline, $x_2 = h0; 1; 0$ is -2 below baseline, and $x_2 = h0; 0; 1$ is -1 below baseline. The only difference between the two panels is that the combined influence of the two predictors equals the sum of the individual influences in the left panel, the combined influence of the two predictors does not equal the sum of the individual influences in the right panel.

An interaction between nominal predictors consists of an additional effect, for each specific combination of categorical values, away from the baseline combination. The magnitude of the interactive effect is whatever is left over after the additive effects have been applied to the baseline. The model that includes an interaction term can be written as

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_{1,2} x_{1,2} \quad (14.9)$$

²Actually, there is a sense in which nominal predictors cannot be multiplied to generate a predictor for interaction. Instead of ordinary multiplication of two scalars, use the outer product of two vectors. The outer product was defined in Section 8.8.1, p. 144.

$$= \beta_0 + \sum_{j=1}^{J_1} \beta_{1;j} x_{1;j} + \sum_{k=1}^{K_2} \beta_{2;k} x_{2;k} + \sum_{j=1}^{J_1} \sum_{k=1}^{K_2} \beta_{1;2;j;k} x_{1;j} x_{2;k}$$

with the constraints

$$\sum_{j=1}^{J_1} \beta_{1;j} = 0 \quad \text{and} \quad \sum_{k=1}^{K_2} \beta_{2;k} = 0 \quad \text{and} \quad \sum_{j=1}^{J_1} \sum_{k=1}^{K_2} \beta_{1;2;j;k} = 0 \quad \text{and} \quad \sum_{k=1}^{K_2} \sum_{j=1}^{J_1} \beta_{1;2;j;k} = 0 \quad (14.10)$$

In these equations, the term $\beta_{1;2;j;k}$ has J_1 times K_2 components, all of which are zero except for a 1 at the particular combination of levels $x_{1;j}$ and $x_{2;k}$. This mysterious and arcane notation will be revealed in all its majestic grandeur in [Chapter 19](#). For now, the main point is to understand that the term “interaction” refers to non-additive influence of the predictors on the predicted, regardless of whether the predictors are measured on a nominal scale or a metric scale.

14.1.7 Linking combined predictors to the predicted

Once the predictor variables are combined, they need to be mapped to the predicted variable. This mathematical mapping is called (inverse) link function, denoted by $f()$ in the following equation:

$$y = f(\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_{1;2} x_1 x_2) \quad (14.11)$$

Until now, we have been assuming that the link function is the identity function, $f(x) = x$. For example, in Equation 14.9, y equals the linear combination of the predictors; there is no transformation of the linear combination before mapping the result to y .

Before describing different link functions, it is important to make some clarifications of terminology and corresponding concepts. First, the f and f^{-1} in Equation 14.11 is usually called the inverse link function, because the link function itself is thought as transforming the values into a form that can be linked to the linear model. I will abuse convention and simply refer to either f or f^{-1} as “the” link function, and rely on context to disambiguate which direction of linkage is intended. The reason for this terminological sadness is that the arrows in hierarchical diagrams of Bayesian models will flow from the linear model toward the data, and therefore it is natural for functions to map toward the data, as in Equation 14.11. But repeatedly referring to f^{-1} as the “inverse” link would strain my patience and violate my aesthetic sensibilities. Second, the value that results from the link function $f(x)$ is not a data value per se. Instead, $f(x)$ is the value of a parameter that expresses some characteristic of the data, usually their mean. Therefore the function $f()$ in Equation 14.11 is sometimes called the mean function, and is written

$= f()$ instead of $y = f()$. I will not use this terminology because most students already think that “mean” means something else, namely the sum divided by N . The fact that in Equation 14.11 is a parameter value and not a data value will become clear in subsequent sections as we encounter specific cases and examples.

There are situations in which a non-identity link function is appropriate. Consider, for example, predicting response time as a function of amount of caffeine consumed. Response time declines as caffeine dosage increases, and therefore a linear prediction of response time from dosage would have a negative slope. This negative slope implies that for a very large dosage of caffeine, response time would become negative, which is impossible (unless caffeine causes precognition (i.e., foreseeing events before they occur)). Therefore a direct

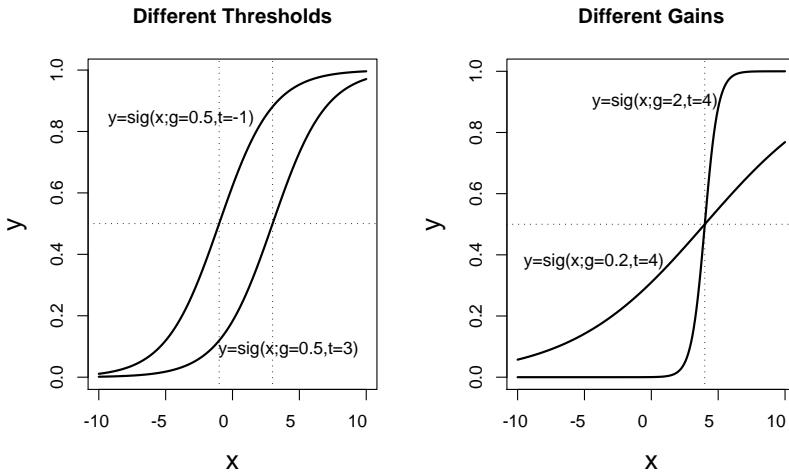


Figure 14.6: Examples of sigmoid functions of a single variable. The left panel shows sigmoids with the same gain but different thresholds. The right panel shows sigmoids with the same threshold but different gains.

linear function cannot be used for extrapolation to large values and we might instead want to use an exponential link function such as $y = \exp(\beta_0 + \beta_1 x)$.

14.1.7.1 The sigmoid (a.k.a. logistic) function

A frequently used link function is the sigmoid, also known as the logistic:

$$y = \text{sig}(x) = \frac{1}{1 + \exp(-x)} \quad (14.12)$$

Notice the negative sign in front of the parameter β_1 . The sigmoid function ranges between 0 and 1. The sigmoid is nearly 0 when x is large negative, and is nearly 1 when x is large positive.

For linear combinations of predictors, the sigmoid link function is most conveniently parameterized in threshold form. For a single predictor variable, the sigmoid function applied to the linear function of the predictor yields

$$y = \text{sig}(x; \beta_0, \beta_1) = \frac{1}{1 + \exp(-(\beta_0 + \beta_1 x))} \quad (14.13)$$

where β_1 , called the gain, corresponds to β_1 in Equation 14.2, and where β_0 , called the threshold, corresponds to $-\beta_0/\beta_1$ in Equation 14.2.

Examples of Equation 14.13, i.e., the sigmoid of a single predictor, are shown in Figure 14.6. Notice that the threshold is the point on the x -axis for which $y = 0.5$. The gain indicates how steeply the sigmoid rises through that point.

Figure 14.7 shows examples of a sigmoid of two predictors. Above each panel is the equation for the corresponding graph. The equations are parameterized in threshold form, as in Equation 14.4. In other words, $y = \text{sig}(\beta_0 + \beta_1 w_1 x_1 + \beta_2 w_2 x_2)$, with $\sqrt{\beta_1^2 + \beta_2^2} = 1$. Notice, in particular, that the coefficients of w_1 and w_2 in the plotted equations do indeed have Euclidean length of 1.0. For example, in the upper-right panel, $\sqrt{0.71^2 + 0.71^2} = 1.0$, except for rounding error.

The coefficients of the variables determine the orientation of the sigmoidal "cliff". For example, compare the two top panels in Figure 14.7, which differ only in the coefficients,

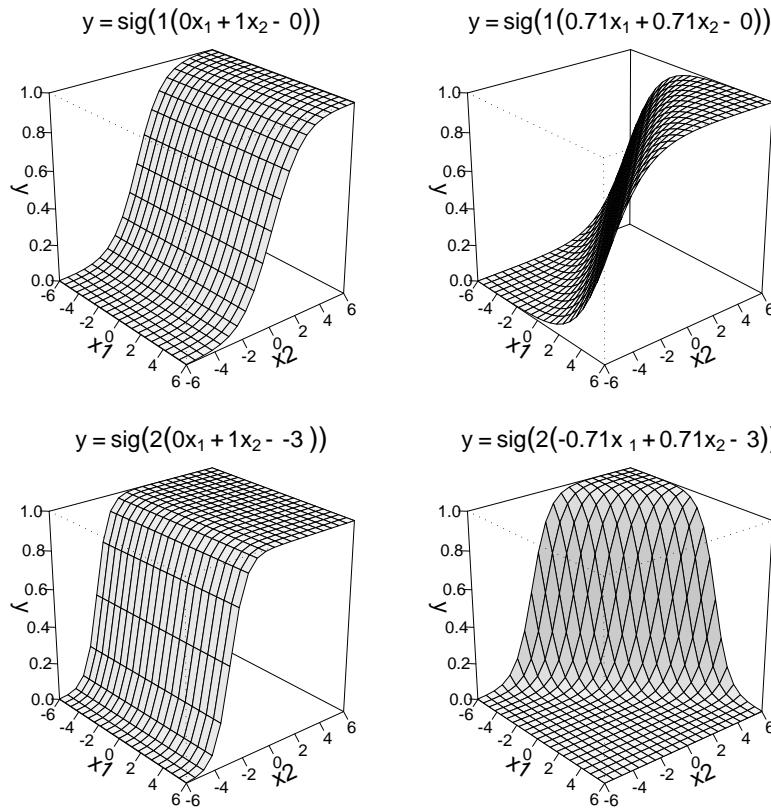


Figure 14.7: Examples of sigmoid functions of two variables. Top two panels show sigmoids with the same gain and threshold, but different coefficients on the predictors. The left two panels show sigmoids with the same coefficients on the predictors, but different gains and thresholds. The lower right panel shows a case with a negative coefficient on the first predictor.

not in gain or threshold. In the top left panel, the coefficients are $w_1 = 0$ and $w_2 = 1$, and the curve rises in the x_2 direction. In the top right panel, the coefficients are $w_1 = 0.71$ and $w_2 = 0.71$, and the curve rises in the positive diagonal direction.

The threshold determines the position of the sigmoidal curve. In other words, the threshold determines the values for which $y = 0.5$. For example, compare the two left panels of Figure 14.7. The coefficients are the same, but the thresholds (and gains) are different. In the upper left panel, the threshold is zero, and therefore the mid-level of the curve is over $x_2 = 0$. In the lower left panel, the threshold is -3, and therefore the mid-level of the curve is over $x_2 = -3$.

The gain determines the steepness of the sigmoidal curve. Again compare the two left panels of Figure 14.7. The gain of the upper left is 1, whereas the gain of the lower left is 2.

Terminology: The logit function. The inverse of the logistic function is called the logit function. For $0 < p < 1$, $\text{logit}(p) = \log(p/(1-p))$. It is easy to show (try it!) that $\text{logit}(\text{sig}(x)) = x$, which is to say that the logit is indeed the inverse of the sigmoid. Some authors, and programmers, prefer to express the connection between predictors and predicted in the opposite direction, by first transforming the predicted variable to match the

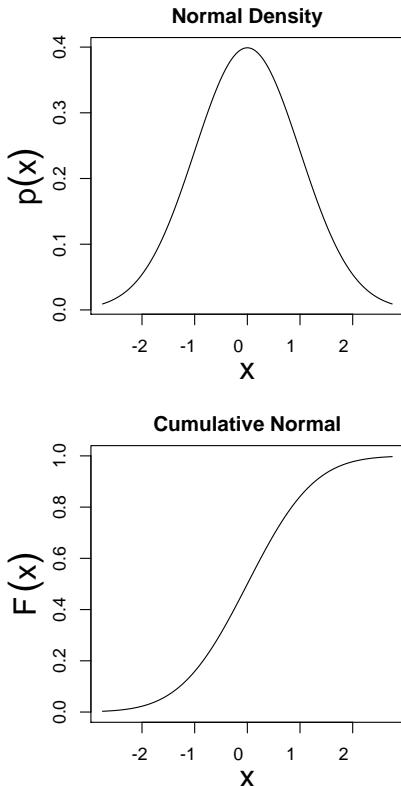


Figure 14.8: Top: A normal density with mean zero and precision one. Bottom: The corresponding cumulative normal function.

linear model. In other words, you may see the link expression of these ways:

$$y = \text{logistic } \beta_0 + \beta_1 x_1 + \dots$$

$$\text{logit}(y) = \beta_0 + \beta_1 x_1 + \dots$$

The two expressions achieve the same result, mathematically. The difference between them is merely a matter of emphasis. In the first expression, the combination of predictors is transformed so it maps onto ~~the~~ expressed in its original scale. In the second expression, transformed onto a new scale, and that transformed value ~~is~~ used as a combination of predictors.

14.1.7.2 The cumulative normal (a.k.a. Phi) function

Another frequently used link function is the cumulative ~~mat~~ distribution. It is qualitatively very similar to the sigmoid or logistic function. ~~Meters~~ will use the logistic or the cumulative normal depending on mathematical convenience or ease of interpretation. For example, when we consider ordinal predicted variables (Chapter 21), it will be natural to model the responses in terms of a continuous underlying variable that has normally distributed variability, which leads to using the cumulative normal as a model of response probabilities.

The cumulative normal is denoted $\Phi(x; \mu, \sigma)$, where x is a real number and where μ and σ are parameter values, called the mean and precision of ~~the~~ distribution. The parameter μ governs the point at which the cumulative normal $\Phi(x)$ equals 0.5. In other words, μ plays the same role as the threshold in the logistic sigmoid parameter σ governs the steepness of the cumulative normal function $\Phi(x)$. The σ parameter plays

the same role as the gain parameter in the logistic sigmoid. A graph of a cumulative normal appears in Figure 14.8. For this example, $\mu = 0$, and notice that $(0) = 0.5$. This means that the area under the normal density to the left of 0 is 0.5.

Terminology: The probit function. The inverse of the cumulative normal is called the probit function. (“Probit” stands for “probability unit” Bliss, 1934). The probit function maps a value p , for $0 < p < 1$, onto the infinite real line, and a graph of the probit function looks very much like the logit function. You may see the link expressed either of these ways:

$$y = \mu_0 + \mu_1 x_1 + \dots$$

$$\text{probit}(y) = \mu_0 + \mu_1 x_1 + \dots$$

Traditionally, the transformation of (in this case, the probit function) is called the link function, and the transformation of the linear combination (in this case, the function) is called the inverse link function. As mentioned before, use the traditional terminology and call either one a link function, relying on context to disambiguate.

14.1.8 Probabilistic prediction

In the real world, there is always variation that we cannot predict from. This unpredictable “noise” may be deterministically caused by sundry factors we have neither measured nor controlled, or the noise might be caused by non-determinism. It does not matter either way because in practice the best we can predict the probability that y will have any particular value, dependent upon x . Therefore we use the deterministic value predicted by Equation 14.11 as the predicted tendency of y as a function of the predictors. We do not predict that y is exactly $\mu_0 + \mu_1 x_1 + \mu_2 x_2 + \dots$ because we would surely be wrong. Instead, we predict that y tends to be near $\mu_0 + \mu_1 x_1 + \mu_2 x_2 + \dots$.

To make this notion of probabilistic tendency precise, we specify a probability distribution for y that depends on $\mu_0 + \mu_1 x_1 + \mu_2 x_2 + \dots$. To keep the notation tractable, first define $f = f(\mu_0 + \mu_1 x_1 + \mu_2 x_2 + \dots)$. Do not confuse this use of f with the unrelated f mentioned in the cumulative normal function. With this f , we then denote the probability distribution of y as some to-be-specified probability density function, abbreviated as “pdf”:

$$y \sim \text{pdf}(\cdot; \mu_0 + \mu_1 x_1 + \mu_2 x_2 + \dots)$$

The pdf might have various additional parameters, denoted b , to specify its shape. Examples are provided in the next section, where all these ideas are brought together.

14.1.9 Formal expression of the GLM

In general, the likelihood function specifies the probability of each possible predicted value y as a function of the predictor values and various parameter values etc. The generalized linear model can be written:

$$= f(\mu_0 + \mu_1 x_1 + \mu_2 x_2 + \dots) \quad (14.14)$$

$$y \sim \text{pdf}(\cdot; \mu_0 + \mu_1 x_1 + \mu_2 x_2 + \dots) \quad (14.15)$$

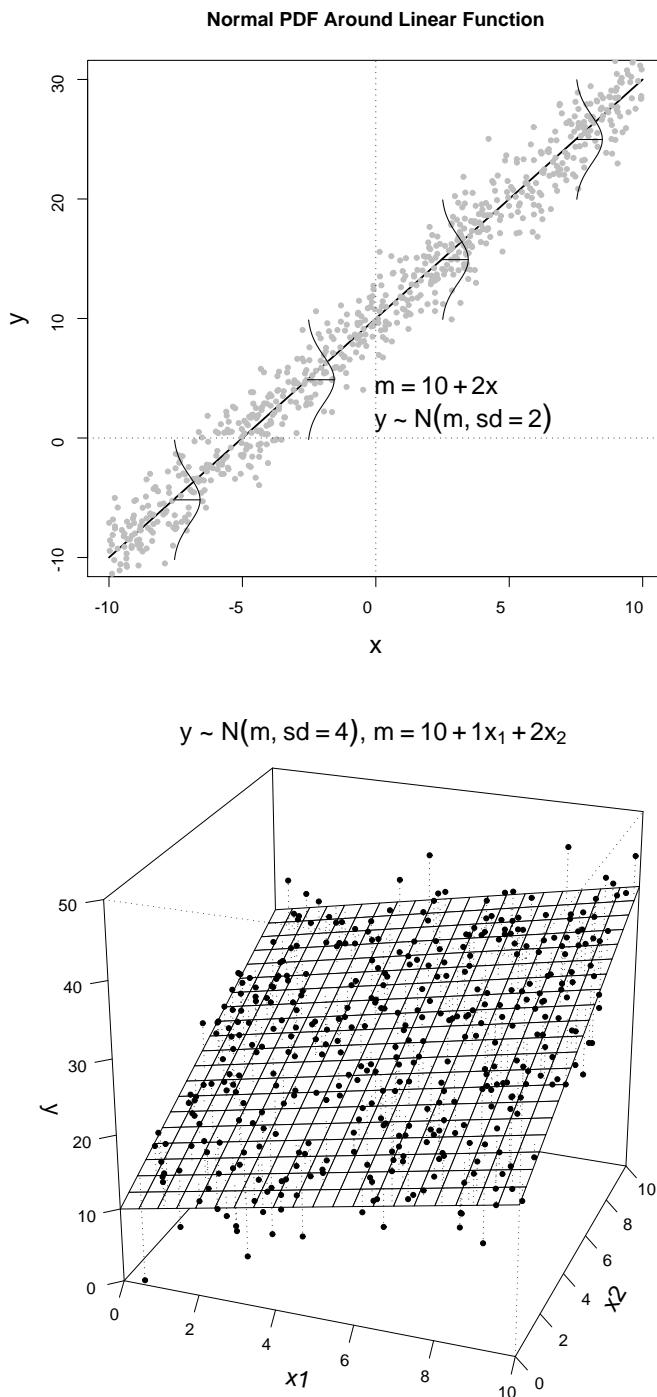


Figure 14.9: Examples of points normally distributed around a linear function.

$$y \sim \text{dbern}(m), m = \text{sig}(1(0.71x_1 + 0.71x_2 - 0))$$

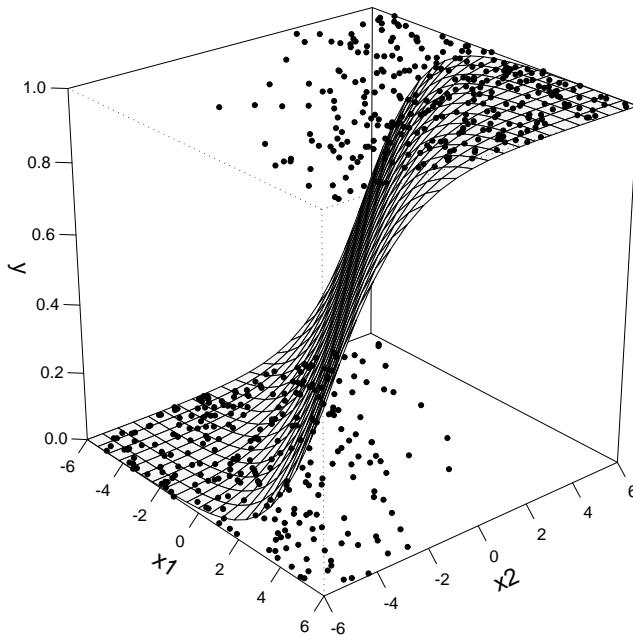


Figure 14.10: Examples of points Bernoulli distributed and a sigmoid function of two predictors. All the points are either $y=1$ or $y=0$ (intermediate values such as $y=.6$ cannot occur).

The function f in Equation 14.14 is called the “link” function, because it links the combination of predictors to the predicted tendency. The ~~optimal~~ parameters; $[; ::]$ in Equation 14.15 may be needed for various types of the probability density function (pdf) that describe the probability distribution \mathbf{y} around the tendency.

Figure 14.9 shows a random sample of points normally distributed around a line or plane. The upper panel illustrates a case of the generalized model of Equations 14.14 and 14.15 in which there is a single predictor with $\beta_0 = 10$ and $\beta_1 = 2$. The link function is simply the identity function $f(\beta_0 + \beta_1 x) = \beta_0 + \beta_1 x$. The probability density function is normal with a standard deviation of 2.0. Profiles of this density are superimposed on the graph to make it explicit. Notice that the normal density is always centered on the line that marks the predicted tendency as a function of the predictor.

The lower panel of Figure 14.9 shows a case with two predictors. The predictors are combined linearly, with no interaction. The link function is the identity. The probability function is normal with a standard deviation of 4.0. Each randomly generated point is connected to the underlying linear core by a vertical line, to explicitly indicate the random variation of the point from the plane. The line marks the predicted tendency as a function of the predictors, and the data values are normally distributed above and below that tendency.

Figure 14.10 shows another case of the GLM. In this case, the points are Bernoulli distributed around a sigmoid function of two predictors annotated at the top of the graph. There is a linear combination of predictors, with a sigmoid function, and a Bernoulli probability function that defines the probability that $y=1$. The graph shows that values of y can only be 0 or 1, and the sigmoid function defines the probability that $y=1$ for particular

predictor values. The sigmoidal surface plots the tendency = 1 as a function of the predictors.

14.2 Cases of the GLM

Table 14.1, p. 312, displays the various cases of the generalized linear model that are considered in this book. Subsequent chapters of the book go through the table in reading order: left to right within rows, then top to bottom across rows.

The first row of Table 14.1 lists cases for which the predicted variable is metric. Moving from left to right within this row, the first column indicates a situation in which there is only a single group, and the predicted value for the group is simply the mean of the group. In this situation, there is no need to explicitly denote a predictor variable, and instead the mean of the group can be denoted by a single parameter. This situation corresponds to what classical null hypothesis significance testing (NHST) calls a single-group-test. This case is described in its Bayesian setting in Chapter 15.

Moving to the next column, there is a single metric prediction. This corresponds to so-called “simple linear regression”, and is explored in Chapter 16. By inspecting the equation for the GLM in the cell, you can see that the only difference from the previous cell is the inclusion of the predictor x_1 and its coefficient β_1 .

Moving rightward to the next column, we come to the scenario involving two or more metric predictors, which corresponds to “multiple regression” and is explored in Chapter 17. By examining the equations for the GLM in the cell, you see that the basic form is the same, but merely with extra terms added for the additional predictors.

The next two columns involve nominal predictors, instead of metric predictors, with the penultimate column devoted to a single predictor and the last column devoted to two or more predictors. The last two columns correspond to what *ANOVA* “oneway ANOVA” and “multifactor ANOVA”. If that terminology is unfamiliar to you, don’t worry, it will be explained in Chapters 18 and 19.

In all the cases in the first row, the link function is the identity and the probability distribution for the metric predicted values is assumed to be normal. When we move to the second row, however, the predicted variable is dichotomous and therefore the probability distribution for y is a Bernoulli distribution. The link function, which connects the predictors to the probability that $y = 1$, is assumed to be the sigmoid, i.e., logistic function. When the predictors are metric, this situation is generically referred to as “logistic regression” and is discussed in Chapter 20. The case of nominal predictors is also discussed.

Finally, the bottom row of Table 14.1 lists cases for which the predicted variable is ordinal. These cases are considered in Chapter 21. Note that the link function is the cumulative normal instead of the sigmoid, and the ordinal values are generated by multiple-category generalization of the Bernoulli function, denoted dcat. Again, this will be explained at length in the forthcoming chapters. The point is for you to see the overall organization of topics, and to see how all these cases arise as of the same underlying structure.

The table can be expanded with additional rows and columns when it gets too big to display easily. Additional columns would include combinations of metric and nominal predictors. But it turns out that it is easy in Bayesian models to combine metric and nominal predictors, once you know how to handle metric and nominal predictors individually. Additional rows would involve different types of predicted variables. In particular, a fourth

Table 14.1: Cases of the generalized linear model (single predictor variable).

Predictor Variable(s) x					
	Single Group	Metric		Nominal	
		Single x Factor	Two or More x Factors	Single x Factor	Two or More x Factors
continuous	$y = \mu + \sigma \mathcal{N}(\cdot; \cdot)$	$y = \mu + \beta_1 x_1 + \sigma \mathcal{N}(\cdot; \cdot)$	$y = \mu + \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_{12} x_1 x_2 + \sigma \mathcal{N}(\cdot; \cdot)$	$y = \mu + \beta_0 + \beta_1 x_1 + \sigma \mathcal{N}(\cdot; \cdot)$	$y = \mu + \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_{12} x_1 x_2 + \beta_{123} x_1^2 + \sigma \mathcal{N}(\cdot; \cdot)$
categorical	$y = \text{sig}_0 + \text{dbern}()$	$y = \text{sig}_0 + \beta_1 x_1 + \text{dbern}()$	$y = \text{sig}_0 + \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_{12} x_1 x_2 + \text{dbern}()$	$y = \text{sig}_0 + \beta_0 + \beta_1 x_1 + \text{dbern}()$	$y = \text{sig}_0 + \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_{12} x_1 x_2 + \text{dbern}()$
nominal		$y = \mu + \beta_1 x_1 + \text{thresh}_k(\cdot) + \text{dcat}(\cdot; \cdot_k; \cdot)$	$y = \mu + \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_{12} x_1 x_2 + \text{thresh}_k(\cdot) + \text{dcat}(\cdot; \cdot_k; \cdot)$		

row would include count data for the predicted values. We will, in fact, cover one such case, as described in the next section. When the predicted data are count values, a natural link function is the exponential, and a natural pdf is the Poisson distribution, which will be defined later in the book (Section 22.1.3). In summary, the rows of the table refer to differently scaled predicted values, with their corresponding link functions and pdf's:

y Scale Type	(Inverse) Link Function	pdf
metric	identity	normal
dichotomous	sigmoid (a.k.a. logistic)	Bernoulli
ordinal	thresholded cumulative normal	categorical
count	exponential	Poisson

14.2.1 Two or more nominal variables predicting frequency

Finally, we will also consider the situation in which there are two or more nominal variables used as predictors of frequency count. A frequency count, i.e., how many times something happened, is a special case of a metric scale, but because counts fall at discrete levels, namely non-negative integers, this situation will have a different sort of likelihood distribution. This type of situation, with nominal predictors and frequency-count predicted values, is often called “contingency table analysis” and a typical R analysis conducts a “chi-square test of independence of attributes”. We explore Bayesian analysis of this situation in Chapter 22.

Here is a brief summary of how contingency tables are analyzed using a model much like those in Table 14.1. In fact, a fourth row could be added to Table 14.1, with the predicted type labeled frequency count, and the model falling in the final column, under two nominal predictors. As a concrete example, suppose we are analyzing a political affiliation and religious affiliation of a set of people, and for a sample of people we count many occurrences there are of each combination. We are interested in analyzing possible relationships between political and religious affiliations. Suppose we conduct a poll for one week. We happen to record 27 people who are Democrats and Unitarians. This observed frequency reflects an underlying rate at which that combination is generated by this sort of poll, i.e., the underlying rate for Unitarian Democrats is roughly 27 people per week. The observed rate (i.e., frequency per unit time) for each combination of nominal values is thought to reflect the true underlying rate at which that combination is generated by the world. We conceive of the observed rate as being a random draw from a true underlying rate denoted by λ . The probability of any particular observed rate, given an underlying rate of λ , is modeled by a Poisson distribution, which is denoted as `ppois()`. The Poisson distribution was smuggled into the text back in Exercise 11.3, p. 235, which I'm sure is still as fresh in your memory as a beached fish. Don't try to understand the Poisson; it will be explained again later (Section 22.1.3). The Poisson distribution specifies a probability for each possible observed rate. The Poisson puts highest ~~probabilities~~ on rates near λ .

Our goal is to estimate the underlying rates at which each combination is produced. But more than that, we would like to know if the attributes occur independently of each other, or instead covary in some way. For example, if political and religious affiliation are independent, then there should be the same proportion Unitarians among Democrats as among Republicans. Mathematically, independence means $p(\text{Unitarian} \& \text{Democrat}) = p(\text{Unitarian}) \cdot p(\text{Democrat})$ and $p(\text{Unitarian} \& \text{Republican}) = p(\text{Unitarian}) \cdot p(\text{Republican})$ and so on for every combination of attribute values. To shorten the expressions, I'll substitute

tute U for Unitarian and D for Democrat, whereby independence means

$$p(U \& D) = p(U) \cdot p(D)$$

and so on for every combination of attributes. That expression for probabilities corresponds to the following expression in terms of frequencies:

$$\text{freq}(U \& D) = N = \text{freq}(U) = N \quad \text{freq}(D) = N$$

which can be re-arranged as

$$\text{freq}(U \& D) = \text{freq}(U) \cdot \text{freq}(D) \cdot 1 = N$$

Notice that independence is expressed as a multiplicative product of attribute frequencies. But all the models we've considered in this chapter used additive sum of predictor in uences. To be able to use our familiar additive models, we'll transform the frequencies by a logarithm, because the logarithm of the product of values is the sum of the logarithms of the values. In other words,

$$\log(\text{freq}(U \& D)) = \log(\text{freq}(U)) + \log(\text{freq}(D)) + \log(1) = N$$

The notation “ $\log(\text{freq}(\text{value}))$ ” gets cumbersome, so we'll substitute the notation. Thus,

$$\log(\text{freq}(U \& D)) = u + d + o$$

where o stands in for the constant $\log(N)$. Finally, it's unintuitive to talk about the logarithms of frequencies, so we'll exponentiate to get rid of the leading logarithm, yielding:

$$\text{freq}(U \& D) = \exp(u + d + o)$$

To summarize, if independence is true, then the expression above should be true, for every combination of attribute values.

But of course the attributes are usually not independent, we would like some measure of lack of independence. We already have such a measure in the context of linear models, namely, the interaction term. Thus, we will include an interaction term that estimates deviation from independence. Thus, our model ends up as follows. For two nominal attributes, we put the observed frequencies in a table with one attribute's values listed down the rows, and the other attribute's values across the columns. The frequency in the r^{th} row and c^{th} column is denoted freq_{rc} , and the underlying rate for that cell is denoted r_c . The model looks like this:

$$\begin{aligned} r_c &= \exp(o + r + c + r_c) \\ \text{freq}_{rc} &\sim \text{dpois}(r_c) \end{aligned} \tag{14.16}$$

with the usual constraints (from Equation 14.10)

$$\sum_r r_c = 0 \quad \sum_c r_c = 0 \quad \sum_r r_{c;r;c} = 0 \quad \sum_c r_{c;r;c} = 0$$

The point of this over-fast prelude to contingency tables is merely to demonstrate that the core of the model we'll be using is the same as the linear model that was mentioned for multifactor ANOVA in Table 14.1. Thus, all the applied analyses we'll see in the remainder of the book are based on the GLM.

14.3 Exercises

Exercise 14.1. [Purpose: For real-world examples of research, identify which statistical model is relevant.] For each of the examples below, identify the predicted ~~variable~~ and its scale type, identify the predictor variable(s) and its scale type, and identify the cell of Table 14.1 to which the example belongs.

(A) Guber (1999) examined average performance by public high school students on the Scholastic Aptitude Test (SAT) as a function of how much ~~money~~ was spent per pupil by the state, and what percentage of eligible students ~~actually~~ took the exam.

(B) Hahn, Chater, and Richardson (2003) were interested in ~~perceived~~ similarity of simple geometric patterns. Human observers rated pairs ~~steps~~ for how similar the patterns appeared, by circling one of the digits 1–7 printed ~~the~~ page, where 1 meant “very dissimilar” and 7 meant “very similar”. The authors ~~presented~~ a theory of perceived similarity, in which patterns are perceived to be dissimilar to the extent that it takes more geometric transformations to produce one pattern from ~~the other~~. The theory specified the exact number of transformations needed to get from one pattern to the other.

(C) R. L. Berger, Boos, and Guess (1988) were interested in ~~the velocity~~ of rats, measured in days, as a function of the rat's diet. One group ~~of fed~~ freely, another group of rats had a very low calorie diet.

(D) McIntyre (1994) was interested in predicting the tar content of a cigarette (measured in milligrams) from the weight of the cigarette.

(E) You are interested in predicting the gender of a person, ~~based~~ on the person's height and weight.

(F) You are interested in predicting whether a respondent ~~agrees~~ or disagrees with the statement, “The United States needs a federal health plan with a public option”, on the basis of the respondent's political party ~~list~~.

Exercise 14.2. [Purpose: Find student-relevant real-world examples of each type of situation in table 14.1.] For each of the twelve cells of Table 14.1 that is filled with ~~questions~~, provide an example of research involving that cell's model structure. Do this by finding published articles that describe research with the corresponding ~~structure~~. The articles ~~do~~ not need to have Bayesian data analysis; the articles ~~do~~ need to report research that involves the corresponding types of predictor and predicted variables. Because it might be overly time consuming to find published examples of all twelve types, find published articles of at least six types. For each example, specify the following:

The full citation to the published article (see the references of this book for examples of how to cite articles),

The predictor and predicted variables. Describe their ~~mean~~ and the type of scale they are. Briefly describe the meaningful context for the ~~variables~~, i.e., what is the goal of the research.

Chapter 15

Metric Predicted Variable on a Single Group

Contents

15.1	Estimating the mean and precision of a normal likelihood	318
15.1.1	Solution by mathematical analysis	318
15.1.2	Approximation by MCMC in BUGS	322
15.1.3	Outliers and robust estimation: The <i>t</i> -distribution	323
15.1.4	When the data are non-normal: Transformations	326
15.2	Repeated measures and individual differences	328
15.2.1	Hierarchical model	330
15.2.2	Implementation in BUGS	331
15.3	Summary	333
15.4	R code	333
15.4.1	Estimating the mean and precision of a normal likelihood	333
15.4.2	Repeated measures: Normal across and normal within	335
15.5	Exercises	338

It's normal to want to fit in with your friends,
Behave by their means and believe all their ends.
But I'll be high tailing it, fast and askew,
Precisely 'cause I can't abide what you do.

In this chapter, we consider a situation in which we have a metric predicted variable that is observed for items from a single group. For example, we could measure the blood pressure (i.e., a metric variable) for people randomly sampled from first-year university students (i.e., a single group). In this case we might be interested in whether the group's typical blood pressure differs from the recommended value for people of that age as published by a federal agency. As another example, we could measure IQ (i.e., a metric variable) of people randomly sampled from everyone self-identified as vegetarian (i.e., a single group). In this case we could be interested in whether the group's IQ differs from the general population's average IQ of 100.

In the context of the generalized linear model (GLM) introduced in the previous chapter, this chapter's situation involves the most trivial case of the linear core of the GLM, namely

Kruschke, J. K. (2010) Doing Bayesian Data Analysis: A Tutorial with R and BUGS. Academic Press Elsevier. Copyright 2010 by John K. Kruschke. Draft of May 11, 2010. Please do not distribute this preliminary draft. If you report Bayesian analyses based on this book, please do cite it! “

= 0, as indicated in the top-left cell of Table 14.1 (p. 312). The “news” of the present chapter is a detailed look at a particular probability density function, namely the normal distribution, denoted by $N(\mu, \sigma^2)$ in the top-left cell of Table 14.1. We will explore particular prior distributions for the parameters of the normal distribution, and see how those same density functions can be used in extended hierarchical models for research designs involving repeated measures. We will also provide examples of transforming data so they are approximately normal, so that “off the shelf” normal-likelihood models can be used to describe the data.

15.1 Estimating the mean and precision of a normal likelihood

The workhorse for the remainder of the book is the normal likelihood function, just as the Bernoulli likelihood was the focus for the previous part of the book. The normal probability density function was introduced in Section 3.3.2.2, p. 30. Please review that section now. The normal distribution specifies the probability density of a value y , given the values of two parameters, the mean μ and standard deviation:

$$p(y; \mu, \sigma^2) = \frac{1}{Z} \exp \left(-\frac{1}{2} \frac{(y - \mu)^2}{\sigma^2} \right) \quad (15.1)$$

where Z is the normalizer, i.e., a constant that makes the probability density integrate to 1. It turns out that $Z = \sqrt{2\pi\sigma^2}$, but we won't need to use this fact.

To get an intuition for the normal likelihood function, consider three data values $y_1 = 85$, $y_2 = 100$, and $y_3 = 115$, as plotted by large dots in Figure 15.1. The probability density of any single datum, given parameter values (μ, σ^2) as specified in Equation 15.1. The probability of the whole set of independent data values is $p(D_j; \mu, \sigma^2) = p(D_j; \mu, \sigma^2)$, where $D = \{y_1, y_2, y_3\}$. Figure 15.1 shows $p(D_j; \mu, \sigma^2)$ for different values of μ and σ^2 . As you can see, there are values of μ and σ^2 that make the data most probable, but other nearby values also accommodate the data reasonably well.

Given a set of data D , we estimate the parameters with Bayes' rule:

$$p(\mu, \sigma^2 | D) = \frac{p(D_j; \mu, \sigma^2) p(\mu, \sigma^2)}{\sum_{\mu, \sigma^2} p(D_j; \mu, \sigma^2) p(\mu, \sigma^2)} \quad (15.2)$$

Figure 15.1 showed examples $p(D_j; \mu, \sigma^2)$ for a particular data set at different values of μ and σ^2 . The prior, $p(\mu, \sigma^2)$, specifies the believability of each combination of values in the two-dimensional conjoint parameter space, without data. Bayes' rule says that the posterior believability of each combination of values is the prior believability times the likelihood, normalized by the evidence. We saw our first examples of Bayes' rule on a two-dimensional parameter space back in Chapter 8, for plain Figure 8.2, p. 134. Our goal now is to evaluate Equation 15.2, using a normal-density likelihood, for reasonable choices of the prior distribution $p(\mu, \sigma^2)$.

15.1.1 Solution by mathematical analysis

Purely for the sake of mathematical elucidation, it is correct to consider the case in which the standard deviation of the likelihood function is fixed at a specific value. In other words, the prior distribution on σ^2 is a spike over that specific value. We'll denote the value

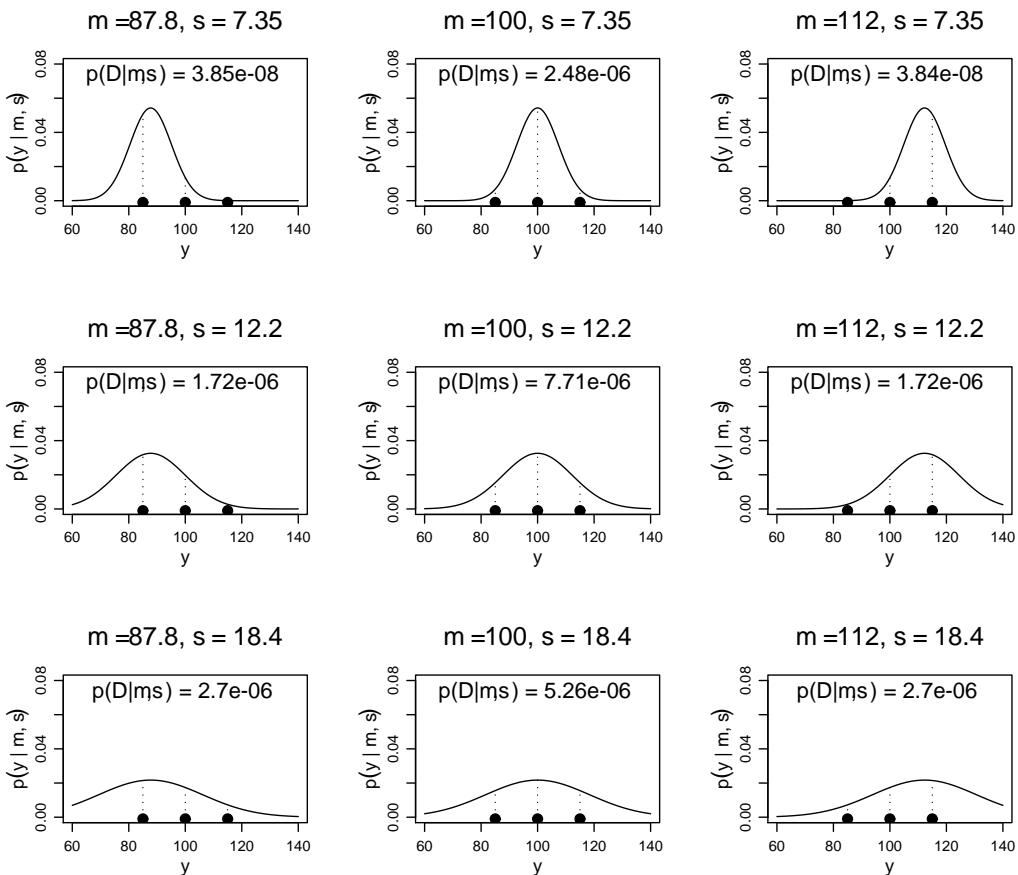


Figure 15.1: The likelihood $p(D_j ; \theta)$ for three data points $D = f85, 100, 115g$ according to a normal likelihood function with different values of m and s . The probability density of an individual datum is the height of a dotted line over the point. The probability of the set of data is the product of the individual probabilities. The middle panel shows m and s that maximize the probability of the data.

as $\theta = S_y$. With this simplifying assumption, we are only estimating θ because we are assuming perfectly certain prior knowledge about θ .

When θ is fixed, then the prior distribution on θ in Equation 15.2 can be easily chosen to be conjugate to the normal likelihood. (The term “conjugate prior” was defined in Section 5.2, p. 67). It turns out that the product of normal distributions is again a normal distribution; in other words, if the prior on θ is normal, then the posterior on θ is normal. It is easy to derive this fact, as we do next.

Let the prior distribution on θ be normal with mean M and standard deviation S . Then the likelihood times the prior is

$$p(y_j ; \theta)p(\theta) = p(y_j ; S_y)p(\theta)$$

$$\propto \exp\left(-\frac{1}{2} \frac{(y_j - M)^2}{S_y^2}\right) \exp\left(-\frac{1}{2} \frac{\theta^2}{S^2}\right)$$

$$\begin{aligned}
 &= \exp \left(-\frac{1}{2} \frac{(y - M)^2}{S_y^2 + S^2} \right) + \frac{M}{S^2} \frac{1}{S_y^2 + S^2} \\
 &= \exp \left(-\frac{1}{2} \frac{S^2(y - M)^2 + S_y^2}{S_y^2 S^2} \right) + \frac{M}{S^2} \frac{1}{S_y^2 + S^2} \\
 &= \exp \left(-\frac{1}{2} \frac{S_y^2 + S^2}{S_y^2 S^2} \right) + \frac{\frac{S_y^2 M_u + S^2 y}{S_y^2 + S^2} + \frac{S_y^2 M_u^2 + S^2 y^2}{S_y^2 + S^2}}{S^2} \\
 &= \exp \left(-\frac{1}{2} \frac{S_y^2 + S^2}{S_y^2 S^2} \right) + \frac{\frac{S_y^2 M_u + S^2 y}{S_y^2 + S^2}}{S^2} \\
 &\quad \exp \left(-\frac{1}{2} \frac{S_y^2 + S^2}{S_y^2 S^2} \right) \\
 &/ \exp \left(-\frac{1}{2} \frac{S_y^2 + S^2}{S_y^2 S^2} \right) + \frac{\frac{S_y^2 M_u + S^2 y}{S_y^2 + S^2}}{S^2} \quad (15.3)
 \end{aligned}$$

where the transition to the last line was valid because the term that was dropped was merely a constant. This result, believe it or not, is progress. Why? Because we ended up, in the innermost parentheses, with a quadratic expression. Notice that the normal prior is also a quadratic expression in . All we have to do is “complete the square” inside the parentheses, and do the same trick that got us to the last line in Equation 15.3:

$$\begin{aligned}
 p(y_j ; S_y) p(\mu) / \exp \left(-\frac{1}{2} \frac{S_y^2 + S^2}{S_y^2 S^2} \right) &+ \frac{\frac{S_y^2 M_u + S^2 y}{S_y^2 + S^2}}{S^2} \\
 &= \exp \left(-\frac{1}{2} \frac{S_y^2 + S^2}{S_y^2 S^2} \right) + \frac{\frac{S_y^2 M_u + S^2 y}{S_y^2 + S^2}}{S^2} \quad (15.4)
 \end{aligned}$$

Equation 15.4 is the numerator of Bayes' rule. When it is divided by the evidence in the denominator, it becomes a probability density function. What's the shape of the function? You can see that Equation 15.4 has the same form as a normal distribution, such that the mean is $\frac{S_y^2 M_u + S^2 y}{S_y^2 + S^2}$ and the standard deviation is $\frac{S_y^2 S^2}{S_y^2 + S^2}$.

That formula is rather unwieldy! It becomes more compact if the normal density is re-expressed in terms of σ^{-2} instead of S^2 . The reciprocal of the squared standard deviation is called the precision of the normal. When the standard deviation goes down, the precision goes up. A very narrow distribution is highly precise. A wide distribution has low precision. Because the posterior standard deviation is $\frac{S_y^2 S^2}{S_y^2 + S^2}$, the posterior precision is

$$\frac{S_y^2 + S^2}{S_y^2 S^2} = \frac{1}{S^2} + \frac{1}{S_y^2}$$

In other words —and this is the punch line— the posterior precision is the sum of the prior precision and the likelihood precision.

The posterior mean can also be re-expressed in terms of precision. The posterior mean

is $\frac{S_y^2 M + S^2 y}{S_y^2 + S^2}$, which becomes

$$\frac{1=S^2}{1=S_y^2 + 1=S^2} M + \frac{1=S_y^2}{1=S_y^2 + 1=S^2} y.$$

In other words, the posterior mean is a weighted average of the prior mean and the datum, with the weighting corresponding to the relative precision of the prior and the likelihood. When the prior is highly precise compared to the likelihood, when S^2 is large compared to S_y^2 , then the prior is weighted heavily and the posterior mean is near the prior mean. But when the prior is imprecise, i.e., very uncertain, the prior does not get much weight and the posterior mean is close to the datum. We have previously seen this sort of relative weighting of prior and data in the posterior. It was up in the case of updating a beta prior, back in Equation 5.8, p. 71.

The formulas for the mean and precision of the posterior can be naturally extended when there are values of y in a sample, instead of only a single value of y . The formulas can be derived from the defining formulas, as was done above, but a short cut can be taken. It is known from mathematical statistics that a set of values are generated from a normal likelihood function, the mean of those values, \bar{y} , is also distributed normally, with the same mean as the generating mean, and standard deviation of $= \sqrt{\frac{S^2}{N}}$. Thus, instead of conceiving of this situation as scores sampled from the likelihood $N(y_i | \mu, S^2)$, we conceive of this as a single score, sampled from the likelihood $N(\bar{y} ; \mu, S^2/N)$. Then we just apply the updating formulas we previously derived. Thus, for N scores y_i generated from a normal likelihood $N(y_i ; \mu, S^2)$ and prior distribution $N(\mu | M, S^2)$, the posterior distribution is also normal with mean

$$\frac{1=S^2}{N=S_y^2 + 1=S^2} M + \frac{N=S_y^2}{N=S_y^2 + 1=S^2} \bar{y}$$

and precision

$$\frac{1}{S^2} + \frac{N}{S_y^2}.$$

Notice that as the sample size increases, the posterior mean is dominated by the data mean.

Instead of estimating the parameter in the likelihood when it is fixed, we can estimate the parameter when it is unknown. The situation is also more conveniently expressed in terms of precision: We want to estimate the precision, τ^2 , when it is unknown. It turns out that when it is unknown, a conjugate prior for the precision is the gamma distribution (e.g., Gelman et al., 2004, p. 50). The gamma distribution was described in Figure 9.8, p. 170. For our purposes, it is not important to review the updating formula for the gamma distribution in this situation. But it is important to gain an intuition for what a gamma prior on precision means, in terms of the beliefs represented. Consider a gamma distribution that is loaded heavily over very small values, but has a long tail extending over large values. This sort of gamma distribution on precision indicates that we believe most strongly in small precisions, but we admit that large precisions are possible. This is a belief about the precision of a normal likelihood function, then this sort of gamma distribution expresses a belief that the data will be very spread out, because small deviations imply large standard deviations. If the gamma distribution is loaded over large values of precision, it expresses a belief that the data will be tightly clustered.

Summary. We have assumed that the data are generated by a normal likelihood function, parameterized by a mean and standard deviation, and denoted $N(y_i; \mu, \tau)$. For purposes of mathematical derivation, we made unrealistic assumptions that the prior distribution is either a spike or a spike on τ , in order to make three main points:

1. A natural way to express a prior on μ is with a normal distribution, because this is conjugate with the normal likelihood when its precision τ is fixed.
2. A natural way to express a prior on the precision τ is with a gamma distribution, because this is conjugate with the normal likelihood when the mean μ is fixed.
3. The formulas for Bayesian updating of the parameter distribution are more conveniently expressed in terms of precision than standard deviation. Normal distributions will be described sometimes in terms of standard deviation and sometimes in terms of precision, so it is important to glean from context which is being referred to.

A joint prior, on the combination of μ and τ parameter values, can also be specified, in such a way that the posterior has the same form as the prior. We will not pursue these mathematical analyses here, because our purpose is merely to identify and motivate typical expressions for the prior distributions on the parameters, so that they can then be used in MCMC sampling in BUGS. Various other sources describe joint priors for the joint parameter space (e.g., Gelman et al., 2004, pp. 78–83).

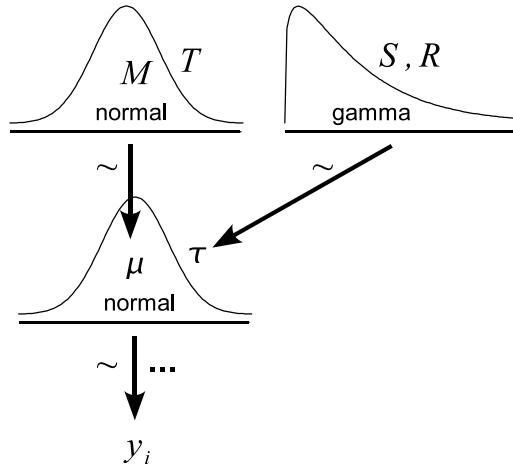


Figure 15.2: A model of dependencies for several metrics drawn from a single group. The normal distributions are parameterized by mean and precision (not standard deviation).

15.1.2 Approximation by MCMC in BUGS

It is easy to estimate the mean and precision of a set of data y_i . Figure 15.2 illustrates the model. The data y_i are assumed to be generated by a normal likelihood function with mean μ and precision τ . The prior on μ is assumed to be normal with mean M and precision T . The prior on τ is assumed to be gamma with shape and rate parameters S and R . We are estimating two parameters, and the prior is over the two-dimensional conjoint parameter space. The prior we have specified assumes that μ and τ are independent.

The model of Figure 15.2 is expressed in BUGS as follows: ([YmetricXsingleBrugs.R](#))

```

9  model {
10    # Likelihood:
11    for( i in 1 : N ) {
12      y[i] ~ dnorm( mu , tau ) # tau is precision, not SD
13    }
14    # Prior:
15    tau ~ dgamma( 0.01 , 0.01 )
16    mu ~ dnorm( 0 , 1.0E-10 )
17  }
```

Notice that each arrow in Figure 15.2 has a corresponding `sp@tion` in the BUGS code. Notice that BUGS parameterizes `norm` by mean and precision, not by mean and standard deviation. The hyperprior constants in this particular BUGS model are `gic` and uninformed; they should be knowledgably specified in real `ap@tions`.

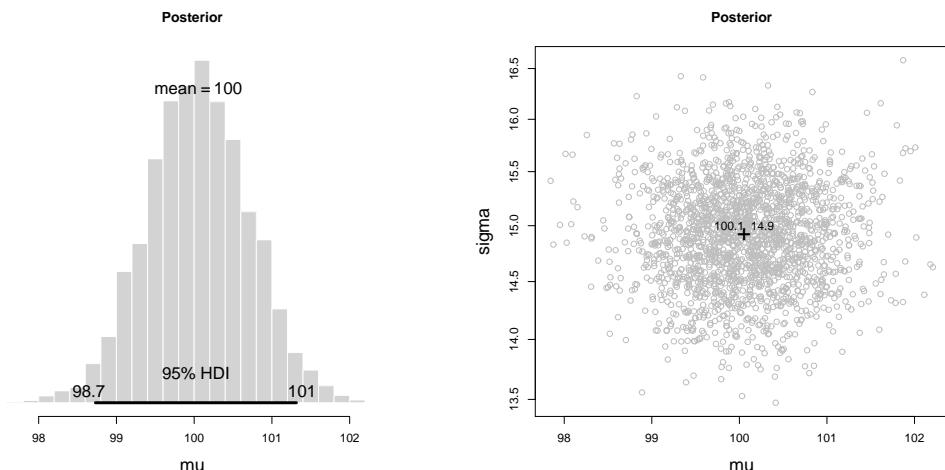


Figure 15.3: Posterior distribution from program in Section 15.4.1 ([YmetricXsingleBrugs.R](#)).

When learning about or debugging a model, it can be useful to generate fictitious data from known parameter values, and then see how well those generated values are estimated by the model's posterior distribution. For this purpose, the program in Section 15.4.1 ([YmetricXsingleBrugs.R](#)) generates some random data from a normal distribution whose true mean is 100 and standard deviation is 15 (just like IQs). The random data themselves have a sample mean and standard deviation somewhat different from those generating values. The result of running the BRugs program is shown in Figure 15.3. The posterior precision (i.e., `tau`) sampled by the BUGS program has been converted to standard deviation (denoted "sigma") via the identity $\sigma = 1/\sqrt{\tau}$. The estimated mean and standard deviation are not far from the generating mean. This reassures us that the program is operating properly.

15.1.3 Outliers and robust estimation: The t distribution

Figure 15.4 shows examples of the t distribution. The t distribution was originally invented by Gosset (1908), who used the pseudonym "Student" because he worked for Guinness

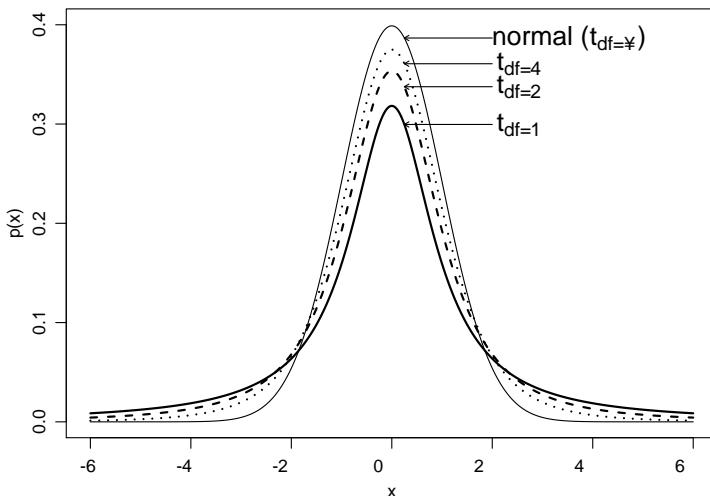


Figure 15.4: t distributions, with a normal distribution superimposed for comparison. The t distribution has tails relative to the normal. The parameter (“degrees of freedom”) controls the relative height of the tails.

Brewery) prohibited publication of any research that might be proprietary or imply problems with their product (such as variability in quality). Therefore the distribution is often referred to as the Student t distribution.

In R, the t density at x is specified by $dt(x, df)$, where df is a parameter called the degrees of freedom. (The R function also has an optional argument called ν , the non-centrality parameter. This option is not implemented in BUGS.) The effect of the df argument is shown in the Figure 15.4. When df is small, the t distribution has tails that are taller than normal. As df approaches infinity, the t distribution becomes normal. The df parameter is a continuous value greater than or equal to one. (Readers might be familiar with the df parameter as related to sample size when sampling distribution, but the t distribution is not used as a sampling distribution here, therefore df is not restricted to being an integer.) In BUGS, the density is specified by $dt(mu, tau, df)$. The mu and tau arguments in BUGS have no equivalent arguments in R. All this linearly shift the t distribution, such that $dt(mu, tau, df)$ in BUGS corresponds to $dt((x-mu)*sqrt(tau), df)$ in R.

Although the t distribution was originally conceived as a sampling distribution for NHST, we will use it instead as a convenient model of data outliers (as is often done; e.g., Damgaard, 2007; Meyer & Yu, 2000; Tsionas, 2002). Outliers are simply data values that fall unusually far from the model’s expected value. Real data often contain outliers, presumably because some extraneous influences have perturbed the data values. Sometimes these extraneous influences can be identified and the affected data values can be corrected. But usually we have no way of knowing whether the suspected outlying value was caused by an extraneous influence, or is a genuine reflection of the target being measured. Instead of deleting suspected outliers from the data according to some arbitrary criterion, we retain all the data but use a likelihood function that is less affected by outliers than is the normal distribution.

Figure 15.5 shows examples of how the distribution is robust against outliers. The curves show the maximum likelihood estimates (MLEs) of the parameters for the

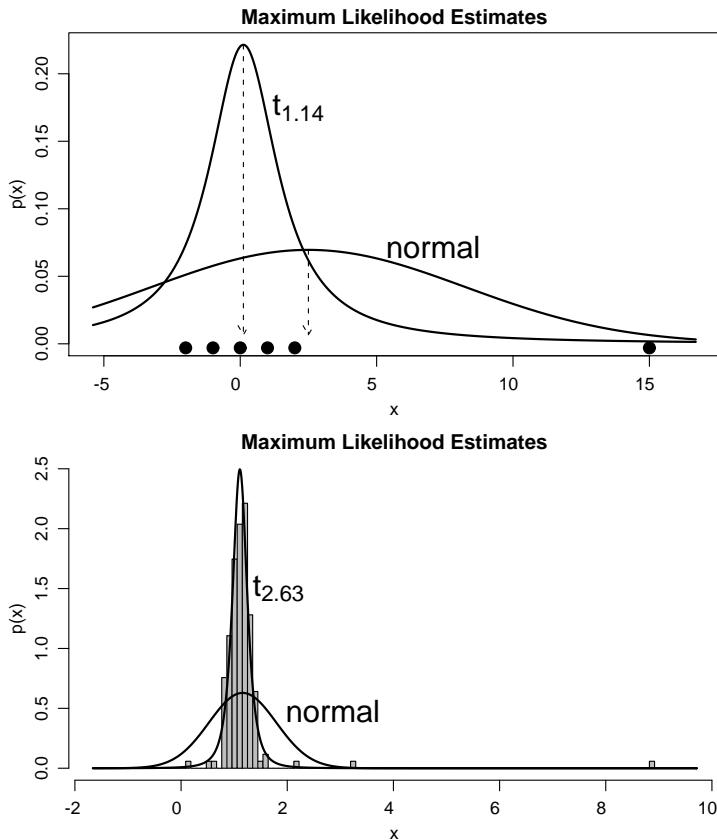


Figure 15.5: The maximum likelihood estimates of normal distributions t to the data shown. Upper panel shows “toy” data to illustrate that the normal accommodates an outlier only by enlarging its standard deviation and, in this case, by shifting its mean. Lower panel shows actual data to illustrate the realistic effect of outliers on estimates of the normal.

normal distributions. More formally, for the given data \mathbf{x} , parameter values were found for the normal that maximized $\phi(D_j; \mu, \sigma^2)$, and parameter values were found for the t distribution that maximized $\phi(D_j; \mu, \sigma^2, df)$, and the curves of those MLEs are plotted with the data. The upper panel of Figure 15.5 shows “toy” data to illustrate that the normal is strongly influenced by an outlier, while the distribution remains centered over the bulk of the data. The lower panel of Figure 15.5 uses realistic data that indicates levels of inorganic phosphorous, measured in milligrams per deciliter, in 177 human subjects aged 65 or older. The authors of the data (Holcomb & Spalsbury, 2000) intentionally altered a few data points to reflect typical transcription errors and illustrate methods for detecting and correcting such errors. We instead assume that we no longer have access to records of the original individual measurements, and must model the corrected data set. The distribution accommodates the outliers and fits the distribution of data much better than the normal.

The t distribution is useful as a likelihood function, for modeling outliers at the level of observed data, but it is also useful for modeling outliers at higher levels in a hierarchical prior. We will encounter such applications when we model individual differences or multiple predictors, e.g., in Section 16.2.

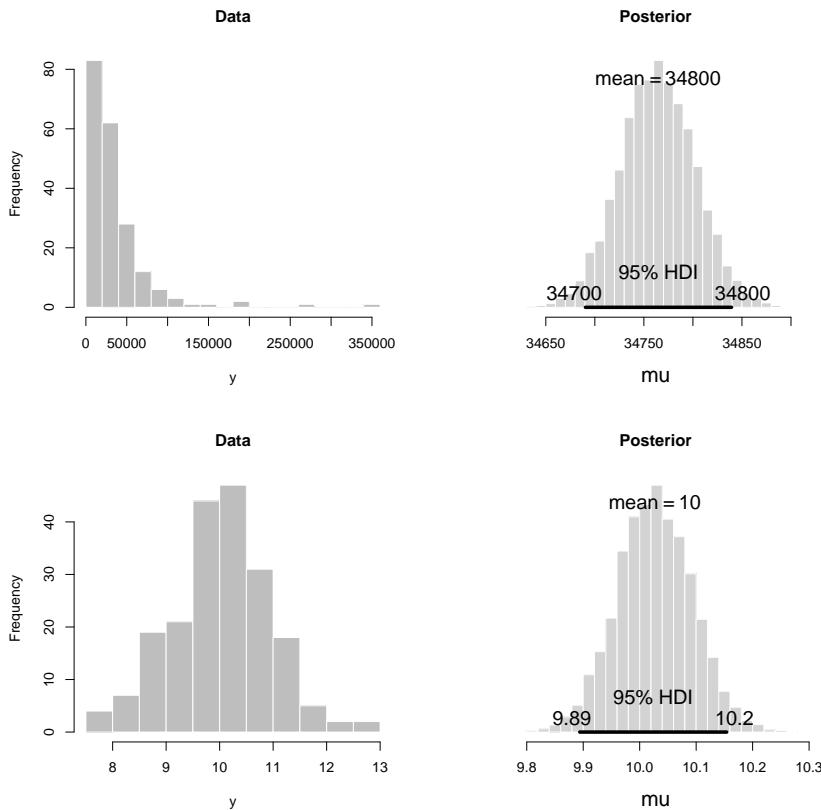


Figure 15.6: Upper left shows skewed, non-normal data. Upper right shows the posterior distribution of μ when the data are modeled (inappropriately!) with a normal likelihood. Bottom left shows the same data transformed, via a logarithm. Bottom right shows the posterior distribution of μ when the transformed data are modeled (appropriately) with a normal likelihood.

15.1.4 When the data are non-normal: Transformations

Suppose there has been discovered a new species of insects in dark caves, and the insects make ultra high frequency sound bursts. One theory says that the sounds should have an average frequency of 22,000 Hz (cycles per second). An entomologist has measured the frequencies of 200 sound bursts from these insects, and would like to know if a mean of 22,000 Hz is credible.

The upper left panel of Figure 15.6 shows a histogram of the data, and the upper right panel shows the posterior distribution of μ from the normal-likelihood model of Figure 15.3. (Assume that the prior was uncontroversial.) According to the posterior, the data exclude a value of 22,000. Therefore the entomologist concludes that the 22,000 Hz theory is not credible.

The proponent of the theory says that the theory was mistimed. Instead of scaling the insect sounds by their raw frequency, they should consider their perceived pitch, which is merely the logarithm of frequency. The natural log of 22,000 is 10. Thus, the theory claims that the sound bursts should have an average (frequency) of 10. The lower left panel of Figure 15.6 shows the logarithm of the data as the upper left panel. When the normal-likelihood model is applied to these data, resulting posterior on

is shown in the lower right of Figure 15.6. (Assume that the prior was uncontroversial.) According to this posterior, the data do not exclude a value of 10. Indeed, the value of 10 is very near the mode of the posterior distribution.

Which analysis do we trust? Is the theory disconfirmed or upheld? In this case, the theory is clearly upheld, because the analysis based on raw frequency is not appropriate. The problem with the analysis based on raw frequency is that the normal-likelihood model assumes that the data are normally distributed, but the frequencies are obviously extremely skewed and non-normal. Graphically, the normal-likelihood model in Figure 15.2 shows that the data are generated by a normal distribution. The parameters are only meaningful to the extent that the normal likelihood actually describes the distribution of the data. Because of the severe skew in the raw-frequency data, the central-tendency parameter is being pulled far to the right by the extreme values far above 20,000, for which there are no symmetrically distributed values below 2000, as is assumed by a symmetric normal distribution. The logarithmically transformed data, on the other hand, are apparently much more normally distributed, and therefore the normal-likelihood model is a viable description of the data.

The moral of the story is this: When using a normal-likelihood model, the data should be at least roughly normally distributed, otherwise the model is a poor description of the data and the model's parameter values may be meaningless. Where there is only moderate non-normality in the data, the misrepresentation by the model parameters is often hardly noticeable, but when there is severe non-normality in the data, then the misrepresentation may be substantial.

There are many real-world situations in which data are not normally distributed. A prominent example is response times: Ask a person to press a button quickly as possible after a light comes on. Repeat this many times. The response times will be strongly skewed toward higher values. This skew makes intuitive sense, because the response time can't be much faster than one or two hundred milliseconds. There are many reasons to be slow on some trials. Another example is shown in Figure 10.5, which reveals that the longevity of lab rats on a restricted diet is skewed to the left.

How should we analyze non-normal data? There are two possibilities. One method perseveres with a normal-likelihood function, transforming the data so that they are very nearly normal. There is nothing wrong with transforming data, as long as all the data are transformed the same way, and the transformation preserves order. A transformation that preserves order is called "monotonic". Examples of monotonic transformations include the exponential function, the logarithm (for positive valued data), and the cubic (i.e., y^3). An example of a non-monotonic transformation is a sine wave function, when the original data span several cycles of the sine wave. Monotonic transformations merely re-scale the data. The measurement scale is arbitrary. Earthquakes could be measured in terms of raw energy, or in terms of the base-10 logarithm of energy, which is the Richter scale. Temperatures can be measured in terms of Fahrenheit or Celsius scales. Distances can be measured in terms of the English or Metric systems. Acoustic vibrations can be measured by frequency or by perceived pitch, which is essentially the logarithm of frequency. The analysis should use whatever transformation makes the data normally distributed, so that the data respect the assumption of the likelihood model and thereby make the parameter estimates meaningful.

The second approach to analyzing non-normal data is to replace the normal likelihood function, and use a more appropriate likelihood function instead. For example, Rouder, Lu, Speckman, Sun, and Jiang (2005) used a Weibull distribution as the likelihood function

for modeling skewed response time distributions. This approach requires familiarity with various non-normal probability distributions. Moreover, use in BUGS, this method also requires that the desired likelihood function can be specified in BUGS. Because of these extra complications, people usually first try the method of transforming data to normal. But a non-normal likelihood function may be preferred when reasonable data transformation succeeds, or when transformed scales are too awkward to work with when a theory posits a specific non-normal distribution for generating data.

There is one more important caveat when transforming data (changing likelihood functions): The prior must be appropriate to the transformed data. For example, when the insect sounds were measured in terms of raw frequency, the prior should be appropriate for data in the broad vicinity of 22,000. But when the insect sounds were measured by $\log(\text{frequency})$, then the prior should be appropriate for data in the general vicinity of 10.

Exercise 15.2 provides a realistic example in which we attempt to estimate the underlying mean that generated a set of scores. It also asks to consider the choice of prior, the robustness of the posterior, and whether a likelihood seems to be an appropriate model of the data. Exercise 15.3 gives another realistic example in which the data are skewed.

The section previous to this one discussed outliers. Can we “transform away”? In simple, one-group data sets, the answer, in principle, is yes. We merely fashion a transformation that arbitrarily compresses the extreme tails of the distribution in just the right way. But then the prior should take into account the transformation. Unfortunately, for more typical complex situations, it is difficult to invent a transformation that appropriately compresses all outliers simultaneously. Therefore, transformations are typically used to adjust large-scale skew in the data, and not for outliersthat, on the other hand, are typically addressed by using tall-tail distributions, that is, the distribution.

15.2 Repeated measures and individual differences

In many real-world applications, we have repeated measurements on the same individual, although we are mainly interested in the central tendency of the group. For example, we might be interested in the typical blood pressure of employees in a company. We can repeatedly measure the blood pressure of each employee at random times of day and on random different days. Our overall goal is to estimate the typical blood pressure of the group as a whole, but now we have repeated measurements from the individual subjects.

We assume that the individual subjects have been randomly selected from the pool of all possible subjects in the population of interest (the pool of all employees in the company). We also assume that the repeated measurements on each subject are mutually independent. This is sometimes a perilous assumption because it is probably only approximately true even if we are careful to design the measurement procedure with the assumption in mind. For example, suppose we are measuring someone's blood pressure, with the goal of getting measurements that represent their typical blood pressure. We could measure the blood pressure several times in rapid succession, but presumably those measurements would be highly correlated because their blood pressure would not change much in only a few seconds duration. Instead, we measure blood pressure at random times across hours and days. The goal is to make repeated measurements far enough apart in time so that we can reasonably assume that the repeated measurements are independent of each other. But we don't want the repeated measurements to be so far apart

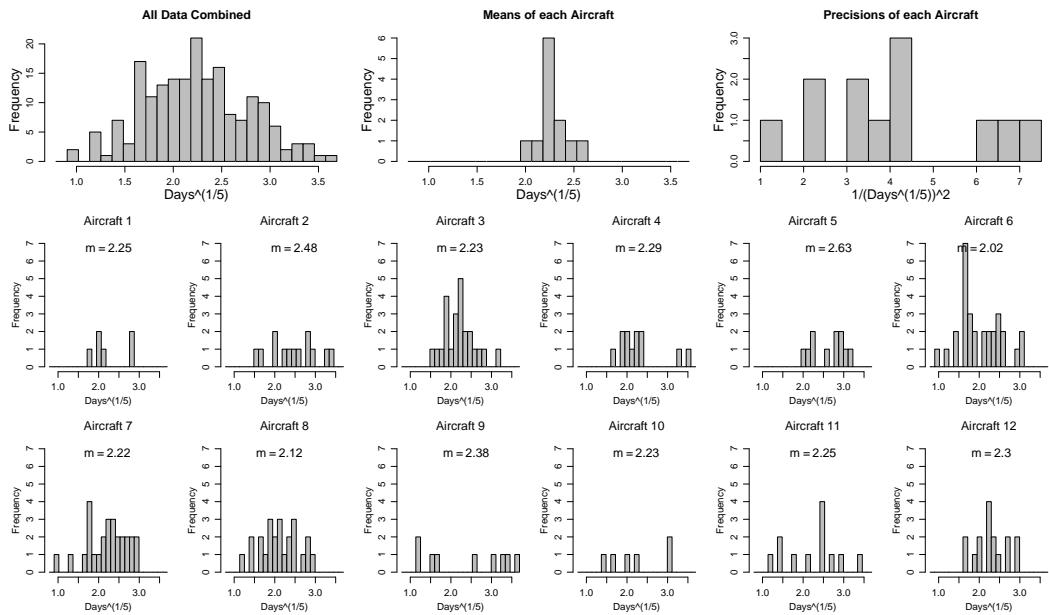


Figure 15.7: Number of days between failures of air-conditioning systems in aircraft. Notice the scale on the x-axis. The lower 12 panels show data from individual aircraft. The upper-left panel shows all those combined into one histogram. The upper-middle panel shows a histogram of the aircraft means.

in time that the person's underlying typical value has changed, e.g., from dramatic change in body fat. Thus, we want the measurements to be far enough apart that they have minimal correlation, but we want the measurements to be close together that they represent a stationary underlying propensity. There is no single correct sampling procedure; it depends on theoretical considerations for the specific application domain.

As another example with real data, consider the time between successive failures of the air conditioning system in each aircraft in a fleet of Boeing 707 jet airplanes (Proschan, 1963). The time between successive failures is of interest that managers can appropriately schedule preventative maintenance and maintain inventory of spare parts. Thus, the emphasis in this situation is estimation of inter-failure duration. We will assume that successive failures of the same air conditioner are independent of each other. One rationale for this assumption is that the repair of a failed system resets the system so that it has no memory for the previous failure. The next failure depends on other factors, such as the basic design of the air conditioner and the typical weather conditions in which the individual aircraft flies.

The data from 12 aircraft are displayed in Figure 15.7 (they omit one aircraft that had only two observations, as was done by Hand et al., 1994, p. 480). Notice that the data have been transformed from their original scale of days⁻⁵, i.e., the 5th root of days. This transformation was necessary because the raw data were extremely skewed and I wanted to use a normal likelihood function, for illustrative purposes (Proschan, 1963, used a more theoretically informed model). A different transformation, such as a logarithm, could have been used instead, but the 5th root yielded nicely symmetric data distributions. Notice that all the data were transformed the same way; there were no different transformations applied to different aircraft. As will be explained when the full model is described, the goal of the transformation was to make the data from each individual aircraft be approximately

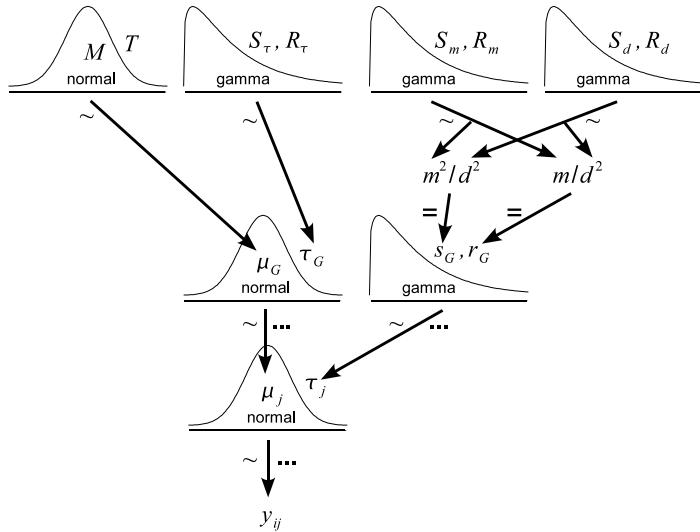


Figure 15.8: A model of hierarchical dependencies for repeated measures from each of several subjects drawn from the same group. The datum from the j^{th} subject is denoted y_{ij} . It is generated by a normal distribution that has a subject mean of μ_j and subject precision of τ_j . The subject means are generated by a normal distribution centered on the overall mean μ_G with a precision of τ_G . These group-descriptive parameters μ_G and τ_G , are themselves being estimated from the data, hence they have prior distributions as shown at the top of the hierarchy. The subject precisions, τ_j , are generated by a gamma distribution which has parameters s_G and r_G . These parameters are themselves being estimated, therefore they have prior distributions as shown in the top right of the diagram.

normal, and to simultaneously make the distribution of all means be approximately normal.

15.2.1 Hierarchical model

To model this situation, we will assume that data from j^{th} individual are normally distributed around an underlying individual mean, denoted μ_j . The dispersion of repeated measurements from that individual is described by the precision τ_j . We will further assume that the individual propensities, μ_j , are normally distributed around a group mean μ_G . The dispersion of the individual means around the group mean is described by the precision τ_G . It is these group-level parameters in which we are primarily interested. Because we want to estimate μ_G and τ_G , we need to specify their prior-belief distributions.

Just as the individual propensities are generated by group-level distributions, the individual precisions τ_j are also generated by group-level distributions. This is to admit the variability of repeated measurements in an individual dependent on the group they are in. For example, a company might engender employees with very stable blood pressures, i.e., all the employees tend to have high precisions in their individual distributions of blood pressure. As another example, an air-conditioner manufacturer might build systems that fail very erratically, i.e., all the systems tend to have low precisions in their individual distributions of inter-failure duration. Therefore the individual μ_j come from a group-level gamma distribution which has shape and rate parameters we can estimate from the

data. Therefore the shape and rate parameters of the gamma distribution must themselves be given higher level prior distributions.

The entire model is illustrated in Figure 15.8. It shows the hierarchical dependency of individual parameters j on the group-level parameters G , and the dependency of the individual parameters j on the group-level parameters r_G . The group-level parameters themselves have higher-level distributions as shown. The higher level distributions are necessary because it is the group-level parameters that we want to estimate, and therefore we must express our prior uncertainty in those group-level parameters by the higher-level distributions. The higher-level distributions for r_G are re-parameterized in terms of the mean (n) and standard deviations (s) of the gamma density for j . This reparameterization is performed merely for convenience in interpreting the values of r_G .

At the top level, the constants for the hyper-priors must be specified values. These values must be appropriate for the domain and data at hand. In other words, the hyper-prior constants should be at least mildly informed by the situation (or richly informed if audience-agreeable prior information is available). For example, we know from common sense that air conditioners will be manufactured to operate without failure for a duration on the order of 10^2 days, with a range of perhaps 10^1 to 10^4 days. Therefore the hyperprior on n should reasonably specify a distribution in this range (and not something like seconds or eons!). Moreover, it is important to remember that the air-conditioner data were transformed to a different scale, and therefore the hyperprior constants need to be appropriate for the transformed scale.

Before proceeding with parameter estimation, it is important to verify that the basic distributional assumptions of the model are not being violated. At the lowest level of the model, the normal likelihood function assumes that the data generated by each individual aircraft are normally distributed. A quick look at Figure 15.7 shows that the twelve individual aircraft have data that are reasonably unimodal and symmetric (unlike the raw, untransformed data). At the next higher level in the model, the normal distribution that generates j assumes that the central tendencies of the aircraft are already distributed. Inspection of the top-middle graph of Figure 15.7 verifies that the 12 aircraft means are unimodal and symmetric. Finally, the gamma distribution that generates r_G assumes that the precisions of the 12 aircraft can be described by a gamma distribution. Inspection of the top-right graph of Figure 15.7 verifies that the 12 aircraft precisions may be reasonably described this way.

If there appear to be prominent outliers in the data, then the data might be better modeled by a t distribution than by a normal distribution. As was discussed in connection with Figure 15.5, the t distribution is far less sensitive to outliers than the normal. The novel point here is that outliers might occur at different levels. There might be outliers of interfailure durations within each individual aircraft. These would be visible in the plots of individual aircraft data. Additionally, there might be outliers at the level of means for each aircraft. These outliers would be visible in the plot of aircraft means. In the hierarchical model of Figure 15.8, these considerations imply that either both of the two lower normal distributions, involving j or G , might be replaced with a t distribution.

15.2.2 Implementation in BUGS

It is straightforward to implement the model in BUGS. Every node in the dependency diagram of Figure 15.8 has a corresponding specification in the BUGS model: (SystemsBrugs.R)

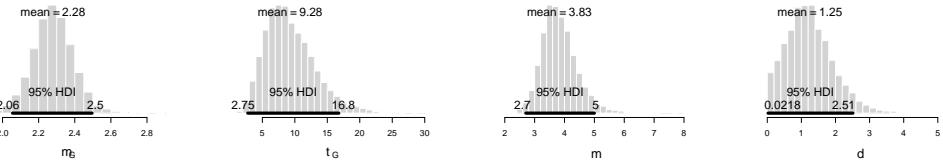


Figure 15.9: Posterior distribution for the top-level hyperparameters in Figure 15.8, for the data in Figure 15.7.

```

9   model {
10    for( i in 1 : Ndata ) {
11      y[i] ~ dnorm( mu[ subj[i] ] , tau[ subj[i] ] )
12    }
13    for ( j in 1 : Nsubj ) {
14      mu[j] ~ dnorm( muG , tauG )
15      tau[j] ~ dgamma( sG , rG )
16    }
17    muG ~ dnorm( 2.3 , 0.1 )
18    tauG ~ dgamma( 1 , .5 )
19    sG <- pow(m,2) / pow(d,2)
20    rG <- m / pow(d,2)
21    m ~ dgamma( 1 , .25 )
22    d ~ dgamma( 1 , .5 )
23  }
```

Notice the use of nested indexing in line 11. In the BUGS `model` cation, the index is the overall row of the data matrix, not the measure within a subject. Nested indexing is used because the data matrix is structured with each row containing a datum and the subject number (i.e., aircraft) from which the datum was obtained. This data structure accommodates the fact that different aircraft contributed different numbers of data points. The complete program is listed in Section 15.4.2 (`systemsBrugs.R`).

Aspects of the resulting posterior distribution are shown in Figure 15.9. The left panel shows the marginal distribution on the group-level parameter, which has a posterior mean of 2.28 days⁵, which corresponds to 61.6 days. In other words, on average, two months elapsed between failures of the air conditioning system on these airplanes. In the second panel, the value of τ_G indicates the standard deviation across aircraft of the average inter-failure duration. The posterior mean is 9.32, corresponding to a standard deviation of $\sqrt{9.32} = 0.33$ days⁵. The parameters μ_G and τ_G describe the distribution of aircraft means μ_j , and thus the posterior is saying that among the most credible descriptions of the data distribution in the top-middle panel of Figure 15.7 is a normal distribution with mean 2.28 and standard deviation of 0.33. The remaining posterior marginals in Figure 15.9 show credible values of m and d , for the gamma distribution that describes how the precisions vary across aircraft. The posterior mean is 3.83 and the posterior median is 1.27, which says that among the most credible descriptions of the data distribution in the top-right panel of Figure 15.7 is a gamma distribution with mean of 3.83 and standard deviation of 1.27.

15.3 Summary

This chapter was intended to explore the normal likelihood function. The purpose was to explore concepts focused on the normal likelihood when σ^2 is simply a constant, not involving additional predictors and parameters. With this foundation in place, subsequent chapters can focus on concepts regarding additional predictors. This chapter made several main points:

When the precision of the normal likelihood function is τ^2 (not estimated), then the “update rules” for the mean and precision of a normal prior are analytically simple. These formulas also reveal why it is convenient for Bayesians to parameterize a normal distribution by its precision instead of by its standard deviation.

It is typical to put a normal prior on the mean of a normal likelihood, and a gamma prior on the precision of a normal likelihood, because of considerations of conjugacy for simple situations. Use of these priors is not necessary, however. It is especially easy to program arbitrary priors in BUGS (although not all will work equally efficiently for sampling).

When there are repeated measures within individuals, either repeated measures within individuals and/or the means across individuals might be modeled as normal distributions. Again, this hierarchical model can be easily programmed in BUGS.

Whenever applying a normal likelihood function, it is important to check that the data are reasonably normally distributed, i.e., at least roughly unimodal and symmetric. Otherwise the parameter values may be misleading. If the data are severely non-normal, they might be transformed to normal, or else a normal likelihood function could be used instead.

Data that have prominent outliers might be better modeled with a distribution than a normal distribution. This idea was introduced in this chapter but examples of applying it in BUGS wait until subsequent chapters.

15.4 R code

15.4.1 Estimating the mean and precision of a normal likelihood

Notice that the data in this program are randomly generated from a normal distribution in R, on line 30. In R, `dnorm` and `rnorm` are parameterized by mean and standard deviation, unlike BUGS, in which `dnorm` is parameterized by mean and precision.

The program also initializes the chains intelligently setting them start amidst the posterior, instead of starting randomly according to the `method` of initializing chains was introduced in Section 9.5 (`conBrugs.R`). Lines 51–52 show that the initial value of μ is the mean of the data, and the initial value of τ^2 is the precision of the data. If the prior is fairly vague, and the data are numerous, then the posteriors will be near the parameter values that maximize the likelihood of the data. Therefore, start the chains near those maximum-likelihood values. We still need to allocate a burn period, however, because the chains needs to independently diverge from each other and settle in to the true posterior.

(`YmetricXsingleBrugs.R`)

```

1  graphics.off()
2  rm(list=ls(all=TRUE))
3  library(BRugs)          # Kruschke, J. K. (2010). Doing Bayesian dat      a analysis:
4                                # A Tutorial with R and BUGS. Academic Press / Elsevier.
5  #-----
6  # THE MODEL.
7  modelstring = "
8  # BUGS model specification begins here...
9  model {
10    # Likelihood:
11    for( i in 1 : N ) {
12      y[i] ~ dnorm( mu , tau ) # tau is precision, not SD
13    }
14    # Prior:
15    tau ~ dgamma( 0.01 , 0.01 )
16    mu ~ dnorm( 0 , 1.0E-10 )
17  }
18  # ... end BUGS model specification
19  " # close quote for modelstring
20  writeLines(modelstring,con="model.txt")
21  modelCheck( "model.txt" )
22
23  #-----
24  # THE DATA.
25
26  # Generate random data from known parameter values:
27  set.seed(47405)
28  trueM = 100
29  trueSD = 15
30  y = round( rnorm( n=500 , mean=trueM , sd=trueSD ) ) # R dnorm uses mean and SD
31
32  dataList = list(
33    y = y ,
34    N = length( y )
35  )
36
37  # Get the data into BRugs: (default filename is data.txt).
38  modelData( bugsData( dataList ) )
39
40  #-----
41  # INTIALIZE THE CHAINS.
42
43  nchain = 3
44  modelCompile( numChains = nchain )
45
46  automaticInit = F # TRUE or FALSE
47  if ( automaticInit ) {
48    modelGenInits() # automatically initialize chains from pr      or
49  } else {
50    genInitList <- function() { # manually initialize chains ne      ar the data
51      list( mu = mean( dataList$y ) ,
52            tau = 1 / sd( dataList$y )^2 )
53    }
54    for ( chainIdx in 1 : nchain ) {
55      modelInits( bugsInits( genInitList ) )
56    }
57  }
58
59  #-----
```

```

60 # RUN THE CHAINS
61
62 # burn in
63 BurnInSteps = 500
64 modelUpdate( BurnInSteps )
65 # actual samples
66 samplesSet( c( "mu" , "tau" ) )
67 stepsPerChain = 2000
68 thinStep = 1
69 modelUpdate( stepsPerChain , thin=thinStep )
70
71 #-----
72 # EXAMINE THE RESULTS
73
74 filenamert = "YmetricXsingleBrugs"
75
76 source("plotChains.R")
77 muSum = plotChains( "mu" , saveplots=F , filenamert )
78 sigmaSum = plotChains( "tau" , saveplots=F , filenamert )
79
80 muSample = samplesSample( "mu" )
81 tauSample = samplesSample( "tau" )
82 sigmaSample <- 1 / sqrt( tauSample ) # Convert precision to SD
83
84 source("plotPost.R")
85 windows()
86 plotPost( muSample , xlab="mu" , breaks=30 , main="Posterior" )
87 dev.copy2eps(file=paste(filenamert,"PostMu.eps",sep=""))
88
89 nPts = length(muSample) ; nPtsForDisplay = min( nPts , 2000 )
90 thinIdx = seq( 1 , nPts , nPts / nPtsForDisplay )
91 windows()
92 plot( muSample[thinIdx] , sigmaSample[thinIdx] , col="gray" ,
93       xlab="mu" , ylab="sigma" , cex.lab=1.5 , main="Posterior" , log="y" )
94 points( mean(muSample) , mean(sigmaSample) , pch= "+" , cex =2 )
95 text( mean(muSample) , mean(sigmaSample) ,
96       bquote( .(round(mean(muSample),1)) ** " " .(round(mean(sigmaSample),1)) ) ,
97       adj=c(.5,-0.5) )
98 dev.copy2eps(file=paste(filenamert,"PostMuSigma.eps",sep=""))
99
100 #-----

```

15.4.2 Repeated measures: Normal across and normal within

Below is the complete program for analyzing the data from Figure 15.7, regarding inter-failure durations of aircraft air conditioners. A glimpse of the posterior was shown in Figure 15.9.

(SystemsBrugs.R)

```

1 graphics.off()
2 rm(list=ls(all=TRUE))
3 library(BRugs)           # Kruschke, J. K. (2010). Doing Bayesian data analysis:
4                           # A Tutorial with R and BUGS. Academic Press / Elsevier.
5
6 # THE MODEL.
7 modelstring =
8 # BUGS model specification begins here...
9 model {

```

```

10   for( i in 1 : Ndata ) {
11     y[i] ~ dnorm( mu[ subj[i] ] , tau[ subj[i] ] )
12   }
13   for ( j in 1 : Nsubj ) {
14     mu[j] ~ dnorm( muG , tauG )
15     tau[j] ~ dgamma( sG , rG )
16   }
17   muG ~ dnorm( 2.3 , 0.1 )
18   tauG ~ dgamma( 1 , .5 )
19   sG <- pow(m,2) / pow(d,2)
20   rG <- m / pow(d,2)
21   m ~ dgamma( 1 , .25 )
22   d ~ dgamma( 1 , .5 )
23 }
24 # ... end BUGS model specification
25 " # close quote for modelstring
26 writeLines(modelstring,con="model.txt")
27 modelCheck( "model.txt" )

28 #
29 #-----#
30 # THE DATA.

31
32 # Load the aircraft data:
33 load( "Systems.Rdata" ) # loads dataMat
34 nSubj = length( unique( dataMat[, "Aircraft"] ) )
35 # Transform the data:
36 DaysTransf = dataMat[, "Days"]^(1/5)
37 dataMat = cbind( dataMat , DaysTransf )
38 colnames( dataMat ) = c( colnames( dataMat )[1:3] , "DaysTra    nsf" )

39
40 # Specify data, as a list.
41 dataList = list(
42   y = dataMat[, "DaysTransf"] ,
43   subj = dataMat[, "Aircraft"] ,
44   Ndata = NROW(dataMat) ,
45   Nsubj = nSubj
46 )
47 # Get the data into BRugs: (default filename is data.txt).
48 modelData( bugsData( dataList ) )

49 #
50 #-----#
51 # INTIALIZE THE CHAINS.

52
53 # First, compile the model:
54 nchain = 10
55 modelCompile( numChains = nchain )

56
57 modelGenInits() # works when the priors are not too flat
58
59 #
60 #-----#
61 # RUN THE CHAINS

62 # burn in
63 BurnInSteps = 1000
64 modelUpdate( BurnInSteps )
65 # actual samples
66 samplesSet( c( "muG" , "taug" , "mu" , "tau" , "m" , "d" ) )
67 stepsPerChain = ceiling(10000/nchain)
68 thinStep = 100

```

```

69 modelUpdate( stepsPerChain , thin=thinStep )
70
71 #-----
72 # EXAMINE THE RESULTS
73
74 source("plotChains.R")
75 source("plotPost.R")
76 filenamert = "SystemsBrugs"
77
78 # Examine chains for convergence and autocorrelation:
79 muSum = plotChains( "muG" , saveplots=F , filenameroot=filenamert )
80 tauSum = plotChains( "tauG" , saveplots=F , filenameroot=filenamert )
81 mSum = plotChains( "m" , saveplots=F , filenameroot=filenamert )
82 dSum = plotChains( "d" , saveplots=F , filenameroot=filenamert )
83 mu1Sum = plotChains( "mu[1]" , saveplots=F , filenameroot=filenamert )
84 tau1Sum = plotChains( "tau[1]" , saveplots=F , filenameroot=filenamert )
85
86 # Extract chains from BUGS into R:
87 muGsample = samplesSample( "muG" )
88 tauGsample = samplesSample( "tauG" )
89 mSample = samplesSample( "m" )
90 dSample = samplesSample( "d" )
91 muSample = NULL
92 tauSample = NULL
93 for ( sIdx in 1:nSubj ) {
94   muSample = rbind( muSample , samplesSample( paste("mu[",s_idx,"]",sep="") ) )
95   tauSample = rbind( tauSample , samplesSample( paste("tau[    ",sIdx,"]",sep="") ) )
96 }
97
98 # Plot the aircraft mu:
99 windows(15,6)
100 layout( matrix( 1:nSubj , nrow=2 , byrow=T ) )
101 for ( sIdx in 1:nSubj ) {
102   plotPost( muSample[sIdx,] , xlab=bquote(mu[.(sIdx)]) )
103 }
104 dev.copy2eps(file=paste(filenamert,"PostMu.eps",sep=""))
105
106 # Plot the aircraft tau:
107 windows(15,6)
108 layout( matrix( 1:nSubj , nrow=2 , byrow=T ) )
109 for ( sIdx in 1:nSubj ) {
110   plotPost( tauSample[sIdx,] , xlab=bquote(tau[.(sIdx)]) , HDItextPlace=.3 )
111 }
112 dev.copy2eps(file=paste(filenamert,"PostTau.eps",sep=""))
113
114 # Plot the hyperdistributions:
115 windows(15,3)
116 layout( matrix(1:4,ncol=4) )
117 plotPost( muGsample , xlab=expression(mu[G]) , breaks=30      )
118 plotPost( tauGsample , xlab=expression(tau[G]) , breaks=      30 )
119 plotPost( mSample , xlab=expression(m) , breaks=30      )
120 plotPost( dSample , xlab=expression(d) , breaks=30 , HDItextPlace=.1 )
121 dev.copy2eps(file=paste(filenamert,"PostHyper.eps",sep=""))
122
123 #-----

```

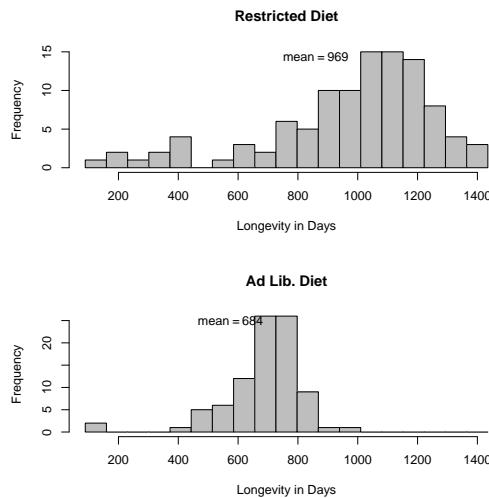


Figure 15.10: For Exercise 15.3.

15.5 Exercises

Exercise 15.1. [Purpose: View the prior from BUGS.] For the program in Section 15.4.1 (`YmetricXsingleBrugs.R`), generate graphs of prior distribution, analogous to Fig 15.3. To do this, comment out the data in the data specification, `xprior`, in Section 8.5.1, p. 141. But also beware of the following: Automatically initialize the chains from the prior, and, when plotting the results, comment out the `plotChains` commands because the MCMC sigma values are too extreme for some of the BUGS graphics.

Exercise 15.2. [Purpose: A realistic example of estimating a single mean, with consideration of priors.] University students who agreed to participate in a problem-solving experiment were also tested for their vocabulary level, using the Wechsler Adult Intelligence Scale Revised (WAIS-R), which is normed to have a general-population mean of 10 and standard deviation of 3. Here are the data from the university students (from Hand et al., 1994, set #392, p. 322):

```
14 11 13 13 13 15 11 16 10 13 14 11 13 12 10 14 10 14 16 14 14 11 11 11 13
12 13 11 11 15 14 16 12 17 9 16 11 19 14 12 12 10 11 12 13 13 14 11 11 15
12 16 15 11
```

(A) Are the data roughly normal, i.e., essentially unimodal and symmetrically distributed? (No formal analysis is required here; only an `eye` assessment is expected.) Therefore, can a normal likelihood function be applied?

(B) Discuss the rationale for a prior distribution. Justify the constants you choose for the hyperprior. Hint: The results of the analysis will be presented to a skeptical audience, who may doubt any claims about how prior knowledge of the general population informs an analysis of university students. Therefore your prior must be very vague and widely dispersed around the general population values.

(C) Is the mean vocabulary score of the university students different from the general population mean of 10? Report the 95% HDI for the various priors you considered in the previous part.

Exercise 15.3. [Purpose: A realistic example of estimating a mean, with consideration of whether a normal likelihood distribution is appropriate.] Suppose we know that the mean life of a rat that eats ad lib is 700 days. This value is, in fact, about right for lab rats. Figure 15.10, which shows data from R. L. Berger et al. (1988), as reported in Had (1994, data set #242), and which are available in the RatLives.Rdata . This is an Rdata file, which stores values in compressed format, not text format. To get the data into R, type `load("RatLives.Rdata")`. It will look like nothing happens when you type that command, but the variables are now loaded into R; you can see the variables that R knows about by typing `ls()`. You will see listed the two vectors `adlibDiet` and `restrictedDiet`. When rats are placed on a restricted diet, their longevity can be affected. Figure 15.10 shows the results.

(A) Using the raw data from the restricted-diet longevities and its 95% HDI. Be sure to report the prior you used. Hint: The HDI extends from about 917 to 1020 days.

(B) Is it appropriate to apply a normal likelihood function to the data? Transform the data by squaring the longevities. Estimate and its 95% HDI (now in days squared). Be careful to use an appropriate prior! Hint: The HDI extends from about 967 to 1053 days squared.

(C) Why is the first estimate lower than the second estimate? Which estimate is more appropriate and why?

Exercise 15.4. [Purpose: Think about specifying an informed prior, and non-normal, non-gamma priors.] In the program of Section 15.4.1 (metricXsingleBrugs.R), the prior on the group mean uses an extremely small precision. This prior is silly, because we know that individual IQ scores from the general population should be near 100, with a standard deviation of about 15. Special populations might have IQ scores consistently above or below 100, but almost certainly within, say, 100 points above and below 100. Consider the following expressions of the prior belief on

- (i) `mu = dnorm(100 , sd=50)` # that's $sd=50$, so $\tau = \text{pow}(50,-2)$
- (ii) `mu = dgamma(4 , 0.04)`
- (iii) `mu = dunif(0 , 200)`
- (iv) `mu = dgamma(3 , 0.03)`

(A) Plot the densities of (i) and (ii) on the same graph, superposed. Hint:
`mu = seq(-50 , 300 , length=501)`

```
plot( mu , dnorm( mu , 100 , 50 ) , type="l" , ylim=c(0,0.01) )
lines( mu , dgamma( mu , 4 , 0.04 ) )
```

What do the densities of (i) and (ii) have in common? (Hint: Compute the mean and sd of the gamma distribution.) Should either of (i) or (ii) be preferred over the other? (Hint: Consider negative IQ scores, which should not be allowed.)

(B) Plot the densities of (iii) and (iv) on the same graph, superposed. What do the densities of (iii) and (iv) have in common? Should one of (iii) or (iv) be preferred over the other? (Hint: Consider IQ scores greater than 200, which should not be allowed.)

In the program of Section 15.4.1 (metricXsingleBrugs.R), the prior on the group precision is a general diuse distribution. This diuse prior is silly, because we know that individual IQ scores have a standard deviation around 15, (a precision of $1/15^2 = 0.004444\ldots$) in the general population. The SD might be a bit smaller in special populations.

(C) Suppose we believe that the smallest SD we would find in a sized population is 5. What is the corresponding precision? (Notice that larger than $1/15^2$.)

(D) We want to create a gamma distribution for precision that has a corresponding

to $1=15^2$, such that most of the gamma distribution is less than the precision corresponding to an IQ precision of 45^2 , because, as mentioned in the previous part, we don't think that precisions much greater than 5^2 are tenable. Therefore, to be sure that our prior encompasses that maximum but still allows for considerable uncertainty, we will set the standard deviation of the gamma distribution half of the difference between the precisions determined in the previous two parts. What is the value of the difference between the precisions determined in the previous two parts? (Hint: The answer is 0.01777... Explain.)

(E) What are the values of the shape and rate parameters for a gamma distribution that has mean of 0.00444 and standard deviation of 0.01777?

(F) Re-run the program of Section 15.4.1 (metricXsingleBrugs.R) using a prior for tau from the previous part, and a prior for mu that is the most appropriate from (i)–(iv) above. Include the posterior histogram of mu with your analysis. Does the conclusion differ noticeably from the diuse priors?

Exercise 15.5. [Purpose: With repeated measures, the group estimate and individual estimates reflect different sources of variation in the data.]

Consider the data, with different assignments to subjects as follows:

y values:	1	2	3	21	22	23	41	42	43
s, large between subj. var.:	1	1	1	2	2	2	3	3	3
s, small between subj. var.:	1	2	3	1	2	3	1	2	3

The first row shows nine data values. The second row indicates a situation in which the first three data points come from subject 1, the second three come from subject 2, and so on. The third row indicates a different situation, in which the first, fourth, and seventh data points come from subject 1, as so forth. The assignment of data to subjects in the second row produces large between-subject variance, while the assignment of data to subjects in the third row produces small between-subject variance. Notice that for both situations, however, the overall mean is the same because the data are the same. We are interested in which situation gives a more precise estimate of the group-level mean.

Modify the program of Section 15.4.2 (stemsBrugs.R) for use with these data. (No need to be too fancy. For example, in the list , just type in `y = c(1, 2, 3, 21, 22, 23, 41, 42, 43)` and `subj = c(1,1,1, 2,2,2, 3,3,3)` .) Be sure to make the prior appropriate for these data. Run the program twice, once for each assignment of data to subjects (using the same prior for both runs). Hint: Your results should look something like those in Figure 15.11. Notice that the mean of the posterior is essentially the same for both runs, but the HDI widths are quite different. Explain why.

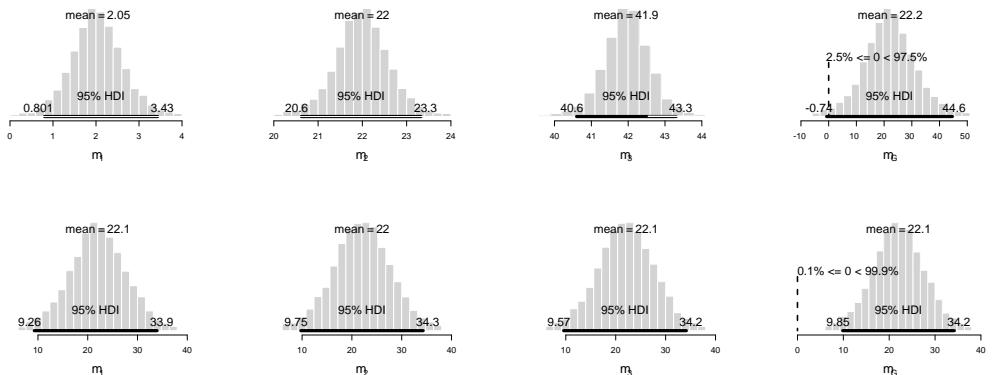


Figure 15.11: For Exercise 15.5, involving the same data arranged across subjects. Top row: Large between-subject variance. Bottom: Small between-subject variance. Notice the width of the HDI on the group parameter, G .

Chapter 16

Metric Predicted Variable with One Metric Predictor

Contents

16.1 Simple linear regression	344
16.1.1 The hierarchical model and BUGS code	46
16.1.1.1 Standardizing the data for MCMC sampling	347
16.1.1.2 Initializing the chains	348
16.1.2 The posterior: How big is the slope?	349
16.1.3 Posterior prediction	355
16.2 Outliers and robust regression	352
16.3 Simple linear regression with repeated measures	354
16.4 Summary	357
16.5 R code	358
16.5.1 Data generator for height and weight	358
16.5.2 BRugs: Robust linear regression	359
16.5.3 BRugs: Simple linear regression with repeated measures	362
16.6 Exercises	366

The agri-bank's threatnin' to revoke my lease
If my eld's production don't rapid increase.
Oh Lord how I wish I could divine the trend,
Will my furrows deepen? and Will my line end?

In this chapter we consider situations such as predicting `graph`'s weight from their height, or predicting their blood pressure from their `weight`, or predicting their income from years of education. In these situations, the predicted variable is metric and the single predictor is also metric. We will describe the relationship between the predicted variable, and predictor `x`, with a simple linear model and normally distributed residual randomness in `y`. In terms of the generalized linear model (GLM), the model assumes that $y \sim N(\mu, \sigma^2)$ with $\mu = \beta_0 + \beta_1 x$. This model appears in Table 14.1 (p. 312) in the first row and second column. This model is often referred to as "simple linear regression".

Kruschke, J. K. (2010) Doing Bayesian Data Analysis: A Tutorial with R and BUGS. Academic Press Elsevier. Copyright 2010 by John K. Kruschke. Draft of May 11, 2010. Please do not distribute this preliminary draft. If you report Bayesian analyses based on this book, please do cite it! —

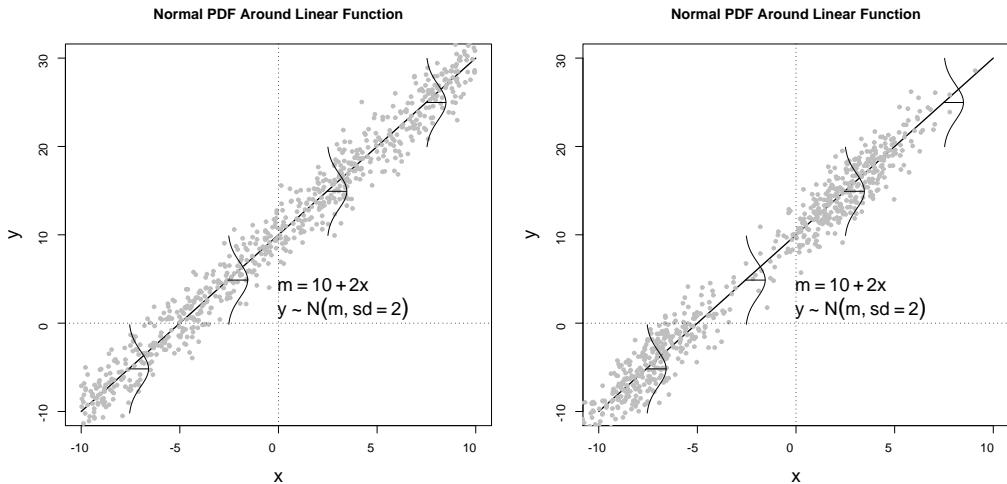


Figure 16.1: Examples of points normally distributed around a linear function. (The left panel repeats Figure 14.9, p. 309.) The model ~~assumes~~ that the data are normally distributed vertically around the line, as shown. Moreover, the variance of y is the same at all values of x . The model puts no constraints on the distribution of x . The right panel shows a case in which x are distributed bimodally, whereas in the left panel they are distributed uniformly. In both panels, there is homogeneity of variance.

16.1 Simple linear regression

Figure 16.1 shows an example of simulated data generated by the assumed model, with parameters displayed in the Figure. At any value of x , the mean predicted value is $= 0 + 2x$, and the data values are normally distributed around that mean. For a review of how to interpret the slope, 2 , and intercept, 0 , see Figure 14.1, p. 295.

Note that the model only specifies the dependency on x . The model does not say anything about what generates x and there is no probability distribution assumed for describing x . The x values in the left panel of Figure 16.1 were sampled randomly from a uniform distribution, merely for purposes of illustration; whereas the x values in the right panel of Figure 16.1 were sampled randomly from a bimodal distribution. Both panels show data from the same model of the dependency on x .

It is important to emphasize that the model assumes homogeneity of variance. At every value of x , the variance of y is the same! This homogeneity of variance is easy to see in the left panel of Figure 16.1 because the x values are uniformly distributed: The smattering of data points in the vertical direction appears visually to be the same at all values of x . Homogeneity of variance is less easy to identify visually when the x values are not uniformly distributed. For example, the right panel of Figure 16.1 displays data that may appear to violate homogeneity of variance, because the apparent spread of the data seems to be larger for $x = 2:5$ than for $x = 7:5$. Despite this deceiving appearance, the data do respect homogeneity of variance. The reason for the appearance is that for regions in which x is sparse, there is far less opportunity for the sample values to come from the tails of the normal distribution. In regions where x is dense, there are many opportunities for y to come from the tails.

In applications, the x and y values are provided by some real-world process. In the

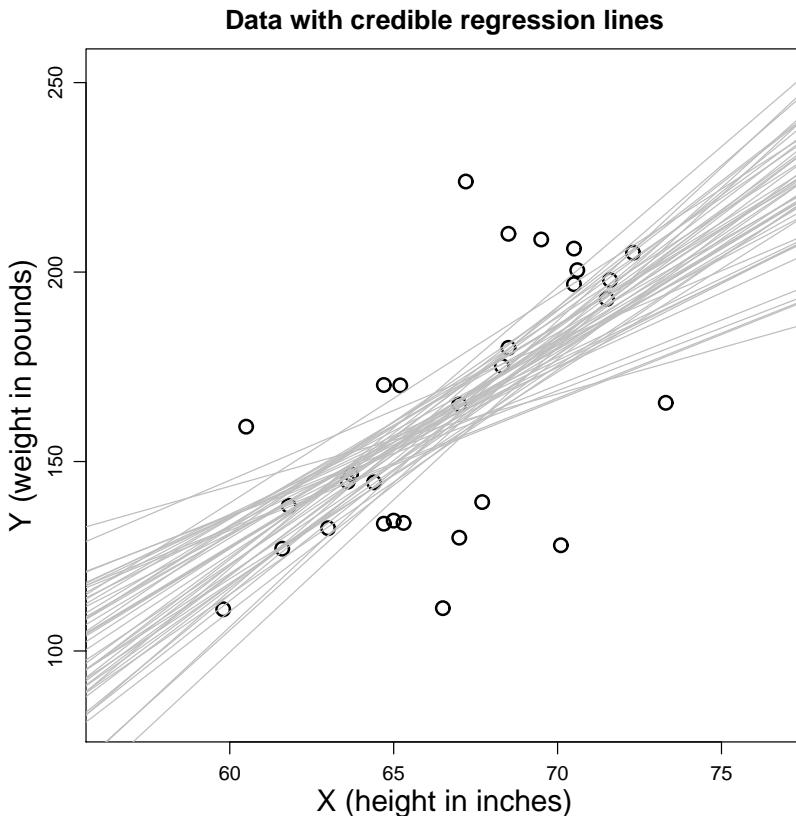


Figure 16.2: Data points, with a smattering of believable regression lines superimposed.

real-world process, there might or might not be any direct causal connection between x and y . It might be that x causes y , or y causes x , or some third factor causes both x and y have no causal connection, or some combination of any of those! The simple linear model makes no claims about causal connections between x and y . The simple linear model merely describes a tendency for x values to be linearly related to y values, hence “predictable” from the x values.

As an example, suppose we have measurements of height, x , and weight, in pounds, for some randomly selected people. Figure 16.2 shows height and weight values of the people. The data were generated from the `proj1` in Section 16.5.1 (`HtWtDataGenerator.R`), which uses realistic population parameters. The data in Figure 16.2 do appear to indicate that as height increases, weight tends to increase. This covariation between height and weight does not imply that one attribute causes the other. When an adult eats a lot of sugary foods, or goes on a diet, changing his or her weight, his or her height does not change. Despite the lack of causal relationship, the two values do covary, and one can be (imperfectly) predicted from the other.

Our goal is to determine what regression lines are most ~~believable~~^{likely}, given the data. In other words, we want to infer what combinations of α_0 and α_1 are most believable, given the data. We use Bayes' rule:

$$p(\alpha_0; \alpha_1; | y) = p(y | \alpha_0; \alpha_1;) p(\alpha_0; \alpha_1;) \quad \text{and} \quad p(y | \alpha_0; \alpha_1;) p(\alpha_0; \alpha_1;)$$

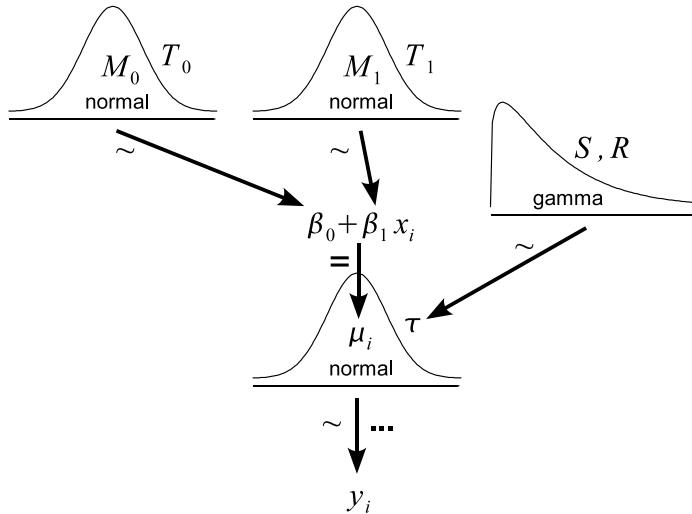


Figure 16.3: A model of dependencies for simple linear regression. The parameter β_1 is of primary interest; it describes the slope of the line relating x to y .

Analytical forms for the posterior can be obtained for appropriate priors. For example, analogous to the derivation in the previous chapter, if the precision is fixed, i.e., if the prior on τ is a spike over a certain value, then normal priors on the slope and intercept yield normal posteriors (e.g., Bolstad, 2007). Fortunately, we do not have to worry much about mathematical derivations because we can let BUGS approximate the posterior. Our job, therefore, is to specify sensible priors and to make sure that BUGS generates a trustworthy posterior sample that is converged and well-mixed.

16.1.1 The hierarchical model and BUGS code

There are three parameters in the model, and therefore we need to establish a prior distribution on the joint three-dimensional parameter space. To keep things simple, we will assume that the three parameters are independent in our prior belief. They would not need to be; for example, we might have a prior belief that the slope and intercept are (anti-)correlated, or we might have a prior belief that the precision and the slope are correlated. For now, however, we assume independence in the prior, and therefore we can express the prior as three marginal distributions on the separate parameters.

Figure 16.3 shows the hierarchy of dependencies that we assume. At the bottom of the figure, we see that the data depend on the parameters μ_i and τ , which describe the mean and precision of a normal distribution. In other words, data are modeled by a normal likelihood function, as was emphasized in the previous chapter. The value of the precision, τ , depends on our prior belief distribution, shown in the upper right of the figure as a gamma distribution with two shape parameters. The shape parameters are constants that express what we think the precision is likely to be and how certain we are in that prior belief.

The mean, μ_i , of the likelihood function is a linear function of x_i as shown in the middle of Figure 16.3. The linear function has two parameters, the intercept β_0 and the slope β_1 . The slope parameter is what we are usually most interested in because it describes the relation between x and y . The values of the slope and intercept depend on the priors shown in the upper layer of the figure. Both the slope and the intercept here have prior

beliefs modeled as normal distributions.

It is useful to compare Figure 16.3 with Figure 15.2, p. 320. The underlying approach is the same for both scenarios, in that there is a normal likelihood for the data, and we set priors on the mean and the precision. The only difference is that Figure 16.3 has one additional prior for the β_1 parameter. Otherwise, the models are identical.

Every arrow in Figure 16.3 has a corresponding line in the BUGS model specification: (SimpleLinearRegressionBrugs.R)

```

9  model {
10    for( i in 1 : Ndata ) {
11      y[i] ~ dnorm( mu[i] , tau )
12      mu[i] <- beta0 + beta1 * x[i]
13    }
14    beta0 ~ dnorm( 0 , 1.0E-12 )
15    beta1 ~ dnorm( 0 , 1.0E-12 )
16    tau ~ dgamma( 0.001 , 0.001 )
17 }
```

Notice that the constants in the hyperpriors are generic for a diffuse prior. In real applications, you would probably want to use better informed constants.

16.1.1.1 Standardizing the data for MCMC sampling

In principle, we could run the BRugs code on the raw x, y data. In practice, however, the attempt often fails. There's nothing wrong with the mathematics or logic, the problem is that believable values of the slope and intercept parameters tend to be tightly correlated, and this narrow diagonal zone of believability is difficult for sampling algorithms to explore.

The right panel of Figure 16.4 shows an example of this type of correlation in believable values. The believable slopes and intercepts are extremely correlated. Sampling from such a tightly correlated distribution is typically very difficult to do directly. It is difficult to discover a point in the narrow zone in the first place. Therefore, if a viable point is discovered, the chain does not move efficiently. Gibbs sampling gets stuck because it keeps "bugged" into the walls" (recall discussion of Section 8.4.2.1, p91). Metropolis algorithms often are not clever enough to automatically tune a proposal distribution to match a diagonal posterior.

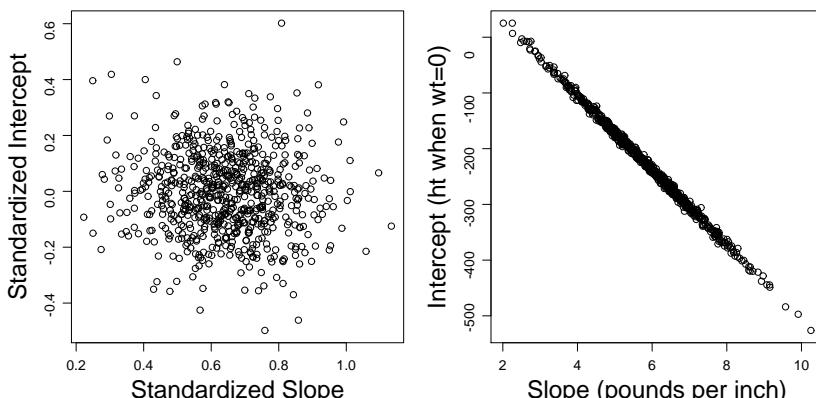


Figure 16.4: Slope and intercept values of believable regression lines, in standardized and raw scales.

The reason for the correlation of slope and intercept is easily examining Figure 16.2. There you can see various believable lines that ~~through~~ fit the scatter of points. Notice that if a line has a steep slope, its intercept must be ~~be~~ low but if a line has a smaller slope, its intercept must be higher. Thus, there is a trade-off between slope and intercept for the believable lines.

One of the main tricks used for successful execution of the MCMC sampling is standardizing the data. Standardizing simply means re-scaling data relative to their mean and standard deviation:

$$z(x) = \frac{(x - M_x)}{SD_x} \quad \text{and} \quad z(y) = \frac{(y - M_y)}{SD_y} \quad (16.1)$$

where M_x is the mean of the data values and SD_x is the standard deviation of the data values. (Do not confuse M_x and M_y with the constants used in the specification of the priors in the hierarchical diagram.) It is easy to prove, using simple algebra, that the mean of the resulting $z(x)$ values is zero, and the standard deviation of the resulting $z(x)$ values is one, for any initial data set.

Having used BUGS to find slope and intercept values for the ~~standardized~~ data, we then need to convert the parameter values back to the original ~~scales~~. Denote the intercept and slope for standardized data as $\hat{\alpha}$ and $\hat{\beta}$ (Greek letter “zeta”), and denote the predicted value of y as \hat{y} . Then:

$$\frac{\hat{y} - M_y}{SD_y} = \frac{(x - M_x)}{SD_x} \quad \text{from Eqn. 16.1}$$

$$\hat{y} = M_y + SD_y \left[\frac{(x - M_x)}{SD_x} \right] \quad (16.2)$$

Thus, for every believable combination of α_1 values, there is a corresponding believable combination of α_0 ; α_1 values specified by Equation 16.2.

The sampled points in Figure 16.4 were generated by running `ols` on the standardized data, then transforming the results according to [Equation 16.2](#). Notice that the standardized slope and intercept, in the left panel of Figure 16.5, show no noticeable correlation. This is because the believable intercepts tend to be near zero and that's because the standardized data have means of zero. Thus, even when the slopes of believable lines differ, the intercepts still hover around zero.

16.1.1.2 Initializing the chains

Figure 16.4 showed an example of how the believable values posterior distribution can occupy a fairly narrow region of parameter space. For MCMC chain to randomly sample from the posterior, the random walk must first get into the modal region of the posterior in the first place. We might simply start the chain at any point in parameter space, randomly selected from the prior distribution, and wait through the burn-in period until the chain randomly wanders into the bulk of the posterior. Unfortunately, for many real-world situations, this burn-in period can be a very long time. Therefore it helps to initialize the chains near the bulk of the posterior if we can.

If the data set is large and dominates the prior, or if the prior is informed from previous results and is reasonably consistent with the new data, then the peak of the likelihood function will be reasonably near the peak of the posterior distribution. Therefore, a reasonable candidate for the initial value of the chain is the maximum-likelihood estimate of the parameters. This heuristic is useful if we have a simple way of determining the maximum-likelihood estimate. Fortunately, for simple linear regression, we do. When the data are standardized, the maximum-likelihood estimate (MLE) of β_0 is zero, the MLE of β_1 is the correlation (denoted r) of x and y , and the MLE of the precision, τ^2 , is $1 - r^2$. To get an intuition for the statement about precision, consider what happens when the correlation approaches 1 (its maximum possible value). As r approaches 1, the $x; y$ data fall very close to a straight line, which implies that the deviation of the data away from the line is very small (recall Figure 16.1). Hence σ^2 is large, τ^2 is small, and hence τ is large, as reflected by the formula $\tau = \sqrt{1 - r^2}$.

16.1.2 The posterior: How big is the slope?

Figure 16.5 shows the posterior distribution of slope values. The standardized and original-scale slopes indicate the same relationship on different scales, and therefore the posterior distributions are identical except for a change of scale. The posterior distribution tells us exactly what we want to know: The believable slopes. We see that weight increases by about 5 or 6 pounds for every 1-inch increase in height. The HDI provides a useful summary of the range of believable slopes.

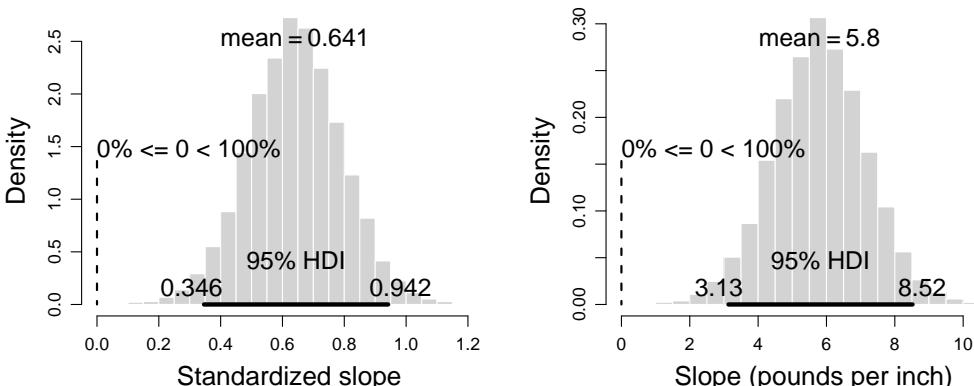


Figure 16.5: Posterior distribution of slopes.

If we were interested in determining whether the predicted value on the predicted variable, we might use the decision rules discussed in Section 12.1.3 (p. 244). We may want to establish a ROPE around zero for the predictor, then check whether the entire 95% HDI excludes the ROPE. (Usually, if the true value is zero, the HDI will overlap the ROPE, thereby reducing false alarms.)

It would be possible to “test” whether the slope is non-zero by doing a Bayesian comparison of two models: One model would have an arbitrarily flat prior on the slope parameter; the other model would have an arbitrarily peaked prior. It would be straightforward, in principle, to set up the model comparison using multidimensional MCMC as described in Section according to the method explained in Section 10.2 (p. 197). But, as argued in Section 12.2.2 (p. 249) and elsewhere, all that model comparison tells us

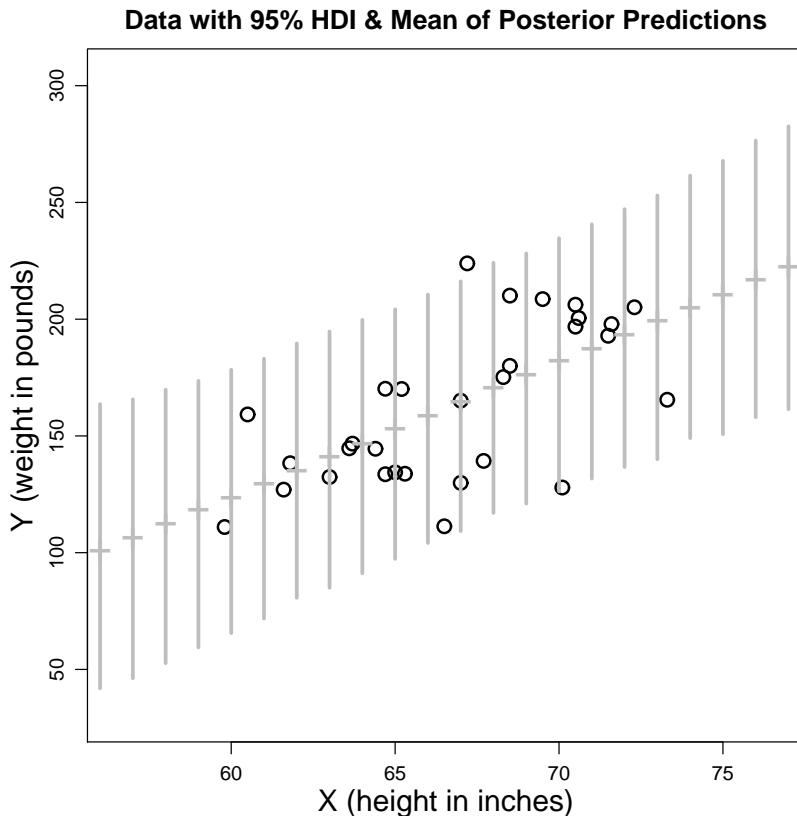


Figure 16.6: Data points, with 95% HDIs and means of posterior predictions in y, at selected values. The HDIs are wider (i.e., taller) for extrapolated values than for x values in the middle of the observed data.

is which of two unbelievable models is less unbelievable. In the present application, the model comparison would favor the model that allows non-zero slopes. But that is not all we want to know, because usually we also want to know what values are credible. That information is provided by the posterior distribution of the parameter values, often starting with an (mildly) informed prior distribution instead of an arbitrarily diffuse prior of mathematical convenience. Thus, in the present application there is little to be gained by doing a Bayesian model comparison of null prior against arbitrary diffuse prior.

16.1.3 Posterior prediction

Linear regression is often used for predicting values from x values. For example, we can use the results of the regression to answer the question: if a person were 75 inches tall, what is the probable weight of the person? The Bayesian answer is the probability of every possible weight, given the height and the previous data $p(y|x; D)$. This distribution of values has uncertainty stemming from the inherent noise and also from uncertainty in the estimated values of the regression coefficients.

A simple way to get a good approximation $p(y|x; D)$ is by generating random values of y for every step in the MCMC sample of credible parameter values. At any step in the chain, there are particular values of α_0 and α_1 , which we use to generate representative predicted values of y according to $N(\mu = \alpha_0 + \alpha_1 x; \sigma^2 = \tau^2)$. We do that at all

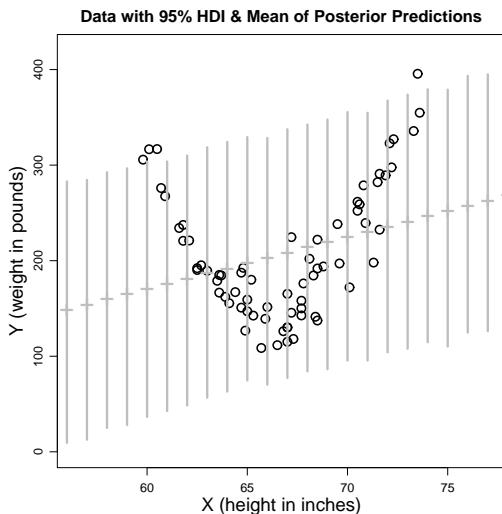


Figure 16.7: Data points, with 95% HDIs and means of posterior predictions at selected x values. The data show systematic discrepancy from the linear-model predictions.

the steps in the chain, and thereby create a large set of predicted values. From this set, we can compute the mean predicted value, the HDI of the predicted values, etc.

A summary of the posterior predictions is displayed in Figure 16.6. Each vertical grey segment, over a particular x value, indicates the extent of the 95% HDI for the distribution of randomly generated values at that x value. The dash across the middle of the grey segment indicates the mean of the posterior predicted values.

Consider the limits of the HDIs at each value of x . Notice that the length of the HDIs is larger for high and low x values than for middling x values. To verify this fact, hold a ruler against the HDIs for small, middling, and large values. This variation in the length of the HDIs makes sense intuitively: As the x value goes farther from the observed data, the predicted y value should become less certain. Another way of understanding the variation in HDI lengths is by considering the credible regression lines in Figure 16.2. All of those differently sloped lines go near the bulk of the data in the middle; the different slopes extend higher or lower at extrapolative values. We average across all those lines to generate predicted y values. Therefore the predicted y values will be more widespread at extrapolative values of x .

If the model of the data is a good model of the data, then the parameter values for the model ought to generate simulated data that “look like” the real data. What it means for simulated data to “look like” the real data is defined by the analyst, and may be anything that is useful for the application at hand. Loosely speaking, a model is a good model if (1) the data fall mostly within the predicted zone, (2) the data are distributed approximately the way the model says they should be, e.g., normally, and (3) discrepancies of data from predictions are random and not “systematic”.

Posterior prediction is also useful for alerting us to non-linear trends in data. Figure 16.7 shows a set of data for which the y values have a clear quadratic (i.e., parabolic) relationship with x , in addition to a linear trend. These data can be entered into a linear regression program, and the resulting posterior prediction HDIs are shown as the vertical grey segments in Figure 16.7. We see that the data do fall mostly within the range of the posterior predictions, but the distribution of the data is systematically discrepant from the linear spine of the model: At high and low values of x , the data fall well above the linear spine, but at middling values of x , the data fall well below the linear spine. This sort of systematic discrepancy from the posterior predictions is a clue that the model could be improved.

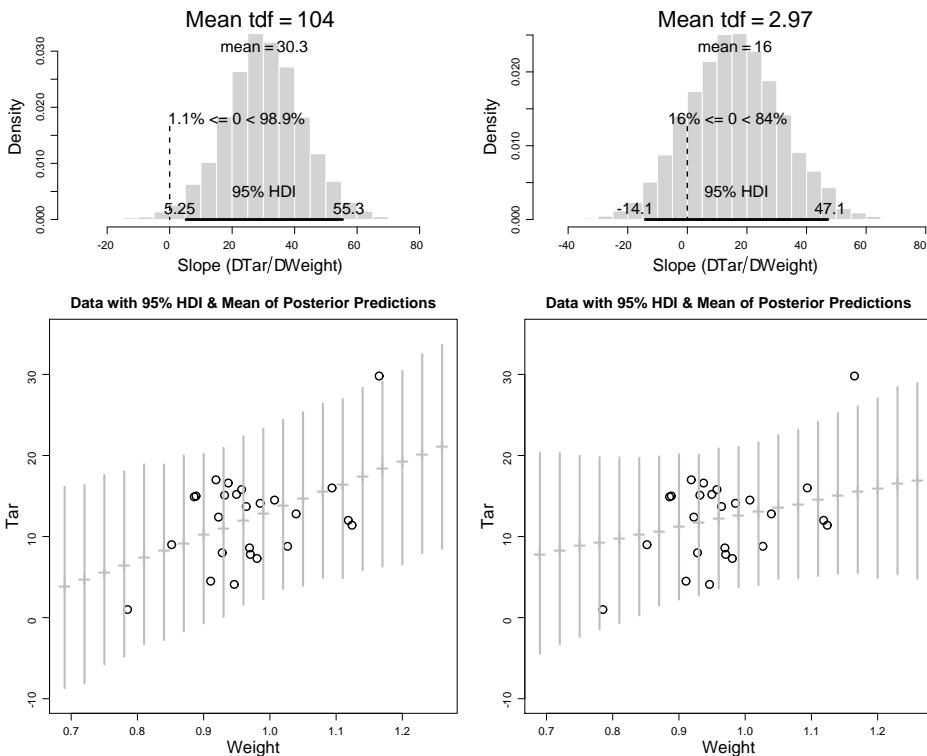


Figure 16.8: Data fit with different priors on the df parameter of the likelihood function. Left column shows results from prior bias on df values, such that a virtually normal likelihood function is assumed. Right column shows “robust” regression resulting from a strong prior bias on df values.

16.2 Outliers and robust regression

Recall from Figure 15.5, p. 325, that the estimated parameters for a normal likelihood function can be greatly distorted by outliers in the data. The sensitivity also occurs when the normal distribution is used in linear regression. The normal likelihood function demands that the regression line is vertically close to all the data points, because the likelihood value is tiny for points more than about three standard deviations from the line. Consequently, outlying data points can have disproportionate leverage on the estimate of the regression coefficients.

As an example, consider some data regarding 25 brands of cigarettes (McIntyre, 1994). For each brand, a cigarette was assessed for its weight, nicotine, amount of tar, and amount of carbon monoxide produced when burned. There are many relationships among these variables might be investigated, but consider predicting the tar content from the weight of the cigarette. These data are shown in the lower panels of Figure 16.8. The points are rather diversely scattered, and do not show an overwhelmingly strong covariation. There is, however, an outlying point at the heaviest weight and highest tar content (in the upper right of the scatterplot). This outlying point can have a disproportionate influence on the estimated slope in the regression model.

The left column of Figure 16.8 shows the results when a normal likelihood function is used (actually, when the likelihood function is a distribution with large df). The histogram in the upper panel shows that the slope is credibly larger than zero. The plot of posterior

predictive HDIs in the lower panel shows that the estimated regression line has been tilted up toward the outlier. The attraction to the outlier is caused by the small tails of the normal likelihood function: For the outlying datum to be “under the umbrella” of the normal likelihood, the line must get fairly close.

The right column of Figure 16.8 shows the results when a t distribution is used as the likelihood function, with a prior on the df parameter that strongly favors small values. This prior implies that the likelihood function has the tall tails of a low-df t distribution, unless the data strongly suggest otherwise. As can be seen in the figure, the 95% HDI of the posterior on the slope contains zero. The lower right panel shows that the estimated regression line is not slanted as far toward the outlier, because the tail of the distribution gives the outlier some modest non-zero probability despite being far from the mean of the distribution. The posterior slope estimate here better reflects the bulk of the points in the scatterplot, and is not so strongly dominated by the outlier.

The BUGS model for robust regression is a simple extension of the model for the normal likelihood. We simply replace the normal distribution with a t distribution, and include the necessary prior specification. Whereas the normal distribution has two parameters, namely the mean and the precision, the t distribution has those two plus a third parameter, namely the degrees of freedom which can take on real values of 1 or greater (for a reminder of how the df parameter works, see Figure 15.4, p. 324). One way to put a prior on the df parameter is by sampling a value, denoted udf , from a uniform distribution, and then transforming that value into the range allowed for df , which extends from 1 to infinity. One such transformation appears on line 18 of this model specification: (SimpleRobustLinearRegressionBrugs.R)

```

9 model {
10   for( i in 1 : Ndata ) {
11     y[i] ~ dt( mu[i] , tau , tdf )
12     mu[i] <- beta0 + beta1 * x[i]
13   }
14   beta0 ~ dnorm( 0 , 1.0E-12 )
15   beta1 ~ dnorm( 0 , 1.0E-12 )
16   tau ~ dgamma( 0.001 , 0.001 )
17   udf ~ dunif(0,1)
18   tdf <- 1 - tdfGain * log(1-udf) # tdf in [1,Inf).
19   # tdfGain specified in data section
20 }
```

Notice that tdf is used as the df parameter in the likelihood function specified on line 11. The value of tdf is determined from udf by the transformation $\text{tdf} \leftarrow 1 - \text{tdfGain} * \log(1 - \text{udf})$, where tdfGain is a constant that expresses the prior belief in large values of df . Figure 16.9 shows a graph of the transformation. You can see the graph that when tdfGain is small, the df values are close to 1 across almost the entire range of udf . But when tdfGain is large, then the df values are large across most of the range of udf .

The prior on the df parameter expresses how much we believe there are outliers in the data. For example, the left column of Figure 16.8 was created by setting $\text{tdfGain}=100$, and the right column of Figure 16.8 was created by setting $\text{tdfGain}=1$. The resulting posteriors are noticeably different.

So, you may ask, which conclusion from Figure 16.8 is correct? Does the tar content of cigarettes increase by about 30 units for every unit increase in weight, as indicated by the nearly-normal likelihood, or does the tar content increase by about only about half that much, as indicated by the long-tailed t distribution? The answer is, in this case, that it depends

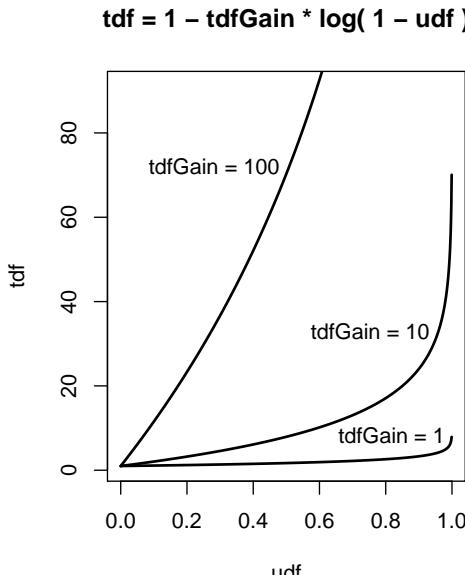


Figure 16.9: Transformation from udf to tdf , used to set the prior of tdf parameter in robust regression.

on your prior beliefs. If you have some very good reason to believe that tar content should be normally distributed (with the same variance) at every level of weight, then the left column of Figure 16.8 describes your posterior beliefs. But the assumption of normality virtually disavows the possibility of outliers in the data, and therefore seems untenable, at least for this application. Therefore the prior should allow low values of tdf . The dependency of the posterior on the prior suggests that either more specific knowledge should be brought to bear, or additional predictors should be included if available, or more data need to be collected.

16.3 Simple linear regression with repeated measures

Suppose that for every individual, we have multiple observations (x_{ij}, y_{ij}) pairs. With these data, we can estimate a linear model for every individual. We also assume that the individuals are mutually representative of a common group. Then we can use the estimates from the individuals to inform estimates of group-level parameters.

One example of this scenario comes from Walker, Gustafson, Farmer (2007), who measured reading ability scores of children across several years. Thus each child contributed several age and reading-score pairs. A regression model can describe each child's reading ability through time, and higher-level distributions describe the distribution of intercepts and slopes across individuals. By virtue of being linked indirectly through the higher level distribution, estimates of individuals are mutually informed by other individuals.

As another concrete example of this situation, consider an experiment in which the investigators were interested in how quickly different organs clear themselves of contaminants (Feldman, 1988, who reported data from an unpublished experiment by S. B. Weinstock and J. D. Brain). The researchers administered iron oxide particles to rats, because the iron oxide remaining in the body could be assayed noninvasively via magnetometry. Four rats were given intravenous injection of iron oxide, particles of which are taken up by liver endothelial cells. Four other rats were given the iron oxide by tracheal instillation, so that the particles were taken up by lung macrophages. All the researchers were

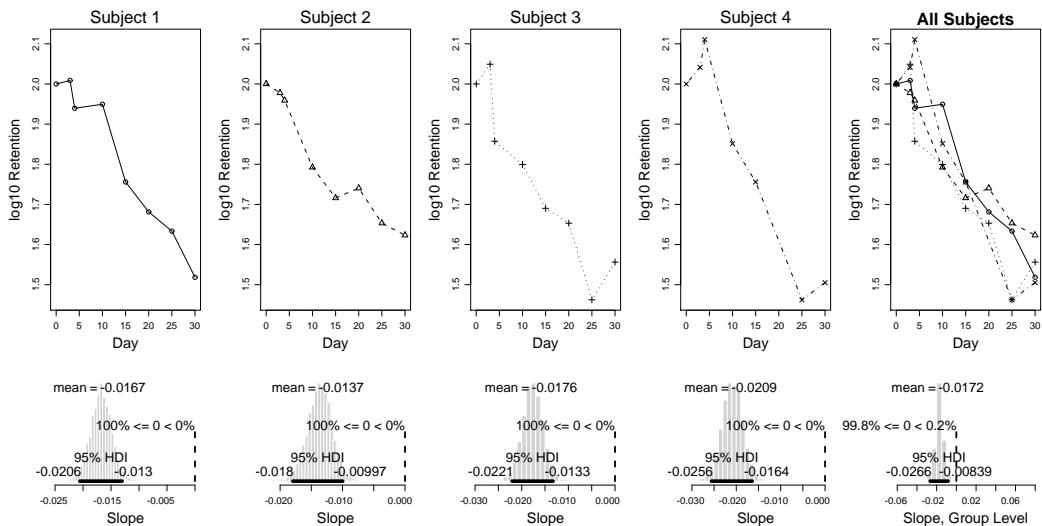


Figure 16.10: Upper row shows data from Feldman (1988). Lower row shows believable values of the slopes.

specifically interested in comparing the clearance rates of the two groups, we will consider only the “lung” group.

The amount of iron oxide retained in the body, as a percentage of the initially assayed amount, was measured at various times during the following days. The data for the lung group are shown in the top row of Figure 16.10. The retention amount is plotted on a logarithmic scale, so that the retention curves are approximately linear. Notice that all the curves start at a value of 2.0, because the first measurement establishes a baseline that defines 100%, and $\log_10(100) = 2.0$. Some curves rise above the initial value, which presumably does not indicate spontaneous creation of iron, but instead indicates measurement noise, either at the initial or subsequent times. The graphs indicate a reduction through time. The goal of our analysis is to determine what the believable reductions are, given the data. We may also want to know if the apparent reduction is believably non-zero.

To model this situation, each subject's data set is described by an individual linear regression, and the regression coefficients of the individuals are, in turn, modeled by group-level distributions. The group-level distributions are controlled by parameters that describe the central tendency and variability of the group, and these group-level parameters are the ones in which we are primarily interested.

Figure 16.11 shows the hierarchy of dependencies. At the individual level, we just see a simple linear regression for each individual, with the same components as Figure 16.3. The regression coefficients for the j^{th} individual, namely α_j , β_j , and σ_j , in turn come from distributions that describe the group. For example, the individual slope coefficients, β_j , are assumed to come from a normal distribution, with mean β_G and precision $\tau_{\beta G}$. We are interested in estimating those group-level parameters, therefore each is given a prior distribution at the top level of the diagram.

Every arrow in Figure 16.11 has a corresponding line in the Stan model specification: (SimpleLinearRegressionRepeatedBrugs.R)

```
8 model {
9   for( r in 1 : Ndata ) {
```

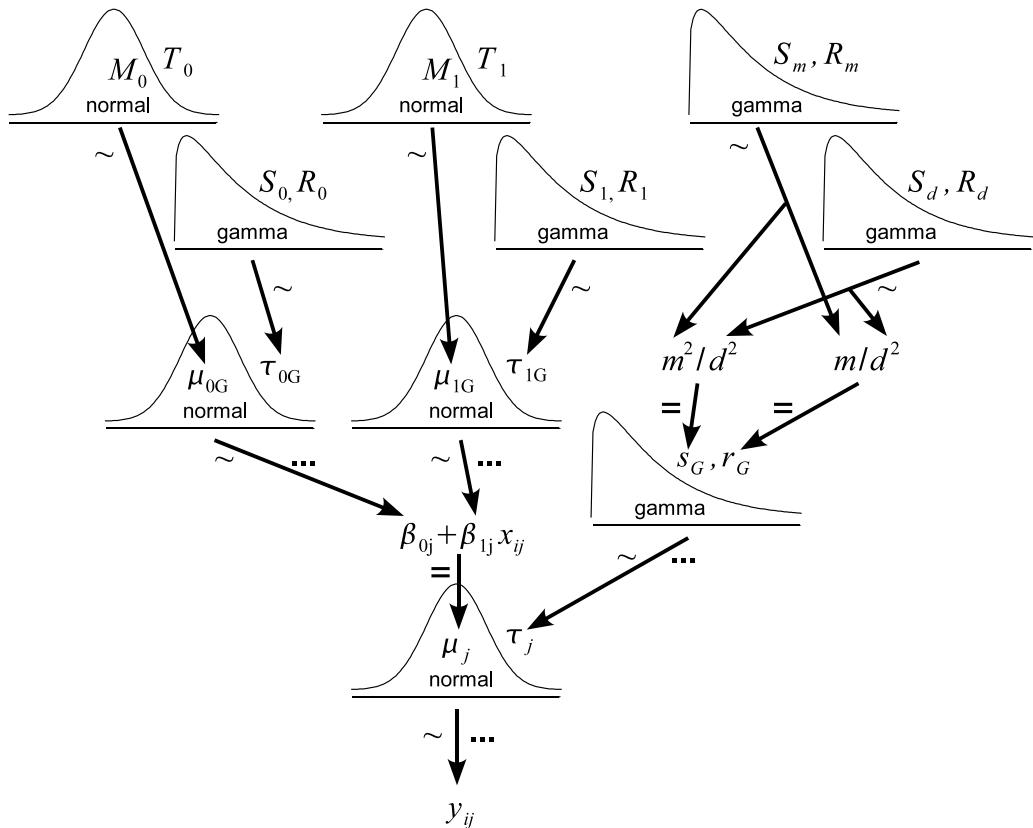


Figure 16.11: A model of dependencies for repeated scores y_{ij} of N subjects drawn independently from the same group. The slope for j^{th} subject is β_{1j} . Across subjects, it is distributed normally, with mean μ_{1G} and precision τ_{1G} , i.e., across subjects the variance of the slopes is τ_{1G} .

```

10      y[r] ~ dnorm( mu[r] , tau[ subj[r] ] )
11      mu[r] <- b0[ subj[r] ] + b1[ subj[r] ] * x[r]
12  }
13  for ( s in 1 : Nsubj ) {
14    b0[s] ~ dnorm( mu0G , tau0G )
15    b1[s] ~ dnorm( mu1G , tau1G )
16    tau[s] ~ dgamma( sG , rG )
17  }
18  mu0G ~ dnorm(0,.01)
19  tau0G ~ dgamma(.1,.1)
20  mu1G ~ dnorm(0,.01)
21  tau1G ~ dgamma(.1,.1)
22  sG <- pow(m,2)/pow(d,2)
23  rG <- m/pow(d,2)
24  m ~ dgamma(1,.1)
25  d ~ dgamma(1,.1)
26 }
```

The constants in the top-level priors are generically vague because of laziness during programming. More thoughtful priors may be desirable.

The results of running the program are shown in the lower rbf^{Figure 16.10}. The results match closely those reported in the non-Bayesian analysis of Feldman (1988), who reported 0.0173 for the group slope. The results also show that every individual slope is

believably different from zero, as is the estimate of the group average slope.

The four rats in this experiment had remarkably similar clearance rates. In many realms of research, there is much more variation between subjects. Suppose that different rats had very different clearance rates, with some faster and some slower than those reported here, such that the group-average clearance rate was ~~the same~~. What aspects of the posterior would be most affected by the increase in variability of individual slopes? Would the certainty of individual slope estimates be much affected? Would the certainty of the group-average slope be much affected? Exercise 16.1 has you systematically alter the data, to answer these questions. It is also a good intuition-stirring exercise to consider the influence of differing individual intercepts. In particular, without changing the individual slopes, but merely by changing the individual intercepts, estimate of the group-average slope can be affected. See Exercise 16.2.

The data in this example were very well behaved, insofar as there were no blatantly outlying data points. Specifically, within each individual all the values descended in a reasonably linear trend. And, across subjects, all the individual slopes were reasonably similar to each other, with no subjects who had clearances much faster or much slower than the others. With a larger sample, however, outlying subjects might be encountered. And there is always the chance that device malfunctions or transcription errors could be inadvertently introduced into individual data points. Whether there is reason to believe that outliers may be lurking in the data distributions can be substituted for the normal distributions at the appropriate level. Exercise 16.3 shows example of outlying slopes and intercepts, involving income and family size.

16.4 Summary

This chapter explained Bayesian simple linear regression, where “simple” means it involves a single predictor. There were several main points:

Prediction of y from x does not imply or assume any particular causal or temporal relation between y and x .

The usual model for linear regression assumes homoscedasticity: The variance of y is the same at all values of x .

It is typical, but not necessary, to put normal priors on the intercept and slope parameters in the linear regression model.

To make MCMC sampling efficient, it is helpful to standardize the y and x data, and to initialize the chains at the MLE values.

Posterior predictions are useful for predicting y values from x values, especially when extrapolating or interpolating to novel values. Posterior predictions are also useful for checking whether there are systematic discrepancies between data from the predictions of the model.

It is straightforward to extend the model to situations with repeated measures, in which every individual has data estimated by simple linear regression, and overarching parameters describe the distribution of slope, intercept, and precision parameters of the group.

When there are outliers in repeated measures from an individual in the distribution of parameters across individuals, a distribution may be used instead of a normal distribution.

In the next chapter, we introduce additional predictors to the model, and explore the concept of interaction across predictors.

16.5 R code

16.5.1 Data generator for height and weight

(HtWtDataGenerator.R)

```

1 HtWtDataGenerator = function( nSubj , rndsd=NULL ) {
2   # Random height, weight generator for males and females. Use  s parameters from
3   # Brainard, J. & Burmester, D. E. (1992). Bivariate distributions for height and
4   # weight of men and women in the United States. Risk Analysis, 12(2), 267-275.
5   # John K. Kruschke, January 2008.
6
7   require(MASS)
8
9   # Specify parameters of multivariate normal (MVN) distributions.
10  # Men:
11  HtMmu = 69.18
12  HtMsd = 2.87
13  InWtMmu = 5.14
14  InWtMsd = 0.17
15  Mrho = 0.42
16  Mmean = c( HtMmu , InWtMmu )
17  Msigma = matrix( c( HtMsd^2 , Mrho * HtMsd * InWtMsd ,
18                   Mrho * HtMsd * InWtMsd , InWtMsd^2 ) , nrow=2 )
19  # Women cluster 1:
20  HtFmu1 = 63.11
21  HtFsd1 = 2.76
22  InWtFmu1 = 5.06
23  InWtFsd1 = 0.24
24  Frho1 = 0.41
25  prop1 = 0.46
26  Fmean1 = c( HtFmu1 , InWtFmu1 )
27  Fsigma1 = matrix( c( HtFsd1^2 , Frho1 * HtFsd1 * InWtFsd1 ,
28                   Frho1 * HtFsd1 * InWtFsd1 , InWtFsd1^2 ) , nrow=2 )
29  # Women cluster 2:
30  HtFmu2 = 64.36
31  HtFsd2 = 2.49
32  InWtFmu2 = 4.86
33  InWtFsd2 = 0.14
34  Frho2 = 0.44
35  prop2 = 1 - prop1
36  Fmean2 = c( HtFmu2 , InWtFmu2 )
37  Fsigma2 = matrix( c( HtFsd2^2 , Frho2 * HtFsd2 * InWtFsd2 ,
38                   Frho2 * HtFsd2 * InWtFsd2 , InWtFsd2^2 ) , nrow=2 )
39
40  # Randomly generate data values from those MVN distribution  s.
41  if ( !is.null( rndsd ) ) { set.seed( rndsd ) }
42  datamatrix = matrix( 0 , nrow=nSubj , ncol=3 )
43  colnames(datamatrix) = c( "male" , "height" , "weight" )
44  maleval = 1 ; femaleval = 0 # arbitrary coding values

```

```

45 for ( i in 1:nSubj ) {
46     # Flip coin to decide sex
47     sex = sample( c(maleval,femaleval) , size=1 , replace=TRUE      , prob=c(.5,.5) )
48     if ( sex == maleval ) {
49         datum = mvtnorm(n = 1, mu=Mmean, Sigma=Msigma ) }
50     if ( sex == femaleval ) {
51         Fclust = sample( c(1,2) , size=1 , replace=TRUE , prob=c(pro    p1,prop2) )
52         if ( Fclust == 1 ) {
53             datum = mvtnorm(n = 1, mu=Fmean1, Sigma=Fsigma1 ) }
54         if ( Fclust == 2 ) {
55             datum = mvtnorm(n = 1, mu=Fmean2, Sigma=Fsigma2 ) }
56     }
57     datamatrix[ i , ] = c( sex , round( c( datum[1] , exp( datum[2] ) )      ) , 1 ) )
58 }
59
60 return( datamatrix )
61 } # end function

```

16.5.2 BRugs: Robust linear regression

(SimpleRobustLinearRegressionBrugs.R)

```

1 graphics.off()
2 rm(list=ls(all=TRUE))
3 library(BRugs)          # Kruschke, J. K. (2010). Doing Bayesian dat      a analysis:
4                               # A Tutorial with R and BUGS. Academic Press / Elsevier.
5 #-----
6 # THE MODEL.
7 modelstring = "
8 # BUGS model specification begins here...
9 model {
10     for( i in 1 : Ndata ) {
11         y[i] ~ dt( mu[i] , tau , tdf )
12         mu[i] <- beta0 + beta1 * x[i]
13     }
14     beta0 ~ dnorm( 0 , 1.0E-12 )
15     beta1 ~ dnorm( 0 , 1.0E-12 )
16     tau ~ dgamma( 0.001 , 0.001 )
17     udf ~ dunif(0,1)
18     tdf <- 1 - tdfGain * log(1-udf) # tdf in [1,Inf).
19     # tdfGain specified in data section
20 }
21 # ... end BUGS model specification
22 " # close quote for modelstring
23 writeLines(modelstring,con="model.txt")
24 modelCheck( "model.txt" )
25
26 #-----
27 # THE DATA.
28
29 cigData = read.csv(file="McIntyre1994data.csv")
30 nSubj = NROW(cigData)
31 x = cigData[, "Wt"]
32 xName="Weight"
33 y = cigData[, "Tar"]
34 yName="Tar"
35
36 # Re-center data at mean, to reduce autocorrelation in MCMC s      ampling.
37 # Standardize (divide by SD) to make initialization easier.

```

```

38 xM = mean( x ) ; xSD = sd( x )
39 yM = mean( y ) ; ySD = sd( y )
40 zx = ( x - xM ) / xSD
41 zy = ( y - yM ) / ySD
42
43 # Specify data, as a list.
44 tdfGain = 1 # 1 for low-baised tdf, 100 for high-baised tdf
45 dataList = list(
46   x = zx ,
47   y = zy ,
48   Ndata = nSubj ,
49   tdfGain = tdfGain
50 )
51 # Get the data into BRugs:
52 modelData( bugsData( dataList ) )
53
54 #-----
55 # INTIALIZE THE CHAINS.
56
57 nchain = 3
58 modelCompile( numChains = nchain )
59
60 genInitList <- function() {
61   r = cor(x,y)
62   list(
63     beta0 = 0 ,      # because data are standardized
64     beta1 = r ,      # because data are standardized
65     tau = 1 / (1-r^2) , # because data are standardized
66     udf = 0.95 # tdf = 4
67   )
68 }
69 for ( chainIdx in 1 : nchain ) {
70   modellnits( bugslnits( genInitList ) )
71 }
72
73 #-----
74 # RUN THE CHAINS
75
76 # burn in
77 BurnInSteps = 100
78 modelUpdate( BurnInSteps )
79 # actual samples
80 samplesSet( c( "beta0" , "beta1" , "tau" , "tdf" ) )
81 stepsPerChain = ceiling(10000/nchain)
82 thinStep = 10
83 modelUpdate( stepsPerChain , thin=thinStep )
84
85 #-----
86 # EXAMINE THE RESULTS
87
88 source("plotChains.R")
89
90 fname = paste("SimpleRobustLinearRegressionBrugsTdfGa in",tdfGain,
91               sep="")
92 #beta0Sum = plotChains( "beta0" , saveplots=F , filenamero      ot=fname )
93 #beta1Sum = plotChains( "beta1"      , saveplots=F , filenamero      ot=fname )
94 #tauSum    = plotChains( "tau"       , saveplots=F , filenameroot=f      name )
95 #tdfSum    = plotChains( "tdf"       , saveplots=F , filenameroot=f      name )
96

```

```

97 # Extract chain values:
98 tdfSamp = samplesSample( "tdf" )
99 tdfM = mean( tdfSamp )
100 z0 = samplesSample( "beta0" )
101 z1 = samplesSample( "beta1" )
102 zTau = samplesSample( "tau" )
103 zSigma = 1 / sqrt( zTau ) # Convert precision to SD
104
105 # Convert to original scale:
106 b1 = z1 * ySD / xSD
107 b0 = ( z0 * ySD + yM - z1 * ySD * xM / xSD )
108 sigma = zSigma * ySD
109
110 # Posterior prediction:
111 # Specify x values for which predicted y's are needed:
112 xRang = max(x)-min(x)
113 yRang = max(y)-min(y)
114 limMult = 0.25
115 xLim= c( min(x)-limMult*xRang , max(x)+limMult*xRang )
116 yLim= c( min(y)-limMult*yRang , max(y)+limMult*yRang )
117 yLim = c(-10,35)
118 xPostPred = seq( xLim[1] , xLim[2] , length=20 )
119 # Define matrix for recording posterior predicted y values at each x value.
120 # One row per x value, with each row holding random predicted y values.
121 postSampSize = length(b1)
122 yPostPred = matrix( 0 , nrow=length(xPostPred) , ncol=postSampSize )
123 # Define matrix for recording HDI limits of posterior predicted y values:
124 yHDIlim = matrix( 0 , nrow=length(xPostPred) , ncol=2 )
125 # Generate posterior predicted y values.
126 # This gets only one y value, at each x, for each step in the chain.
127 for ( chainIdx in 1:postSampSize ) {
128   yPostPred[,chainIdx] = rnorm( length(xPostPred) ,
129                               mean = b0[chainIdx] + b1[chainIdx] * xPostPred ,
130                               sd = rep( sigma[chainIdx] , length(xPostPred) ) )
131 }
132 source("HDIofMCMC.R")
133 for ( xIdx in 1:length(xPostPred) ) {
134   yHDIlim[xIdx,] = HDIofMCMC( yPostPred[xIdx,] )
135 }
136
137 # Display believable beta0 and b1 values
138 windows()
139 par( mar=c(4,4,1,1)+0.1 , mgp=c(2.5,0.8,0) )
140 #layout( matrix(1:2,nrow=1) )
141 thinIdx = seq(1,length(b0),length=700)
142 #plot( z1[thinIdx] , z0[thinIdx] , cex.lab=1.75 ,
143 #       ylab="Standardized Intercept" , xlab="Standardized Slope" )
144 plot( b1[thinIdx] , b0[thinIdx] , cex.lab=1.75 ,
145       ylab="Intercept" , xlab="Slope" )
146 dev.copy2eps(file=paste(fname,"SlopeIntercept.eps",sep=""))
147
148 # Display the posterior of the b1:
149 source("plotPost.R")
150 windows(7,4)
151 par( mar=c(4,4,1,1)+0.1 , mgp=c(2.5,0.8,0) )
152 #layout( matrix(1:2,nrow=1) )
153 #histInfo = plotPost( z1 , xlab="Standardized slope" , compVal=0.0 ,
154 #                      breaks=30 )
155 histInfo = plotPost( b1 , main=bquote("Mean tdf"==.(signif(tdfM,3))) , cex.main=2 ,

```

```

156           xlab=bquote("Slope (" * Delta * .(yName) / Delta * .(xName)
157                           * ")") , compVal=0.0 , breaks=30 )
158 dev.copy2eps(file=paste(fname,"PostSlope.eps",sep="      "))
159
160 # Display data with believable regression lines and posteri      or predictions.
161 windows()
162 par( mar=c(3,3,2,1)+0.5 , mgp=c(2.1,0.8,0) )
163 # Plot data values:
164 plot( x , y , cex=1.5 , lwd=2 , col="black" , xlim=xLim , ylim=y      Lim ,
165       xlab=xName , ylab=yName , cex.lab=1.5 ,
166       main="Data with credible regression lines" , cex.main=1.3      3 )
167 # Superimpose a smattering of believable regression lines:
168 for ( i in seq(from=1,to=length(b0),length=50) ) {
169   abline( b0[i] , b1[i] , col="grey" )
170 }
171 dev.copy2eps(file=paste(fname,"DataLines.eps",sep="      "))
172
173 # Display data with HDIs of posterior predictions.
174 windows()
175 par( mar=c(3,3,2,1)+0.5 , mgp=c(2.1,0.8,0) )
176 # Plot data values:
177 #yLim= c( min(c(yHDIlim,y)) , max(c(yHDIlim,y)) )
178 plot( x , y , cex=1.5 , lwd=2 , col="black" , xlim=xLim , ylim=y      Lim ,
179       xlab=xName , ylab=yName , cex.lab=1.5 ,
180       main="Data with 95% HDI & Mean of Posterior Predictions" , ce   x.main=1.33 )
181 # Superimpose posterior predicted 95% HDIs:
182 segments( xPostPred, yHDIlim[1] , xPostPred, yHDIlim[2      ] , lwd=3, col="grey" )
183 points( xPostPred , rowMeans( yPostPred ) , pch="+" , cex=2 ,      col="grey" )
184 dev.copy2eps(file=paste(fname,"DataPred.eps",sep="      "))
185
186 #-----
```

16.5.3 BRugs: Simple linear regression with repeated measures

(SimpleLinearRegressionRepeatedBrugs.R)

```

1 graphics.off()
2 rm(list=ls(all=TRUE))
3 library(BRugs)      # Kruschke, J. K. (2010). Doing Bayesian dat      a analysis:
4                      # A Tutorial with R and BUGS. Academic Press / Elsevier.
5 #-----
6 # THE MODEL.
7 modelstring = "
8 model {
9   for( r in 1 : Ndata ) {
10     y[r] ~ dnorm( mu[r] , tau[ subj[r] ] )
11     mu[r] <- b0[ subj[r] ] + b1[ subj[r] ] * x[r]
12   }
13   for ( s in 1 : Nsubj ) {
14     b0[s] ~ dnorm( mu0G , tau0G )
15     b1[s] ~ dnorm( mu1G , tau1G )
16     tau[s] ~ dgamma( sG , rG )
17   }
18   mu0G ~ dnorm(0,.01)
19   tau0G ~ dgamma(.1,.1)
20   mu1G ~ dnorm(0,.01)
21   tau1G ~ dgamma(.1,.1)
22   sG <- pow(m,2)/pow(d,2)
23   rG <- m/pow(d,2)
```

```

24     m ~ dgamma(1,1)
25     d ~ dgamma(1,.1)
26 }
27 "# close quote for modelstring
28 writeLines(modelstring,con="model.txt")
29 modelCheck( "model.txt" )
30
31 #-----
32 # THE DATA.
33
34 # Data from H. A. Feldman, 1988, Table 4, p. 1731.
35 # Columns are "group" , "subjID" , "time" , "retention"
36 source("Feldman1988Table4data.R")
37 # Remove missing data:
38 includeRowVec = is.finite( Feldman1988Table4data[, "retention"] )
39 dataMat = Feldman1988Table4data[ includeRowVec , ]
40 # Retain only the Group 1 (lung) data:
41 dataMat = dataMat[ dataMat[, "group"]==1 , ]
42 # Convert to log10(retention):
43 dataMat[, "retention"] = log10( dataMat[, "retention"] )
44 # Column names and plot labels
45 yColName = "retention" ; yPlotLab = "log10 Retention"
46 xColName = "time" ; xPlotLab = "Day"
47 subjColName = "subjID" ; subjPlotLab = "Subject"
48 fname = "SimpleLinearRegressionRepeatedBrugs"
49
50 if ( F ) { # change to T to use income data instead of contam.rete ntion data.
51 # Data from http://www.census.gov/hhes/www/income/stat medfaminc.html
52 # Downloaded Dec. 06, 2009.
53 load("IncomeFamszState.Rdata") # loads IncomeFamszStat e
54 dataMat = IncomeFamszState
55 yColName="Income" ; yPlotLab = "Income"
56 xColName="Famsz" ; xPlotLab="Family Size"
57 subjColName="State" ; subjPlotLab="State"
58 fname = "IncomeFamszState"
59 }
60
61 # Extract data info to pass to BUGS:
62 Ndata = NROW(dataMat)
63 subj = as.integer(factor(dataMat[,subjColName]))
64 Nsubj = length(unique(subj))
65 x = as.numeric(dataMat[,xColName])
66 y = as.numeric(dataMat[,yColName])
67
68 # Re-center data at mean, to reduce autocorrelation in MCMC s ampling.
69 # Standardize (divide by SD) to make initialization easier.
70 xM = mean( x ) ; xSD = sd( x )
71 yM = mean( y ) ; ySD = sd( y )
72 zx = ( x - xM ) / xSD
73 zy = ( y - yM ) / ySD
74
75 # Specify data, as a list.
76 dataList = list(
77   Ndata = Ndata ,
78   Nsubj = Nsubj ,
79   subj = subj ,
80   x = zx ,
81   y = zy
82 )

```

```

83 # Get the data into BRugs:
84 modelData( bugsData( dataList ) )
85
86 #-----
87 # INTIALIZE THE CHAINS.
88
89 nchain = 3
90 modelCompile( numChains = nchain )
91
92 genInitList <- function() {
93   b0 = b1 = tau = rep(0,length=Nsubj)
94   for ( sIdx in 1:Nsubj ) {
95     yVec = dataList$y[dataList$subj==sIdx]
96     xVec = dataList$x[dataList$subj==sIdx]
97     lmlInfo = lm( yVec ~ xVec )
98     b0[sIdx] = lmlInfo$coeff[1]
99     b1[sIdx] = lmlInfo$coeff[2]
100    tau[sIdx] = length(yVec) / sum(lmlInfo$res^2)
101  }
102  mu0G = mean(b0)
103  tau0G = 1/sd(b0)^2
104  mu1G = mean(b1)
105  tau1G = 1/sd(b1)^2
106  m = mean(tau)
107  d = sd(tau)
108  list( b0=b0 , b1=b1 , tau=tau ,
109        mu0G=mu0G , tau0G=tau0G ,
110        mu1G=mu1G , tau1G=tau1G ,
111        m=m , d=d )
112 }
113 for ( chainIdx in 1 : nchain ) {
114   modelInits( bugsInits( genInitList ) )
115 }
116
117 #-----
118 # RUN THE CHAINS
119
120 # burn in
121 BurnInSteps = 500
122 modelUpdate( BurnInSteps )
123 # actual samples
124 samplesSet( c( "b0","b1","tau" , "mu0G","tau0G", "mu1G",      "tau1G", "m","d" ) )
125 stepsPerChain = ceiling(5000/nchain)
126 thinStep = 100 # 40 or more
127 modelUpdate( stepsPerChain , thin=thinStep )
128
129 #-----
130 # EXAMINE THE RESULTS
131
132 source("plotChains.R")
133 source("plotPost.R")
134
135 # Check convergence and autocorrelation:
136 checkConvergence = T # check this first time through, examine m,d,tau0G,tau1G
137 if ( checkConvergence ) {
138   # check a few selected chains
139   b01Sum = plotChains( "b0[1]" , saveplots=F , filenameroot=      fname )
140   b11Sum = plotChains( "b1[1]" , saveplots=F , filenameroot=      fname )
141   tau1Sum = plotChains( "tau[1]" , saveplots=F , filenameroot= t=fname )

```

```

142 mu0GSum = plotChains( "mu0G" , saveplots=F , filenameroot= fname )
143 tau0GSum = plotChains( "tau0G" , saveplots=F , filenameroot= t fname )
144 mu1GSum = plotChains( "mu1G" , saveplots=F , filenameroot= fname )
145 tau1GSum = plotChains( "tau1G" , saveplots=F , filenameroot= t fname )
146 mSum = plotChains( "m" , saveplots=F , filenameroot= fname )
147 dSum = plotChains( "d" , saveplots=F , filenameroot= fname )
148 }
149
150 # Extract chain values for subsequent examination:
151 zmu0Gsamp = samplesSample( "mu0G" )
152 zmu1Gsamp = samplesSample( "mu1G" )
153 zb0samp = NULL
154 zb1samp = NULL
155 for ( subjIdx in 1:Nsubj ) {
156   zb0samp = rbind( zb0samp , samplesSample( paste("b0[",sub_jIdx,"]",sep="") ) )
157   zb1samp = rbind( zb1samp , samplesSample( paste("b1[",sub_jIdx,"]",sep="") ) )
158 }
159
160 # Convert to original scale:
161 mu0Gsamp = zmu0Gsamp * ySD + yM - zmu1Gsamp * ySD * xM / xSD
162 mu1Gsamp = zmu1Gsamp * ySD / xSD
163 b1samp = zb1samp * ySD / xSD
164
165 # Display believable intercept and slope values
166 windows(10,5.5)
167 par( mar=c(4,4,1.75,1)+0.1 , mgp=c(2.5,0.8,0) )
168 layout( matrix(1:2,nrow=1) )
169 thinIdx = round(seq(1,length(mu0Gsamp),length=700))
170 plot( zmu1Gsamp[thinIdx] , mu0Gsamp[thinIdx] , cex.lab= 1.75 ,
171       ylab="Standardized Intercept" , xlab="Standardized Slope" )
172 plot( mu1Gsamp[thinIdx] , mu0Gsamp[thinIdx] , cex.lab=1.0 ,
173       ylab=paste("Intercept (",yPlotLab," when ",xPlotLab," = 0)",sep="" ) ,
174       xlab=paste("Slope (change in ",yPlotLab," per unit ",xPlotLab," )") )
175 dev.copy2eps(file=paste(fname,"SlopeIntercept.eps", sep=""))
176
177 # Make graphs of data and corresponding believable slopes:
178 windows(12,6)
179 par( mar=c(4,4,1.75,1)+0.1 , mgp=c(2.5,0.8,0) )
180 layout(matrix(c(1:5,1:5,6:10),nrow=3,byrow=T))
181 xlims = c( min( dataMat[,xColName] ) , max( dataMat[,xColName] ) )
182 ylims = c( min( dataMat[,yColName] ) , max( dataMat[,yColName] ) )
183 sldVec = unique( dataMat[,subjColName] )
184 # Plot data of individual subjects:
185 nSubjPlots = 4 # number of representative subject plots to make
186 subjIdxVec = round(seq(1,length(sldVec),length=nSubjPlots))
187 for ( sIdx in subjIdxVec ) {
188   rVec = ( dataMat[,subjColName] == sldVec[sIdx] )
189   plot( dataMat[rVec,xColName] , dataMat[rVec,yColName] , type="o" ,
190         ylim=ylims , ylab=yPlotLab , xlab=xPlotLab , cex.lab=1.5 ,
191         pch=sIdx%%26 , lty=sIdx , main=bquote(.subjPlotLab) "*" * .(sIdx) ,
192         cex.main=1.75 )
193 }
194 # Plot data of all subjects superimposed
195 plot( NULL,NULL , xlab=xPlotLab,xlim=xlims , ylab=yPlotLab , ab,ylim=ylims ,
196       cex.lab=1.5 , main=paste("All ",subjPlotLab,"s",sep="") , cex.main=1.75 )
197 for ( sIdx in 1:length(sldVec) ) {
198   rVec = ( dataMat[,subjColName] == sldVec[sIdx] )
199   lines( dataMat[rVec,xColName] , dataMat[rVec,yColName] ,
200         lty=sIdx , pch=sIdx%%26 , type="o" )

```

```

201 }
202 # Plot histograms of corresponding posterior slopes:
203 for ( sldx in subjIdxVec ) {
204   histInfo = plotPost( b1samp[sldx,] , xlab="Slope" , compVa      l=0.0 , breaks=30 ,
205                      HDItextPlace=0.9 )
206 }
207 histInfo = plotPost( mu1Gsamp , xlab="Slope, Group Level" ,      compVal=0.0 ,
208                      breaks=30 , HDItextPlace=0.9 )
209 dev.copy2eps(file=paste(fname,"Data.eps",sep=""))
210
211 #-----
```

16.6 Exercises

Exercise 16.1. [Purpose: See the influence of individual slope differences on the estimate of the group-average slope.]

The data shown in Figure 16.10 indicate that all subjects had rates of decline in retention, and therefore the estimate of the group average is fairly certain. In this exercise we change the data so that the individual slopes are more dramatically, and examine the effect on the estimate of the group average.

(A) Alter the data as follows: In the program listed in Section 5.16 (SimpleLinearRegressionRepeatedBrugs.R), just before the data are renamed from `dataMat` to `x` and `y` (at about line 65), subtract 0.30x from subject 1, subtract 0.15x from subject 2, do nothing to subject 3, and add 0.15x to subject 4. Here is an example of code for subject 1:

```

subjRowVec = ( dataMat[,subjColName] == 1 )
dataMat[ subjRowVec , yColName ] = ( dataMat[ subjRowVec , yC olName ]
- .030 * dataMat[ subjRowVec , xColName ] )
```

Just repeat and modify for the remaining subjects. Run `therm` and include the plot of the data. (See Figure 16.12.) The data curves for the four subjects should have four very different slopes.

(B) Relative to the original data, has the posterior mean of the slopes gotten farther away from zero or closer to zero? Include the histogram of the posterior in your write-up. Are all the individual slopes believably different from zero (according to the 95% HDI)? Is the group slope believably different from zero (according to the 95% HDI)? Why is the group-level slope, which is now farther away from zero or more, less believably different from zero than in the original data?

Exercise 16.2. [Purpose: See the influence of differences in individual intercepts on the estimate of group-average slope]

The data shown in Figure 16.10 all start at 2.0 because that is (100), and the data were measured as percentage of original value. Suppose that the data were kept in their raw magnitudes, instead of converted to percentage of magnitude. This would merely change the intercepts of the individual data curves, without changing their slopes, because $\log(ky) = \log(k) + \log(y)$. In this exercise we find out whether this change would have an effect on the estimate of the group slope.

(A) Alter the data as follows: In the program listed in Section 5.16 (SimpleLinearRegressionRepeatedBrugs.R), right before the data are converted to `log10` on line 42, insert this code:

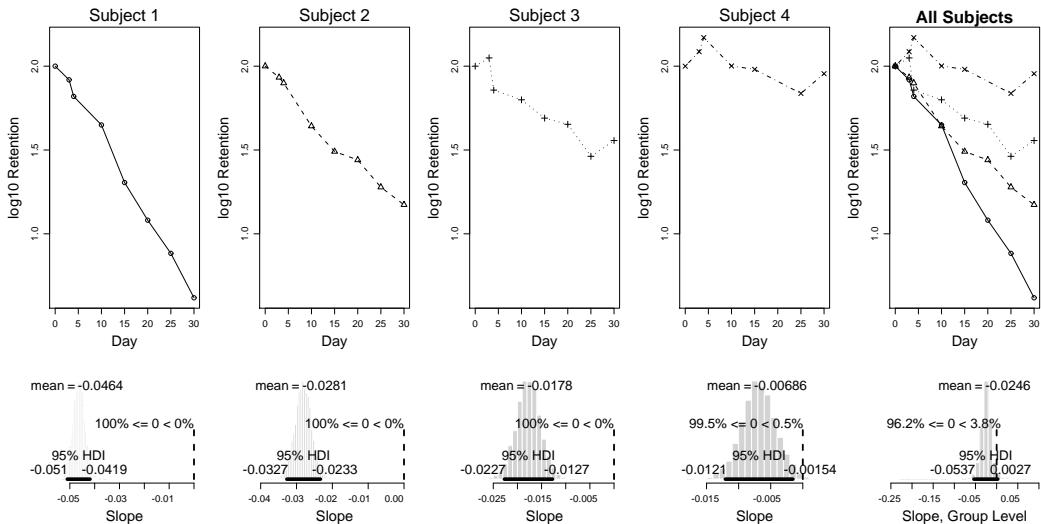


Figure 16.12: For exercise 16.1. Contaminant-retention data with greater variation of individual slopes. Notice the posterior on the group-level slope parameter. Compare with Figure 16.10.

```
for ( subjIdx in 1:4 )      f
  rowIdx = ( dataMat[, "subjID"] == subjIdx )
  dataMat[rowIdx, "retention"] = dataMat[rowIdx, "retenti
                                on"] * 10^(subjIdx-1)
g
```

Run the program and include the plot of the data. (See Figure 16.13.) The data curves for the four subjects should have four very different intercepts.

(B) Relative to the original data, are the estimates of the individual slopes different? Relative to the original data, is the posterior mean of the group slope different? Include the histogram of the posterior of the group slope. Is the group slope believably different from zero (according to the 95% HDI)? Why is the group-level slope less believably different from zero, compared to the original data? Hint: The individual intercepts affect the certainty of the group-average intercept. The group average intercept trades off with the group-average slope; consider scatterplots of $\text{intercept} \times \text{slope}$ and $\mu_1 \times \mu_2$.

Exercise 16.3. [Purpose: Real data for repeated measures of individual regression, with an outlying individual and non-linear trend.] Suppose we are interested in whether families with more members have higher incomes. The U.S. Census Bureau has provided data that indicate the median family income as a function of number of persons per family, for all 50 states and the District of Columbia and Puerto Rico. The data appear in Figure 16.14, p. 369.

(A) Run the program of Section 16.5.3 (`SimpleLinearRegressionRepeatedBrugs.R`) with the income data. Examine the data section of the program and you will find that the necessary lines of code are already available. The program generates a figure much like Figure 16.14, p. 369.

(B) There is suggestion of outliers in these data. One curve (Puerto Rico), falls barely above 20,000, which is far lower than all the others. This suggests an outlier for the distribution of intercepts. Some single data points fall off from the individual linear trends. For example, the 6-person family in the District of Columbia has an income of only about

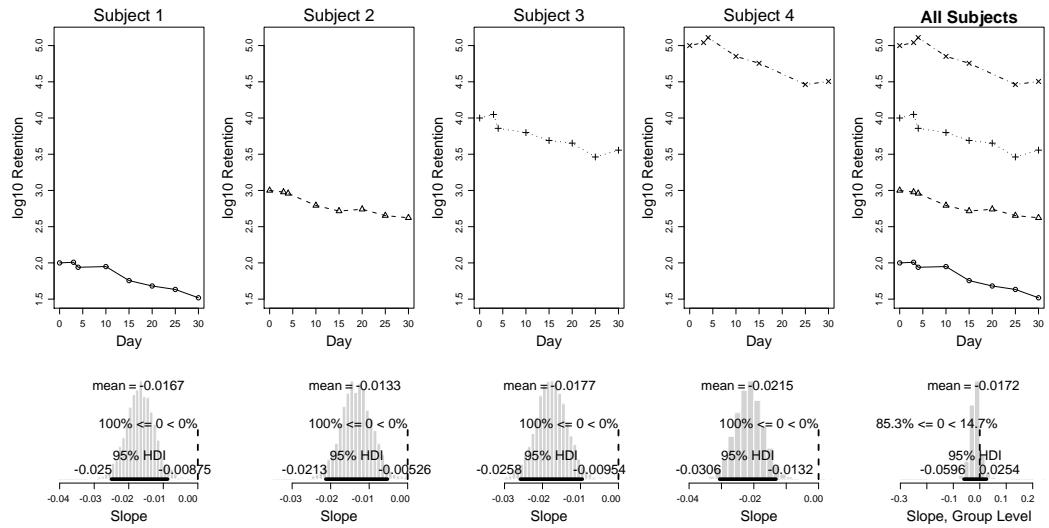


Figure 16.13: For exercise 16.2. Contaminant-retention data with greater variation of individual intercepts Notice the posterior on the group-level slope parameter. Compare with Figure 16.10.

30,000, whereas the 5-person family has an income of over \$60,000. This suggests an outlier for points around linear trends. Finally, some individual slopes seem quite steep compared to others. For example, the income in Hawaii rises at \$50,000 from 2-person to 7-person families. This increase might or might not be an outlier relative to other states. Which distributions in the hierarchical model of Figure 16.6 should be changed to distributions to address these outliers? Change the model specification in your write-up. Show the posterior estimate of the intercept of Puerto Rico for small-t bias and for large-t bias. Show the posterior estimate of the slope of Hawaii for small-t bias and for large-t bias.

(C) The data also suggest a non-linear trend in the data. Income appears to rise for 2, 3, and 4-person families, but then level and decline as family size gets larger. Include in the original (non-) model another term that can capture “quadratic curvature” at the income level: $y = \beta_0 + \beta_1 x + \beta_2 x^2$. The prior on β_2 is analogous to the prior on β_1 . Is the group-average estimate of curvature credibly non-zero?

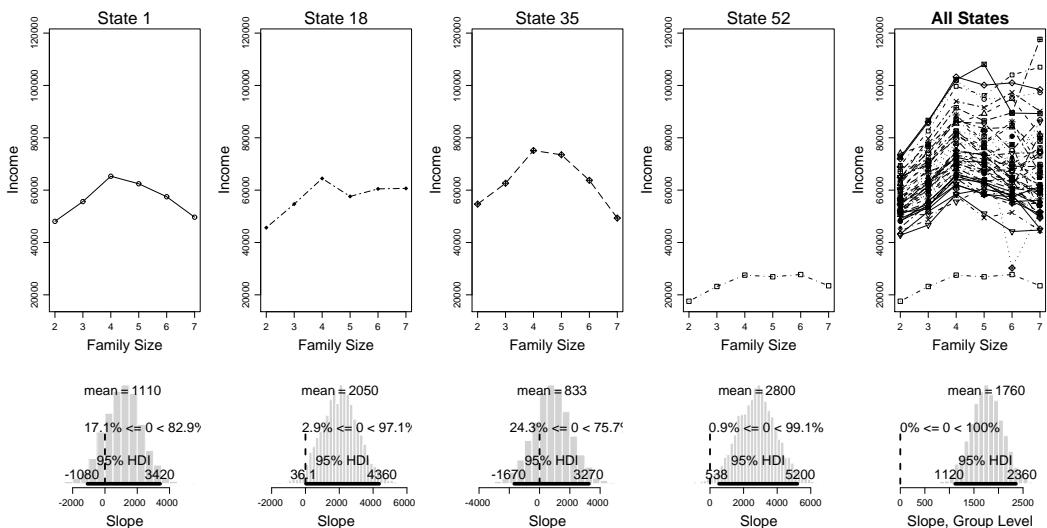


Figure 16.14: For exercise 16.3. The top right graph shows data for 50 states plus Puerto Rico and the District of Columbia, and the first four plots show data from four randomly selected individual states. Lower row shows marginal posterior distribution on the slope parameters.

Chapter 17

Metric Predicted Variable with Multiple Metric Predictors

Contents

17.1	Multiple linear regression	372
17.1.1	The perils of correlated predictors	372
17.1.2	The model and BUGS program	375
17.1.2.1	MCMC efficiency: Standardizing and initializing	376
17.1.3	The posterior: How big are the slopes?	376
17.1.4	Posterior prediction	387
17.2	Hyperpriors and shrinkage of regression coefficients	378
17.2.1	Informative priors, sparse data, and correlated predictors	382
17.3	Multiplicative interaction of metric predictors	383
17.3.1	The hierarchical model and BUGS code	384
17.3.1.1	Standardizing the data and initializing the chain .	385
17.3.2	Interpreting the posterior	385
17.4	Which predictors should be included?	388
17.5	R code	390
17.5.1	Multiple linear regression	903
17.5.2	Multiple linear regression with hyperprior on coefficients	394
17.6	Exercises	399

When I was young two plus two equaled four, but
Since I met you things don't add up no more.
My keel was even before I was kissed, but
Now my predictions all come with a twist.

In this chapter we are concerned with situations such as predicting a person's college grade point average (GPA) from his or her high school GPA and/or a Scholastic aptitude test (SAT) score. Another such situation is predicting a person's blood pressure from his or her height and weight. In these situations, the value to be predicted is on a metric scale, and there is more than one predictor, each of which is also on a metric scale.

We will consider models in which the predicted variable is an additive combination of predictors, all of which have proportional influence on the prediction. This kind of model is called "multiple linear regression", and is listed in Table 4.1 in its first row and third

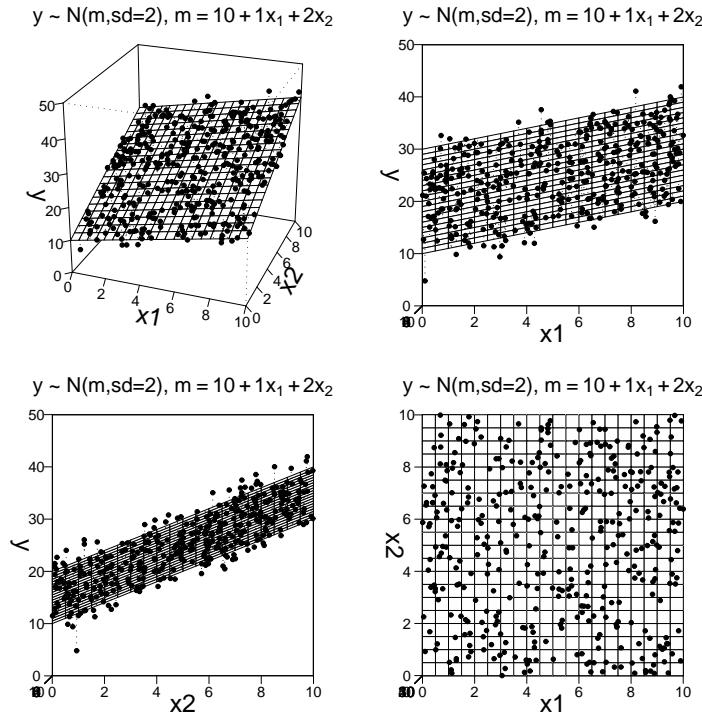


Figure 17.1: Data, that are normally distributed around the values in the plane. The $x_1; x_2$ values are sampled uniformly and independently of each other, as shown in the lower-right panel. The panels show different perspectives on the same plane and data. Notice that when the data are plotted against x_1 (marginalized across x_2), the points appear to rise along with the linear function that generated them. Compare with Figure 17.2.

column. We will also consider non-additive combinations of predictors, which are called “interactions”.

17.1 Multiple linear regression

Figures 17.1 and 17.2 shows examples of data generated by `dplyr` multiple linear regression. The model specifies the dependency on $x_1; x_2$, but does not specify the distribution of $x_1; x_2$. At any position, $x_1; x_2$, the values of y are normally distributed in a vertical direction, centered on the height of the plane at that position. The height of the plane is a linear combination of the $x_1; x_2$ values. Formally, $y \sim N(\mu, \sigma^2)$, and $\mu = \beta_0 + \beta_1 x_1 + \beta_2 x_2$. For a review of how to interpret the coefficients, β_0 , β_1 , and β_2 , see Figure 14.2, p. 297. The model assumes homogeneity of variance: At all values of $x_1; x_2$, the variance of y is the same.

17.1.1 The perils of correlated predictors

Figures 17.1 and 17.2 show data generated from the same model. What differs between them is the distribution of $x_1; x_2$, which is not specified by the model. In Figure 17.1, the $x_1; x_2$ values are distributed uniformly. In Figure 17.2, the x_2 values are negatively

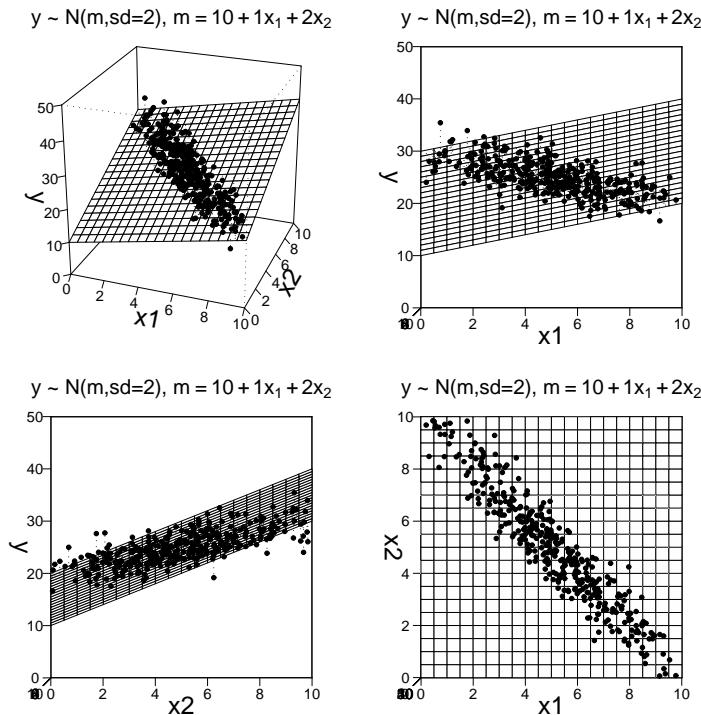


Figure 17.2: Data, that are normally distributed around the values in the plane. The $x_1; x_2$ values are correlated, as shown in the lower-right panel. The panels show different perspectives on the same plane and data. Notice that the data are plotted against x_1 (marginalized across x_2), the points appear to drop, contrary to the linear function that generated them. Compare with Figure 7.1.

correlated: When x_1 is small, x_2 tends to be large, and when x_1 is large, x_2 tends to be small. The correlation of x_1 and x_2 can lead to misinterpretations of their individual influence on y . For instance, notice in Figure 17.2 that when x_1 is near zero, then the data values are near 30, but when x_1 is near 10, then the data values are near 20. This observation that y declines from 30 to 20 might leave them impression that x_1 predicts a decrease in y . But such an impression is wrong, because the data were generated by a function that increases y as x_1 increases; i.e., the coefficient 1 on x_1 is $+1$. The reason that they values appear to decline as x_1 increases is that x_2 decreases when x_1 increases, and x_2 has an even bigger influence on y than x_1 does.

It is not unusual for predictors to be correlated in real data. For example, consider trying to predict a state's average high school SAT scores on the basis of the amount of money the state spends per pupil. If you plot only mean SAT against money spent, there is actually a decreasing trend, as can be seen in the lower left panel of Figure 17.2 (Adam Guber, 1999). In other words, SAT scores tend to go down as spending goes up. Guber (1999) explains how some political commentators have used this evidence to argue against funding public education.

This negative influence of spending on SAT scores seems counter-intuitive. It turns out that the trend is an illusion caused by the influence of another factor, along with the correlation of spending with that other factor. The other factor is the proportion of students who take the SAT. Not all students at a high school take the SAT, because the test

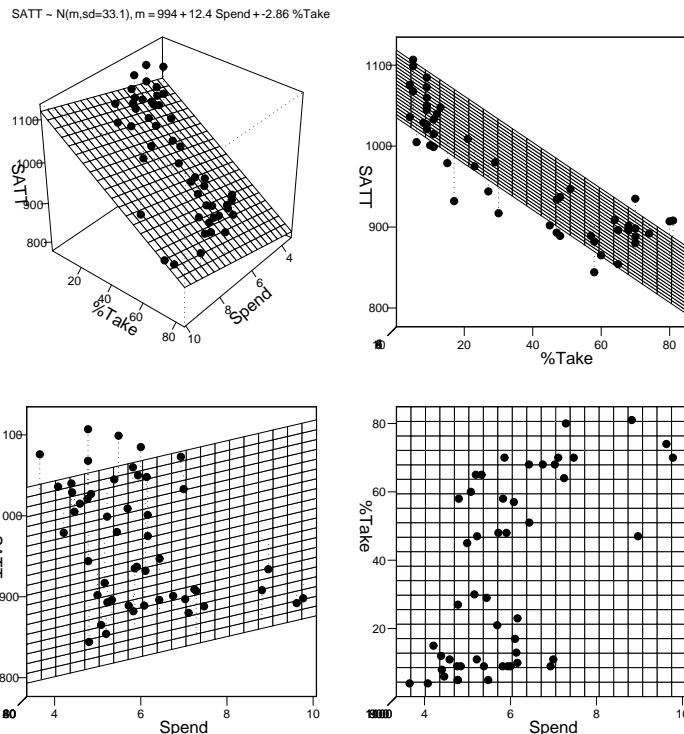


Figure 17.3: An example of multiple linear regression with `nlldata`. The data are plotted as dots, and the gridded plane shows the mean posterior slopes and intercept. The four panels show different perspectives on the same data and plane. “SATT” is the average total SAT score in a state. “%Take” is the percentage of students in the state who take the SAT. “Spend” is the spending per pupil, in thousands of dollars.

is used primarily for college entrance applications, and therefore it is primarily students who intend to apply to college who take the SAT. Most of the students at a high school will take the SAT, because most of the top students will apply to college. But students who are weaker academically may be less likely to take the SAT, because they are less likely to apply to college. Therefore, the more that a high school enrolls mediocre students to take the SAT, lower will be its average SAT score. It turns out that high schools that spend more money per pupil also have a much higher proportion of students who take the SAT. This correlation can be seen in the lower-right panel of Figure 317.

When both predictors, i.e., spending per pupil and percentage of students taking the SAT, are taken into account, the influence of spending on SAT scores is seen to be positive, not negative. This positive influence of spending can be seen in the positive slope of the plane along the “Spend” direction in Figure 17.3. The negative influence of percentage of students taking the SAT, is also clearly shown. To reiterate the main point of this example: It seems that the apparent drop in SAT due to spending is ~~fact~~ spending being correlated with the percentage of students taking the SAT, the latter having a whoppingly negative influence on SAT scores.

The separate influences of the two predictors could be addressed in this example because the predictors had only mild correlation with each other. There was enough independent variation of the two predictors that their distinct relationships to the outcome variable could

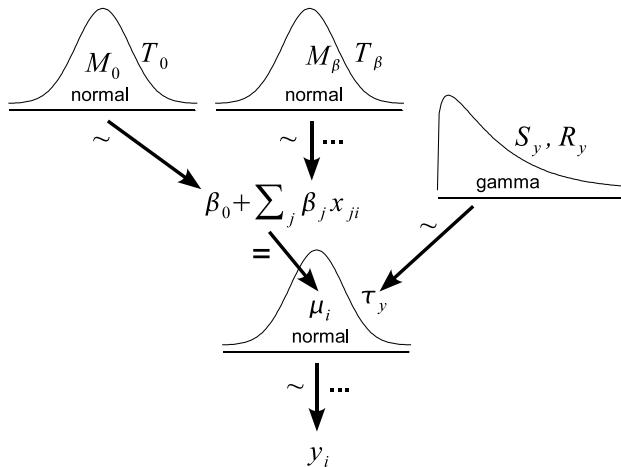


Figure 17.4: Hierarchical diagram for multiple linear regression. (In general, the slope parameters β_j , may have different priors, instead of the same prior repeated for every j as is shown here for simplicity.)

be detected. In some situations, however, the predictors ~~are~~ are tightly correlated that their distinct effects are difficult to tease apart. Correlation of predictors causes ~~the~~ ~~models~~ of their regression coefficients to trade off, as we will see when examine the model and its posterior estimates.

17.1.2 The model and BUGS program

The hierarchical diagram for multiple linear regression is shown in Figure 17.4. It is merely a direct expansion of the one for simple linear regression (which appeared in Figure 16.3, p. 346). Instead of just one slope coefficient for a single predictor, there is another slope coefficient for every one of the multiple regressors. For every coefficient, the prior is normal, just as shown in Figure 16.3.

As usual, the BUGS model specification has a line of code for every variable in the hierarchical diagram. The model specification looks like this: (MultipleLinearRegressionBrugs.R)

```

11 model {
12   for( i in 1 : nData ) {
13     y[i] ~ dnorm( mu[i] , tau )
14     mu[i] <- b0 + inprod( b[] , x[i,] )
15   }
16   tau ~ dgamma(.01,.01)
17   b0 ~ dnorm(0,1.0E-12)
18   for ( j in 1:nPredictors ) {
19     b[j] ~ dnorm(0,1.0E-12)
20   }
21 }
```

Notice that what is written in mathematical notation as $\beta_0 + \sum_j \beta_j x_{ij}$ is expressed in BUGS on line 14 by using the `inprod` function. The `inprod` function is named for the mathematical “inner product” of two vectors. The `function` takes two vectors and returns the sum of their component-by-component products. Thus, for two vectors v and w , the inner product `inprod(v[1:n],w[1:n])` is $\sum_{j=1}^n v_j w_j$. In BUGS, we do not need

to specify the range of indices for the vector when it is on the right side of an assignment operator, and we intend to use all components of the variable. But in BUGS we do need to include the square brackets so that BUGS knows that `variable` is a vector. If we wanted to be completely explicit, we could write the inner product in line 14 as `inprod(b[1:nPredictors] , x[i,1:nPredictors])`. The complete program is listed in Section 17.5.1 (`MultipleLinearRegressionBrugs.R`).

The BUGS model specification happens to put the same prior very slope parameter. This equivalence is applied as a generic convenience, `b` is not required. Indeed, if there is prior information that suggests different priors on different predictors, then the prior knowledge should be respected.

17.1.2.1 MCMC efficiency: Standardizing and initializing

As described previously in Section 16.1.1.1 (p. 347), the MCMC sampling can be made much more efficient if the data are standardized. Standardizing each variable is straightforward. The MCMC sampling then finds regression coefficients that are appropriate for the standardized data. We would like to transform the parameters to the corresponding values that are appropriate to the original, nonstandardized scores. This can be done by generalizing Equation 16.2 to multiple predictors:

$$\begin{aligned} Z_y &= \beta_0 + \sum_j Z_{x_j} \\ \frac{\hat{y} - M_y}{SD_y} &= \beta_0 + \sum_j \frac{x_j - M_{x_j}}{SD_{x_j}} \\ \hat{y} &= \beta_0 SD_y + M_y + \sum_j \frac{SD_y M_{x_j} - SD_{x_j}}{SD_{x_j}} + \sum_j \frac{SD_y}{\{Z\}} x_j \quad (17.1) \end{aligned}$$

The estimate of β_0 is merely $\hat{y}_0 - \hat{M}_y$.

Even after standardizing, it can also help to start the chains at their posterior credible values. To do this, we use the built-in linear model function R called `lm`. There is no need to delve into the inner workings of `lm`, but suffice it to say that it returns the maximum-likelihood estimate (MLE) of the intercept and slope coefficients, in raw scales. These raw-scale slope coefficients can be easily converted to standardized scales, `sd`, to initialize the chains. In conclusion, the standardization of the data makes the chain efficient and less autocorrelated once it reaches the modal region of the posterior. The initialization at the MLE implies that the burn-in period is minimal.

17.1.3 The posterior: How big are the slopes?

Figure 17.5 shows the posterior distribution the results of the SAT data in Figure 17.3 and model in Figure 17.4. You can see that the slope on `Spending` is credibly well above zero, with a mean slope of about 12.3, which suggests that SAT scores rise by about 12.3 points for every extra \$1,000 spent per pupil. The slope on `TestPrep` taking the SAT is also credibly non-zero, with a mean of -2.85, which suggests that SAT scores fall by about 2.85 points for every additional 1% of students who take the test.

The scatter plots in the bottom of Figure 17.5 show correlations among the credible parameter values. In particular, the lower-right scatter plot shows that the coefficient for

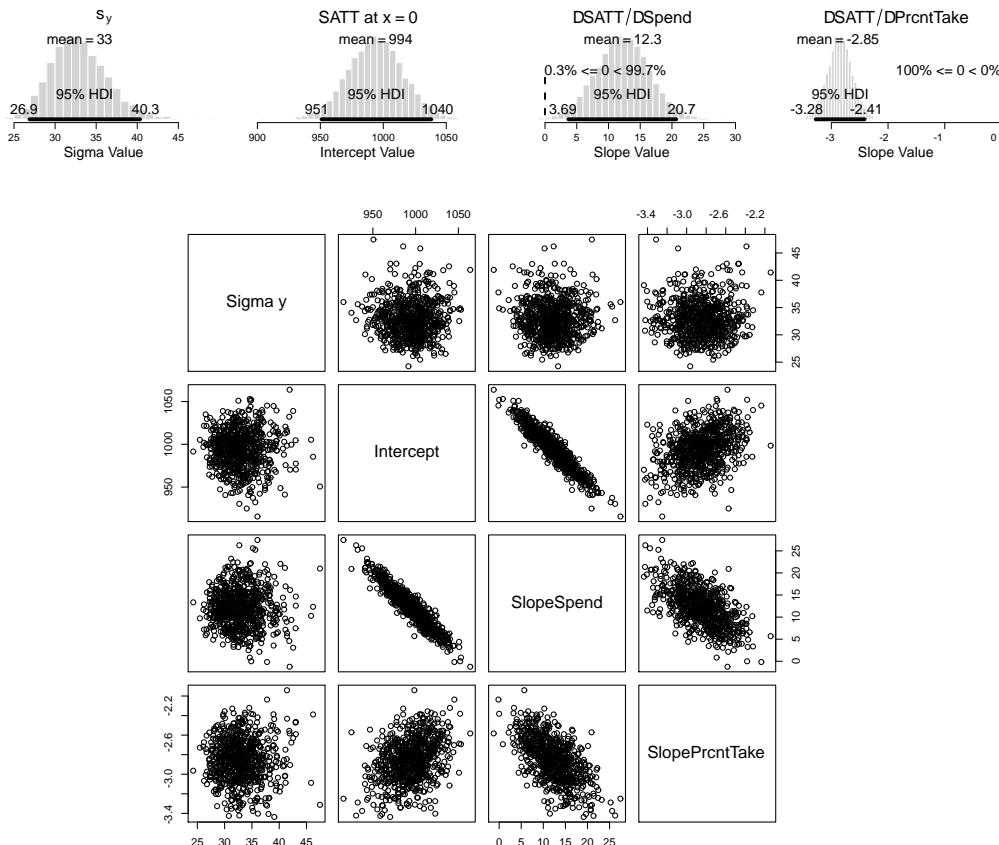


Figure 17.5: Posterior distribution for data in Figure 17.3 and model in Figure 17.4. Upper row indicates HDI's on regression coe cients. Scatterplots reveal correlations among credible parameter values; in particular, the coe cient on Spending (“SlopeSpend”) trades o with the coe cient on Percentage taking the SAT (“SlopePrctTake”), because those predictors are correlated in the data.

Spending trades o with the coe cient on Percentage taking the SAT. The correlation means that if we believe that the in uence of Spending is large then we must believe that the in uence of Percentage Taking is smaller, to stay consistent with the data. Conversely, if we believe that the in uence of Spending is small then we must believe that the in uence of Percentage Taking is larger. This makes sense because those two predictors are correlated in the data, and therefore the two predictors are not di erentially constraining the regression coe cients. Think of this simple example for two data points: $y_1 = 1$ for $x_1 = 1$ and $x_2 = 1$, and $y_2 = 2$ for $x_1 = 2$ and $x_2 = 2$. The linear model $y = \beta_1 x_1 + \beta_2 x_2$, is supposed to satisfy both data points, which implies that $\beta_1 + \beta_2 = 1$. Therefore, to satisfy the data, it could be that $\beta_1 = 2$ and $\beta_2 = -1$, or $\beta_1 = 0.5$ and $\beta_2 = 0.5$, or $\beta_1 = 0$ and $\beta_2 = 1$, etc. In other words, the credible values of β_1 and β_2 are (anti-)correlated.

One of the bene ts of Bayesian analysis is that correlations of credible parameter values are explicit in the posterior distribution. Traditional statistical methods provide only a single “best” (e.g., MLE) parameter value, without indicating the trade os among parameter values. The Bayesian posterior, however, naturally reveals tradeos and redundancies among parameters. It is up to the user to actually look for and interpret the correlations of

parameters, of course. Another benefit of Bayesian analysis is that the model doesn't "explode" when predictors are correlated. If predictors are correlated, the joint uncertainty in the regression coefficients is evident in the posterior, but the model happily generates a posterior regardless of correlations in the predictors. This is a classical, one-best-solution method is much less robust in the presence of strongly correlated predictors.

17.1.4 Posterior prediction

Often we are interested in using the linear model to predict values for various x values. It is straightforward to generate a large sample of credible values for specified x values. From the distribution of y values we can compute the mean and HDI to summarize the centrally predicted value and the uncertainty of the prediction. As was the case in simple linear regression, illustrated back in Figure 16.6, the uncertainty in predicted y is greater for x values outside the bulk of the data. In other words, extrapolation is more uncertain than interpolation.

The last part of the code in Section 17.5 (`MultipleLinearRegressionBrugs.R`) carries out these computations. At every step in the MCMC chain of prior parameter values, the program randomly generates a value based on the linear model. These values, for which posterior predictions are desired, are stored in `xPostPred`, one row per point to be predicted. Thus `xPostPred` has as many columns as there are predictors. The matrix `yPostPred` also has one row per point to be predicted, with each row containing randomly generated values for the corresponding point. The program generates one value per step in the MCMC chain, hence `yPostPred` has as many columns as the steps in the chain. The vector `bSamp` contains posterior values of the intercept, as each step in the MCMC chain. The vector `bSamp[chainIdx,]` contains the slope coefficients, β_j , for the predictors, at one step in the chain. The slope coefficients are forced to be a column vector by passing them through the `cbind` function in R. Then the slopes can be multiplied by the corresponding values and summed together, simultaneously for all the points in one matrix operation: `xPostPred %*% cbind(bSamp[chainIdx,])`. This matrix operation appears on line 235 in the following code: (`MultipleLinearRegressionBrugs.R`)

```

232 for ( chainIdx in 1:chainLength ) {
233   yPostPred[,chainIdx] = rnorm( NROW(xPostPred) ,
234                               mean = b0Samp[chainIdx]
235                               + xPostPred %*% cbind(bSamp[chainIdx,]) ,
236                               sd = rep( sigmaSamp[chainIdx] , NROW(xPostPred) ) )
237 }
```

Notice that the code loops through every step in the MCMC chain in the `yPostPred` matrix one column at a time. The values are generated from a normal distribution, using the `rnorm` function, with a standard deviation `sigmaSamp[chainIdx]`.

17.2 Hyperpriors and shrinkage of regression coefficients

In some research, there are many candidate predictors which could possibly be informative about the predicted variable. For example, when predicting college GPA, we might include high school SAT, high school GPA, income, student, income of parents, years of education of the parents, spending per student's high school, student IQ, student height, weight, shoe size, hours of sleep per night, distance from home

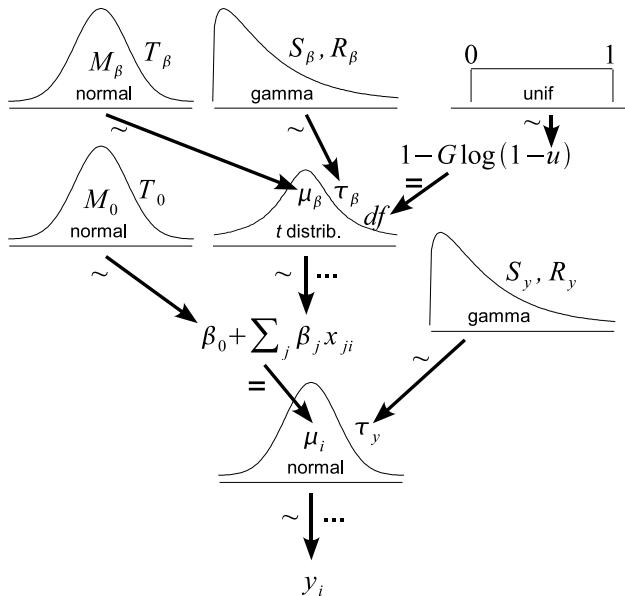


Figure 17.6: Hierarchical diagram for multiple linear regression, with hyperprior on slope coefficients across predictors. Compare with Figure 17.4.

to school, amount of caffeine consumed, hours spent studying, hours spent earning, etc. We can include all the candidate predictors in the model with a regression coefficient for every predictor. Should all those regression coefficients be estimated in isolation from the others as in the model of Figure 17.4? Probably not, because we probably have some prior knowledge that relates the influences of the predictors to each other. If nothing else, we can at least say that the candidate predictors all come from the set of remotely plausible predictors. What do we know about this set of remotely plausible predictors? Most candidate predictors probably have a very small relationship to the predicted variable, but a few candidate predictors may have sizable covariation with the predicted variable. In other words, the regression coefficients are probably distributed something like a distribution, with lots near a mean of zero, but with a few in the extended tails. We therefore put this prior knowledge into the model structure, as shown in Figure 17.6. (This method is mentioned in passing by Gelman et al., 2004, p. 405)

Figure 17.6 indicates that the regression coefficients, β_j , are distributed as a t -distribution. The parameters of the t -distribution are estimated from the data. Presumably, many of the credible regression coefficients will be near zero, but a few will depart a lot from zero, and the distribution will have credible and df values that reflect the distribution of regression coefficients in the data. Exercise 17.1 has you generate the prior in BUGS, so you get a better intuition for how the constants in the hyperprior distribution affect the implied prior on the regression coefficients.

A desirable side-effect of incorporating this prior structure is that the estimates of the regression coefficients experience shrinkage. If many regression coefficients are near zero, then their distribution will have a high precision (parameter), which in turn will shrink the estimates of the regression coefficients. The regression coefficients are mutually informing each other, via the prior knowledge that they should be distributed according to a total distribution.

The shrinkage is desirable not only because it expresses prior knowledge, but also

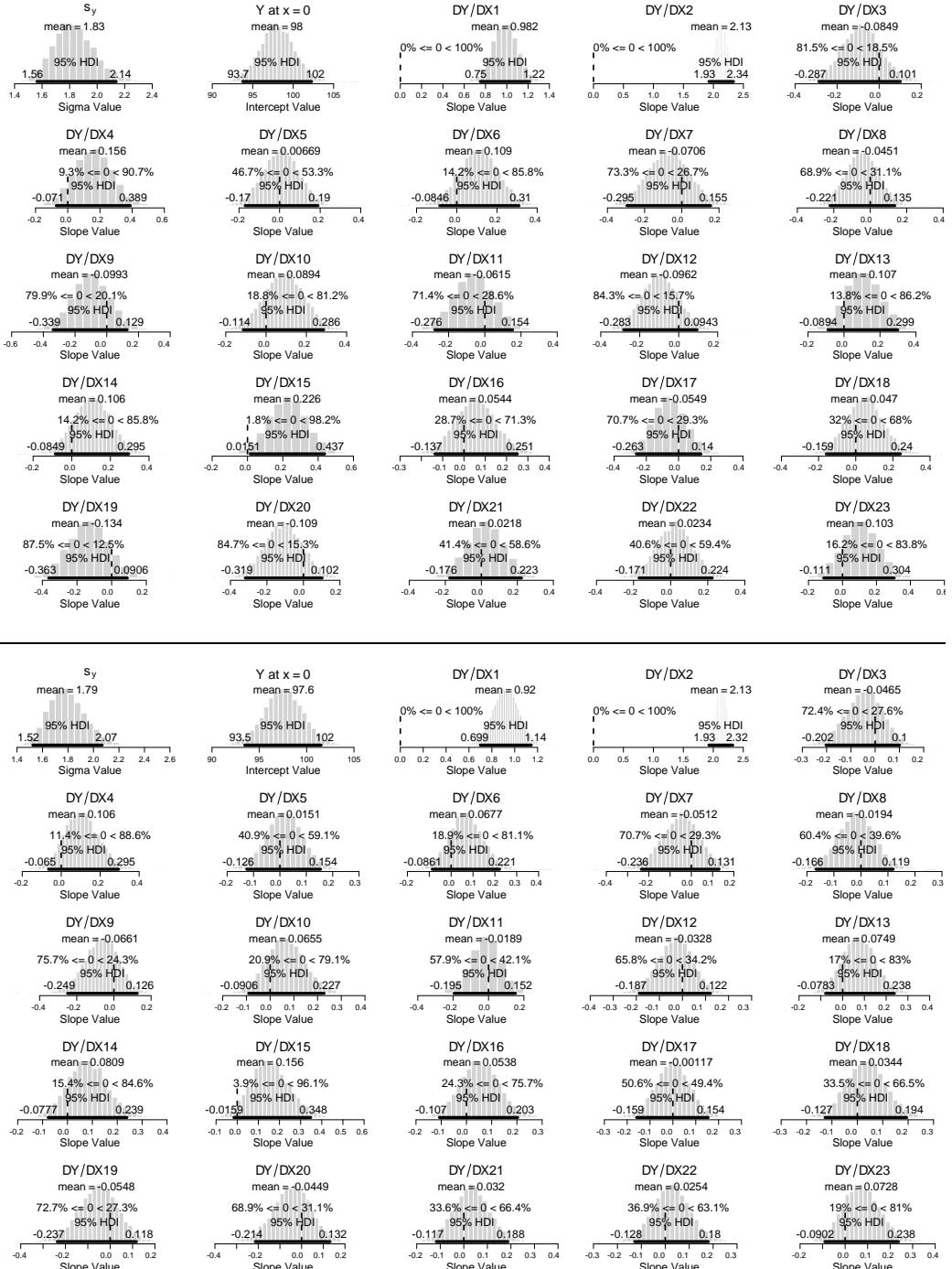


Figure 17.7: Posterior distribution for 100 data points, randomly generated with true parameter values $\beta_0 = 2$, $\beta_1 = 100$, $\beta_2 = 1$, $\beta_3 = 2$, and all other $\beta_j = 0$. Upper half: Model is from Figure 17.4, wherein each slope coefficient is insulated from the other. Notice the “false alarm” for β_5 . Lower half: The model is shown in Figure 17.6, wherein the slope coefficients mutually inform each other. Shrinkage prevents the “false alarm” for β_5 .

because it rationally helps control for “false alarms” including that a predictor has a non-zero regression coefficient. When there are many candidate predictors, some of the may spuriously appear to have credibly non-zero regression coefficients, even when the true coefficient is zero. This sort of false alarm is unavoidable because data are randomly sampled, and there will be occasional coincidences of data that are unrepresentative. By letting the regression coefficients be mutually informed by other predictors, and not only by the data of the single predictor each multiplies, the coefficients are less likely to be spuriously distorted by a rogue sample.

As an example in which we can specify the true regression coefficients, let's randomly generate 100 data points from a linear regression model with parameter values of $\beta_0 = 2$, $\beta_1 = 100$, $\beta_2 = 2$, and all other $\beta_j = 0$ for 21 other predictors. When using insulated regression coefficients for each predictor, as in the hierarchical diagram Figure 17.4, the resulting posterior estimates are shown in the upper half of Figure 17.7. Notice that the estimate of β_{15} (denoted $Y = X_{15}$) suggests that this predictor may have a non-zero regression coefficient. (If we could specify a ROPE of some non-zero width, we might decide that β_{15} is non-zero, but ROPEs are best defined in meaningful context in a generic example like this.) In other words, this apparent non-zero value of β_{15} is a “false alarm”, produced by quirks in the random sample of data.

The estimates in the upper half of Figure 17.7 did not use any knowledge about how the regression coefficients might be related. If we instead use the model of Figure 17.6, with the same data, the resulting posterior is shown in the lower half of Figure 17.7. Notice that the estimate of β_{15} now includes zero inside the HDI. All the estimates of the regression coefficients are reduced a little, relative to the upper half of Figure 17.7. The reason for this shrinkage is that the many regression coefficients near zero mutually inform each other via the overarching distribution.

When using the model of Figure 17.6, the shrinkage of the regression coefficients is toward their mean, μ , not necessarily toward zero. If your prior knowledge suggests that the true mean of the regression coefficients is very close to zero, then that knowledge should be expressed in the model of Figure 17.6 by setting the μ to zero and setting the precision T to a very large value such as 10,000. This setting implies that μ is already fairly certain, and the real prior uncertainty is in the ability of regression coefficients across predictors.

The hyperprior over regression coefficients, shown in Figure 17.6, is meant to express a genuine belief, namely, that all the regression coefficients can be reasonably described by a t distribution. This assumption may be more or less tenable in different situations. For example, a region's agricultural crop yield might depend on rainfall amounts measured at 39 randomly selected locations. The predictions from each of these measuring locations are probably fairly similar, except for perhaps a few outliers, and therefore it seems quite reasonable to model the regression coefficients by a t distribution. In other situations, it may be less reasonable to treat all the regression coefficients as coming from a shared distribution. At the least, the predictors were selected from some implicit set of reasonably likely predictors, and we can think of the overall distribution as reflecting that set. We might still use the hyperprior model, but only as a convenience to impose some degree of shrinkage on the regression coefficients. It should be interpreted carefully. Beware of convenience priors that are used in routinized way.

17.2.1 Informative priors, sparse data, and correlated predictors

This book has emphasized the use of mildly informed priors opposed to conventionalized noninformed priors or strongly informed priors. On the other hand, this book has also mentioned that a benefit of Bayesian analysis is the path for cumulative scientific progress through the use of priors that have been informed by previous research. Informed priors should be used whenever the skeptical scientist ~~cautious~~ of the analysis deems it appropriate, especially if the analysis is accompanied by a check of posterior robustness. If the conclusion from the posterior is strong, even with a ~~highly~~ informed prior, then the mildly informed prior should be used, because a broader ~~range~~ may find the analysis to be convincing. But there are situations in which the use of a ~~strongly~~ informed prior is appropriate.

In general, when the data are sparse (i.e., when the sample size is small), the posterior will be imprecise if the prior is imprecise. But if the prior is ~~not~~ rightly constrained, then a small amount of data can lead to a more decisive posterior. An example of this phenomenon appeared in Exercise 5.7, p. 81, in which the prior allowed for two very different interpretations of the data. In this case, even a small amount of data shifted the posterior toward one or the other interpretation.

Sparse data can also lead to usefully precise posteriors in the context of multiple linear regression, if some of the regression coefficients use informed priors, and the predictors are correlated. To develop this idea, it is important to remember that when predictors are correlated, their regression coefficients are also (anti-)correlated. For example, recall the SAT data from Figure 17.3, p. 374, in which the predictors are correlated in the data, i.e., spending-per-pupil and percent-taking-the-exam are related. Consequently, the posterior estimates of the regression coefficients had a negative correlation, as shown in Figure 17.5, p. 377. The correlation of credible regression coefficients implies that a strong belief about the value of one regression coefficient constrains the value of the other coefficient. Look carefully at the scatterplot of the two slopes shown in Figure 17.5. It can be seen that if we believe that the slope on percent-taking-the-exam is 32, then credible values of the slope on spending-per-pupil are around 18, with an HDI of roughly 22. Notice that this HDI is much smaller than the marginal HDI on spending-per-pupil, which goes from 3.7 to 20.7. Thus, constraining the beliefs about one slope also constrains the beliefs about the other slope, because estimates of the slopes are correlated.

That influence of one slope estimate on another can be used to our advantage when we have prior knowledge about one of the slopes. If some previous or auxiliary research informs the prior of one regression coefficient, that constraint can propagate to the estimates of regression coefficients on other predictors that are correlated with the first. This is especially useful when the sample size is small, and a ~~more~~ mildly informed prior would not yield a very precise posterior. Of course, an informed prior on the first coefficient must be cogently justified for the skeptical audience. A robustness check also may be useful, to show how strong the prior must be to draw strong conclusions. If the information used for the prior is compelling, then this technique can be very useful for leveraging novel implications from small samples. An aside discussion and example from political science is provided by Western and Jackman (1994), and a mathematical discussion is provided by Learner (1978, p. 475).

17.3 Multiplicative interaction of metric predictors

In some situations, the predicted value is not merely an additive combination of the predictors. For example, the effects of drugs are often non-additive. As the dosage of one drug increases, there might be a moderate increase in effectiveness. And as the dosage of another drug increases (when administered by itself), there may be a moderate increase in effectiveness. But when the two drugs are administered together, when the dosages of both drugs are high, they might interact to greatly enhance or greatly reduce effectiveness, beyond a mere addition of the two separate effects. As another example, consider trying to predict subjective happiness from income and health. If either income or health is low, subjective happiness is probably also fairly low. But if both income and health are high, then happiness is more likely to be high. In other words, happiness might increase additively with both income and health; instead, it happiness may be moderated by an interaction of income and health. In general, interaction means that the effect of one predictor varies, depending on the value of the other predictor.

Formally, interaction can be expressed in different ways. We will consider multiplicative interaction. For two metric predictors, regression with multiplicative interaction has these algebraically-equivalent expressions:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_{12} x_1 x_2 \quad (17.2)$$

$$= \beta_0 + \left| \begin{array}{l} \beta_1 + \beta_{12} x_2 \\ \text{slope of } x_1 \end{array} \right\} x_1 + \beta_2 x_2 \quad (17.3)$$

$$= \beta_0 + \beta_1 x_1 + \left| \begin{array}{l} \beta_2 + \beta_{12} x_1 \\ \text{slope of } x_2 \end{array} \right\} x_2 \quad (17.4)$$

These three expressions emphasize different interpretations of interaction, as illustrated in Figure 17.8.

The form of Equation 17.2 is illustrated in the left panel of Figure 17.8. The vertical arrows show that the curved-surface interaction is created by adding the product, $\beta_{12} x_1 x_2$, to the planar linear combination.

The form of Equation 17.3 is illustrated in the middle panel of Figure 17.8. Its dark lines show that the slope in the x_1 direction depends on the value of x_2 . In particular, when $x_2 = 0$, then the slope along x_1 is β_1 , which in the graphed example is 1. But when $x_2 = 10$, then the slope along x_1 is $\beta_1 + \beta_{12} x_2$, which in the graphed example is 1. Again, the slope in the x_1 direction changes when x_2 changes, and only indicates the slope along x_1 when $x_2 = 0$.

The form of Equation 17.4 is illustrated in the right panel of Figure 17.8. It shows that the interaction can be expressed as the slope in the x_2 direction changing when x_1 changes. This is exactly analogous to the middle panel of Figure 17.8, it is important to realize, and visualize, that the interaction can be expressed instead of the slopes on either predictor.

Great care must be taken when interpreting the coefficients of a model that includes interaction terms (Braumoeller, 2004). In particular, lower-order terms are especially difficult to interpret when higher-order interactions are present. In the simple two-predictor case, the coefficient β_1 describes the influence of predictor x_1 only at $x_2 = 0$, because the slope on x_1 is $\beta_1 + \beta_{12} x_2$, as was shown in Equation 17.3 and graphed in the middle panel of Figure 17.8. In other words, it is not appropriate to say that β_1 indicates the overall influence of x_1 on y . Indeed, in many applications, the value of β_1 never realistically gets close to

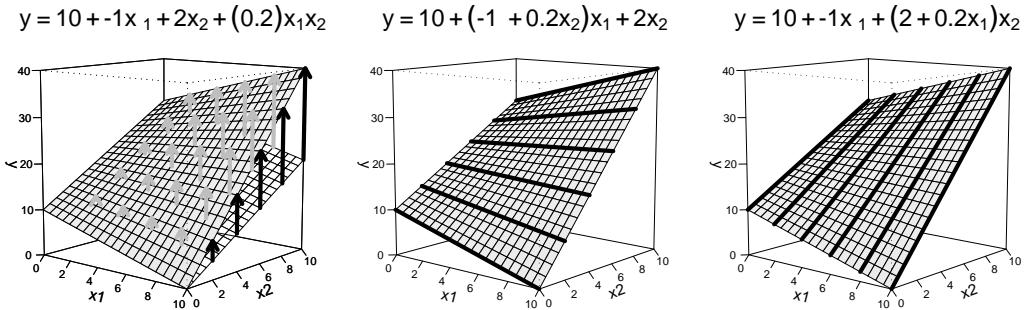


Figure 17.8: A multiplicative interaction of x_1 and x_2 , parsed three ways. The left panel emphasizes that the interaction involves a ~~planar~~ component that adds a vertical amount to the planar additive model, as indicated by the arrows. The middle panel shows the same function, but with the terms algebraically regrouped to emphasize that the slope in the x_1 direction depends on the value x_2 , as shown by the darkened lines. The right panel again shows the same function, but with the terms algebraically regrouped to emphasize that the slope in the x_2 direction depends on the value x_1 , as shown by the darkened lines.

zero, and therefore x_1 has no realistic interpretation at all. For example, suppose we are predicting college GPA y from parental income x_1 and high school GPA x_2 . If there is interaction, then the regression coefficient, b_1 , on parental income, only indicates the slope on x_1 when x_2 (GPA) is zero. Of course, there are no GPAs of zero, and therefore b_1 itself is not very informative.

17.3.1 The hierarchical model and BUGS code

The model for regression with multiplicative interaction is the same as for linear regression but with an added term for the interaction. Because we are fitting the coefficient on the multiplication of the predictors, that coefficient must have a prior, analogous to the priors on the linear coefficients. For the case of two predictors, with an interaction, the BUGS model specification looks like this: (MultiLinRegressModelBugs.R)

```

11 model {
12   for( i in 1 : nData ) {
13     y[i] ~ dnorm( mu[i] , tau )
14     mu[i] <- b0 + b1 * x[i,1] + b2 * x[i,2] + b12 * x[i,1] * x[i,2]
15   }
16   tau ~ dgamma(.001,.001)
17   b0 ~ dnorm(0,1.0E-12)
18   b1 ~ dnorm(0,1.0E-12)
19   b2 ~ dnorm(0,1.0E-12)
20   b12 ~ dnorm(0,1.0E-12)
21 }
```

Notice that the interaction coefficient, b_{12} , has a normal prior analogous to the priors on the slope coefficients.

In models that involve many predictors, there can be many interaction terms. Just as we can put a hyperprior on the slopes, we can also put a prior on the interaction coefficients. The idea is that most interaction terms are near zero, but a few might depart

from zero. This distribution over interaction terms could be modeled as a *t*-distribution. We will not pursue examples that involve numerous interactions.

17.3.1.1 Standardizing the data and initializing the chains

We will, as before, standardize the data before entering them into the BUGS model. This helps reduce correlations in the parameters, but does not eliminate correlations. When initializing the chains, it may suffice to set the interaction coefficient(s) to zero, and start the slopes at their MLE values for a non-interactive model, especially when the interactions are small, as is often the case.

Transforming the standardized estimates back to the *algebraics* is conceptually simple but algebraically much messier when interaction terms are involved. The expression, when there are merely two predictors, turns into this unwieldy form:

$$\begin{aligned}
 z(y) &= \beta_0 + \beta_1 z(x_1) + \beta_2 z(x_2) + \beta_{12} z(x_1)z(x_2) \\
 \frac{y - m_y}{s_y} &= \beta_0 + \beta_1 \frac{x_1 - m_1}{s_1} + \beta_2 \frac{x_2 - m_2}{s_2} + \beta_{12} \frac{x_1 - m_1}{s_1} \frac{x_2 - m_2}{s_2} \\
 y &= m_y + s_y \left(\beta_0 + \beta_1 \frac{m_1}{s_1} + \beta_2 \frac{m_2}{s_2} + \beta_{12} \frac{m_1 m_2}{s_1 s_2} \right) \\
 &\quad + \beta_1 s_y \frac{1}{s_1} \left(\frac{m_1}{s_1} \frac{1}{s_2} \right)^0 x_1 \\
 &\quad + \beta_2 s_y \frac{1}{s_2} \left(\frac{1}{s_1} \frac{m_2}{s_2} \right)^1 x_2 \\
 &\quad + \beta_{12} s_y \frac{1}{s_1 s_2} \left(\frac{1}{s_1} \frac{1}{s_2} \right)^2 x_1 x_2
 \end{aligned} \tag{17.5}$$

When there are more predictors involved, with their interactions, the expression becomes quite protracted. In those situations, there is no avoiding matrix notation and matrix algebra, which greatly facilitates manipulating the forms. We will venture into matrix expressions for the GLM in this book (although we did splurge a bit on matrices in Section 7.1.5, p. 104). Matrix operations are not easily expressed in BUGS, unfortunately.

17.3.2 Interpreting the posterior

To illustrate some of the issues involved in interpreting the parameters of a model with interaction, consider again the SAT data from Figure 17.8. Recall that the mean SAT score in a state was predicted from the spending per pupil and the percentage of students who took the test. When no interaction term was included in the model, the posterior looked like Figure 17.5, which indicated a positive influence of spending and a negative influence of percentage of students.

Would we want to include an interaction term in the model? The meaning of interaction is that the effect of one predictor changes when the level of the other predictor changes.

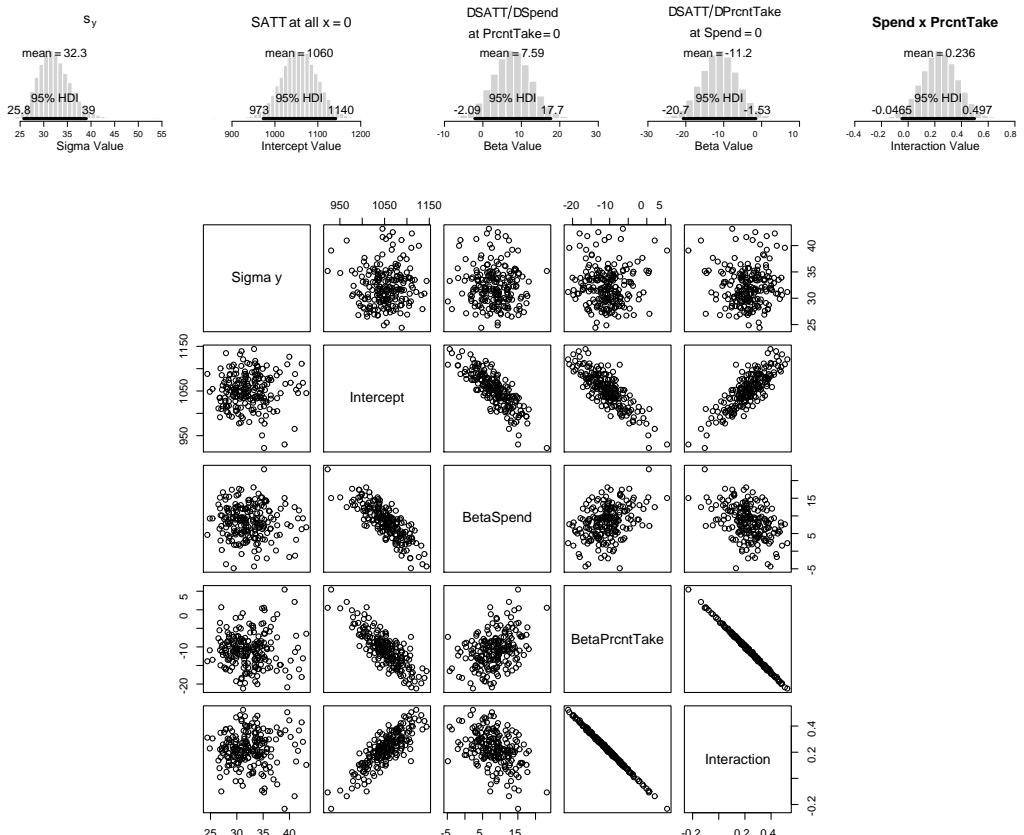


Figure 17.9: Posterior for SAT data in Figure 17.3. Compared to the posterior, which includes an interaction term, with the posterior that excludes the median interaction, shown in Figure 17.5.

Does it make sense in this application that the effect of spending would depend on the percentage of students taking the test? Perhaps yes, because few students are taking the test, they are probably already at the top of the class and therefore might not have as much head-room for increasing their scores. In other words, it is plausible that the effect of spending is larger when the percentage of students taking the test is larger, and we would not be surprised if there were a positive interaction between these predictors. Therefore, it is theoretically meaningful to include an interaction term in the model.

When we incorporate a multiplicative interaction into the model, the posterior looks like Figure 17.9. The top right histogram indicates that the credible interval for the interaction coefficient is indeed positive, as we anticipated it could be. The 95% HDI includes zero, however, which indicates that we do not have very strong evidence in our estimate of the magnitude of the interaction. The scatterplots in the lower part of Figure 17.9 indicate that the interaction coefficient is very strongly correlated with the beta coefficient on percentage of students taking the test. If the MCMC sampling were taking place in these original scales, it would be very inefficient.

A cursory look at the middle histogram of Figure 17.9 might lead a person, inappropriately, to conclude that there is not a credible influence of spending on SAT scores, because zero is among the credible values of the spend. This conclusion is inappropriate because spend

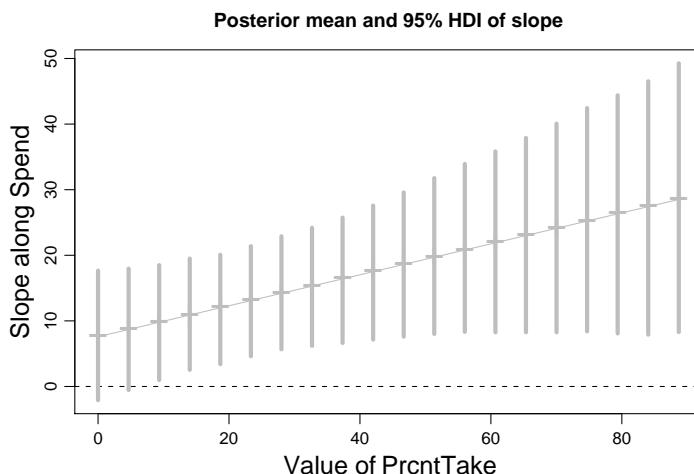


Figure 17.10: The slope in the x_1 direction is $\beta_0 + \beta_1 x_1 + \beta_2 x_2$. Shown here are the means and 95% HDIs for the slope in the x_1 direction, for various values of x_2 .

only indicates the slope on spending when the percentage of students taking the test is zero. The slope changes when the percentage of students changes. Because the interaction tends to be positive, the spending increases when the percentage of students taking the test increases.

Figure 17.10 shows the increase in slope on spending as the percentage taking the test increases. Also plotted is the extent of the 95% HDI of the estimated slope. Notice that the HDI at $PrcntTake = 0$ matches the HDI shown in the middle histogram of Figure 17.9.

Interestingly and importantly, the extent of the HDI is not constant, but also depends on the percentage of students taking the test. Mathematically, the change in extent of the HDI stems for two sources. First, the slope along x_1 is $\beta_0 + \beta_1 + \beta_2 x_2$, which means that the uncertainty in β_2 is being multiplied by x_2 , and therefore the uncertainty in the slope depends on x_2 . But the uncertainty in the slope does not necessarily always increase when x_2 increases because β_1 and β_2 are negatively correlated (see the scatterplot in Figure 17.9): When x_2 has a modest size, its negative correlation removes a bit of uncertainty in the slope. This relationship can be seen more clearly in Figure 17.11, which shows an idealized subset of credible values for β_0 and β_1 . The left panel shows that when $x_2 = 0$, the credible parameters (i.e., the dots) span a range of slopes from 0 to 15. The middle panel shows that when $x_2 = 25$, the parameters span a range of x_1 slopes from 10 to 15, i.e., a much smaller range. The right panel shows that when $x_2 = 50$, the parameters span a range of x_1 slopes from 10 to 25, again a larger range. Thus, because of the negative correlation of credible values of β_0 and β_1 , the narrowest range of x_1 slopes is at an intermediate value of x_2 .

In summary, when there is interaction, then the influence of individual predictors cannot be summarized by their individual regression coefficients alone, because those coefficients only describe the influence when the other variables are zero. A careful analyst considers credible slopes across a variety of values of the other predictors, as in Figure 17.10. Notice that this is true even when the interaction coefficient did not exclude zero from its 95% HDI: Even though the estimate of the interaction was not very precise,

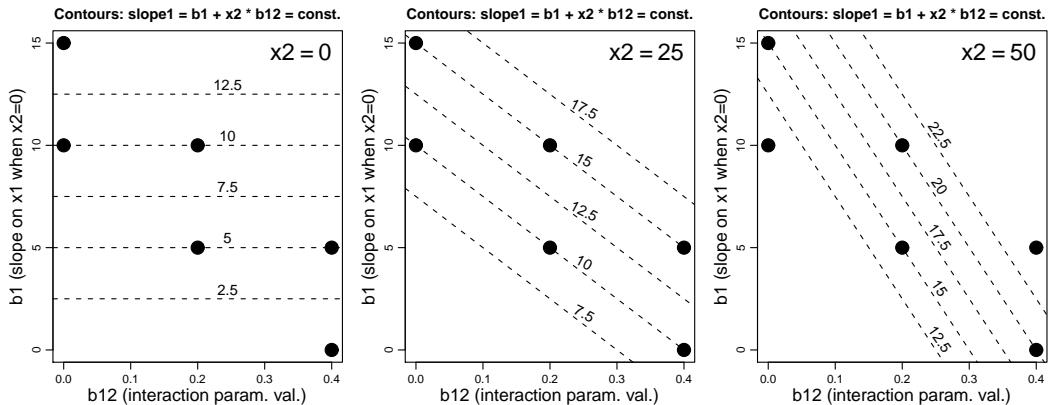


Figure 17.11: Dots show idealized parameter-value combinations, the same in all three panels. Contour lines mark parameter value combinations for slopes in the x_1 direction when x_2 has that value indicated in the panel.

the interaction did have a considerable influence on the interpretation of the predictors.

17.4 Which predictors should be included?

In many research projects, there are specific theoretical justifications for measuring particular predictors and the predicted variable. In these situations we know in advance which predictors and interactions we are interested in modeling. But in other situations, we come to a pre-existing database that has many variables measured for other purposes, and we are curious to know which of the variables might be good predictors of a particular variable of interest. Or, we might have some hybrid situation, in which we have particular predictors that we think may be relevant, but we measure several other variables on the chance that they might be relevant, or because they are needed for some study, or because they are trivial to measure so we go ahead and measure them.

When there are many potential predictors, which ones should be included in our model? And which interactions? A reasonable answer is to include those predictors and interactions that you think would have any chance of providing useful predictive information. If you do not include them, then you have essentially set the prior for those variable's regression coefficients to zero with complete certainty.

Whether or not to include candidate predictors depends on the purpose of the regression analysis. If the goal is to predict the outcome as accurately as possible, using any predictors at all, regardless of how those predictors might be causally related to the predicted variable, then all reasonable predictors could be included. On the other hand, if the purpose of the analysis is to explain the outcome on the basis of the predictors, then only the ones that can be meaningfully related to the outcome should be included (see, e.g., Keith, 2005, p. 70).

There are several costs of including a lot of candidate predictors. One cost is unwieldy interpretation of the results. With too many predictors and interactions, the complexity of the mathematical description may provide little meaningful insight into the data. Unless you really believe that the effect of income on happiness depends on shoe size, and that such interaction of income with shoe size would be theoretically meaningful, don't include shoe

size among the predictors. But, if you think there is some ~~chance~~ that a variable would be informative, include it.

Another cost of including numerous predictors is that “false alarms” may occur more often, such that predictors that really have no predictive value spuriously appear to have credibly non-zero regression coefficients. This problem can be addressed by using prior knowledge about relations among regression coefficients, as described in Section 17.2. By letting the predictors inform a hyperprior, they constrain each other’s estimates, thereby producing shrinkage and attenuation of false alarms. If you have prior knowledge about the particular predictors, that your skeptical audience would tend to, then try to express it in the mathematical form of the prior.

Another cost of including numerous candidate predictors is the noise and parameter instability introduced by the extra predictors can produce loss of precision in the estimation of the coefficients. Even when the true regression coefficients on additional predictors are all zero, the uncertainty in their values can introduce uncertainty in the estimates of the original predictors. This bleeding of uncertainty is ~~never~~ always large, especially when the predictors are uncorrelated and there are many data points. You can experiment with this issue by using the random data generator `mpptbgram` of Section 17.5.1 (`MultipleLinearRegressionBrugs.R`), and including different numbers of predictors with zero coefficients.

All of these issues are distinct from the peril of correlated predictors discussed at the beginning of the chapter. To reiterate, a predictor may ~~appear~~ have a particular relation with the predicted variable, when other candidate predictors are left out of the analysis. But the apparent relation between the predictor and the outcome ~~might~~ be an illusory artifact, produced instead by the predictor being correlated with other confounding factors, while the putative predictor itself has zero or opposite influence on the outcome. Therefore it is important to include all reasonable predictors in the analysis, so that each has a better opportunity for being correctly interpreted, while also taking in mind the costs mentioned in preceding paragraphs.

When considering the inclusion of interaction terms, and the goal of the analysis is explanation, then the main criterion is whether it is theoretically meaningful that the effect of one predictor should depend on the level of another predictor. Omission of an interaction term can also cause loss of precision in the estimates of other-order terms. Moreover, interpretation of interactions and their lower-order terms can be subtle, as we saw, for example, in Figure 17.10. If the goal of the analysis is prediction, without emphasis on explanation, then interaction terms may be included to the extent that they enhance predictability without loss of parsimony. Bayesian model comparison can be useful in this case (see references cited at the end of this section).

Whenever an interaction term is included in a model, it is important to also include all lower-order terms. For example, if an $x_i \times x_j$ interaction is included, then both x_i and x_j should also be included in the model. Although we did not discuss them, it is also possible to include three-way interactions such as $x_i \times x_j \times x_k$, if it is theoretically meaningful to do so. When this is done, it is important to include all the lower-order interactions and single predictors, including x_i , $x_i \times x_k$, $x_j \times x_k$, x_i , x_j , and x_k . When the lower-order terms are omitted, this is artificially setting their regression coefficients to zero, and distorting the posterior estimates on the other terms. For clear discussions and examples of this issue, see Braumoeller (2004) and Brambor, Clark, and Golder (2006).

In some situations, it may be theoretically tenable to suppose that the candidate predictors have exactly zero influence on the variable to be predicted, and our goal is to iden-

tify which predictors have exactly zero influence. In this section we can establish a model-comparison framework in which different models have different regression coefficients fixed at zero. A large-scale Bayesian model comparison then ~~searches~~ ~~for~~ ~~which~~ combinations of zero coefficients are most credible. There are numerous variations on this approach, and it is an active area of research (e.g., Casella & Moreno, 2006; George, 2004; E. I. George, 2000; Greenland, 2008; Liang, Paulo, Molina, Clyde, & Berger, 2008; Scott & Berger, 2006).

17.5 R code

17.5.1 Multiple linear regression

This program was used to Figure 17.5, among others. Its ~~data~~ includes three different data sets. The rest of the program is designed to be ~~general~~ applicable, so that the user can substitute other data sets without modifying ~~the~~ remainder of the program.

The use of the BUGS function `inprod()` allows the model specification to remain unchanged when the number of predictors changes. Unfortunately `inprod()` is processed slowly by BUGS. If you adapt this program for a large application and it is running too slowly, try changing `inprod(b[],x[i,])` to an explicit sum of products `b[1]*x[i,1] + b[2]*x[i,2] + ...`.

(`MultipleLinearRegressionBrugs.R`)

```

1 graphics.off()
2 rm(list=ls(all=TRUE))
3 fname = "MultipleLinearRegressionBrugs"
4 library(BRugs)      # Kruschke, J. K. (2010). Doing Bayesian data analysis:
5                      # A Tutorial with R and BUGS. Academic Press / Elsevier.
6 #-----
7 # THE MODEL.
8
9 modelstring = "
10 # BUGS model specification begins here...
11 model {
12     for( i in 1 : nData ) {
13         y[i] ~ dnorm( mu[i] , tau )
14         mu[i] <- b0 + inprod( b[] , x[i,] )
15     }
16     tau ~ dgamma(.01,.01)
17     b0 ~ dnorm(0,1.0E-12)
18     for ( j in 1:nPredictors ) {
19         b[j] ~ dnorm(0,1.0E-12)
20     }
21 }
22 # ... end BUGS model specification
23 "# close quote for modelstring
24 writeLines(modelstring,con="model.txt")
25 modelCheck( "model.txt" )
26
27 #-----
28 # THE DATA.
29
30 dataSource = c("Guber1999","McIntyre1994","random")[1      ]
31
32 if ( dataSource=="Guber1999" ) {
```

```

33 fname = "Guber1999" # file name for saved graphs
34 dataMat = read.table( file="Guber1999data.txt" ,
35                         col.names = c( "State","Spend","StuTchRat","Salary",
36                                         "PrcntTake","SATV","SATM","SATT" ) )
37 # Specify variables to be used in BUGS analysis:
38 predictedName = "SATT"
39 predictorNames = c( "Spend" , "PrcntTake" )
40 #predictorNames = c( "Spend" , "PrcntTake" , "Salary" , "Stu      TchRat" )
41 nData = NROW( dataMat )
42 y = as.matrix( dataMat[,predictedName] )
43 x = as.matrix( dataMat[,predictorNames] )
44 nPredictors = NCOL( x )
45 }
46
47 if ( dataSource=="McIntyre1994" ) {
48   fname = "McIntyre1994" # file name for saved graphs
49   dataMat = read.csv(file="McIntyre1994data.csv")
50   predictedName = "CO"
51   predictorNames = c("Tar","Nic","Wt")
52   nData = NROW( dataMat )
53   y = as.matrix( dataMat[,predictedName] )
54   x = as.matrix( dataMat[,predictorNames] )
55   nPredictors = NCOL( x )
56 }
57
58 if ( dataSource=="random" ) {
59   fname = "Random" # file name for saved graphs
60   # Generate random data.
61   # True parameter values:
62   betaTrue = c( 100 , 1 , 2 , rep(0,21) ) # beta0 is first component
63   nPredictors = length( betaTrue ) - 1
64   sdTrue = 2
65   tauTrue = 1/sdTrue^2
66   # Random X values:
67   set.seed(47405)
68   xM = 5 ; xSD = 2
69   nData = 100
70   x = matrix( rnorm( nPredictors*nData , xM , xSD ) , nrow=nData      )
71   predictorNames = colnames(x) = paste("X",1:nPredictors,      sep="")
72   # Random Y values generated from linear model with true parameter values:
73   y = x %*% matrix(betaTrue[-1],ncol=1) + betaTrue[1] + rmnorm( nData,0,sdTrue )
74   predictedName = "Y"
75   # Select which predictors to include
76   includeOnly = 1:nPredictors # default is to include all
77   #includeOnly = 1:10 # subset of predictors overwrites default
78   x = x[,includeOnly]
79   predictorNames = predictorNames[includeOnly]
80   nPredictors = NCOL(x)
81 }
82
83 # Prepare data for BUGS:
84 # Re-center data at mean, to reduce autocorrelation in MCMC sampling.
85 # Standardize (divide by SD) to make prior specification easier.
86 standardizeCols = function( dataMat ) {
87   zDataMat = dataMat
88   for ( colIdx in 1:NCOL( dataMat ) ) {
89     mCol = mean( dataMat[,colIdx] )
90     sdCol = sd( dataMat[,colIdx] )
91     zDataMat[,colIdx] = ( dataMat[,colIdx] - mCol ) / sdCol

```

```

92      }
93      return( zDataMat )
94  }
95 zx = standardizeCols( x )
96 zy = standardizeCols( y )
97
98 # Get the data into BUGS:
99 dataList = list(
100   x = zx ,
101   y = as.vector( zy ) , # BUGS does not treat 1-column mat as vecto r
102   nPredictors = nPredictors ,
103   nData = nData
104 )
105 modelData( bugsData( dataList ) )
106
107 #-----
108 # INTIALIZE THE CHAINS.
109
110 nChain = 3
111 modelCompile( numChains = nChain )
112
113 genInitList <- function(nPred=nPredictors) {
114   lmInfo = lm( dataList$y ~ dataList$x ) # R function returns ML      E
115   bInit = lmInfo$coeff[-1]
116   tauInit = length(dataList$y) / sum(lmInfo$res^2)
117   list(
118     b0 = 0 ,
119     b = bInit ,
120     tau = tauInit
121   )
122 }
123 for ( chainIdx in 1 : nChain ) {
124   modelInits( bugsInits( genInitList ) )
125 }
126
127 #-----
128 # RUN THE CHAINS
129
130 # burn in
131 BurnInSteps = 100
132 modelUpdate( BurnInSteps )
133 # actual samples
134 samplesSet( c( "b0" , "b" , "tau" ) )
135 stepsPerChain = ceiling(10000/nChain)
136 thinStep = 2
137 modelUpdate( stepsPerChain , thin=thinStep )
138
139 #-----
140 # EXAMINE THE RESULTS
141
142 source("plotChains.R")
143 source("plotPost.R")
144
145 checkConvergence = F
146 if ( checkConvergence ) {
147   b0Sum = plotChains( "b0" , saveplots=F , filenameroot=fnam    e )
148   bSum = plotChains( "b" , saveplots=F , filenameroot=fname )
149   tauSum = plotChains( "tau" , saveplots=F , filenameroot=fn      ame )
150 }
```

```

151 # Extract chain values:
152 zb0Samp = matrix( samplesSample( "b0" ) )
153 zbSamp = NULL
154 for ( j in 1:nPredictors ) {
155   zbSamp = cbind( zbSamp , samplesSample( paste("b[,j,]",      sep="" ) ) )
156 }
157
158 zTauSamp = matrix( samplesSample( "tau" ) )
159 zSigmaSamp = 1 / sqrt( zTauSamp ) # Convert precision to SD
160 chainLength = length(zTauSamp)
161
162 # Convert to original scale:
163 bSamp = zbSamp * matrix( sd(y)/apply(x,2,sd) , byrow=TRUE ,
164                           ncol=nPredictors , nrow=NROW(zbSamp) )
165 b0Samp = ( zb0Samp * sd(y)
166             + mean(y)
167             - rowSums( zbSamp
168                         * matrix( sd(y)/apply(x,2,sd) , byrow=TRUE ,
169                                       ncol=nPredictors , nrow=NROW(zbSamp) )
170                         * matrix( apply(x,2,mean) , byrow=TRUE ,
171                                       ncol=nPredictors , nrow=NROW(zbSamp) ) ) )
171
172 sigmaSamp = zSigmaSamp * sd(y)
173
174 # Save MCMC sample:
175 save( b0Samp , bSamp , sigmaSamp ,
176       file="MultipleLinearRegressionBrugsGuber1999.Rdata"      )
177
178 # Scatter plots of parameter values, pairwise:
179 if ( nPredictors <= 6 ) { # don't display if too many predictors
180   windows()
181   thinIdx = round(seq(1,length(zb0Samp),length=200))
182   pairs( cbind( zSigmaSamp[thinIdx] , zb0Samp[thinIdx] , zb      Samp[thinIdx,] ) ,
183         labels=c("Sigma zy","zIntercept",paste("zSlope",pred      ictorNames,sep="")))
184   windows()
185   thinIdx = seq(1,length(b0Samp),length=700)
186   pairs( cbind( sigmaSamp[thinIdx] , b0Samp[thinIdx] , bSam      p[thinIdx,] ) ,
187         labels=c( "Sigma y" , "Intercept" , paste("Slope",predict      orNames,sep="")))
188   dev.copy2eps(file=paste(fname,"PostPairs.eps",sep=""))
189 }
190 # Show correlation matrix on console:
191 cat("\nCorrelations of posterior sigma, b0, and bs:\n")
192 show( cor( cbind( sigmaSamp , b0Samp , bSamp ) ) )
193
194 # Display the posterior:
195 nPlotPerRow = 5
196 nPlotRow = ceiling((2+nPredictors)/nPlotPerRow)
197 nPlotCol = ceiling((2+nPredictors)/nPlotRow)
198 windows(3.5*nPlotCol,2.25*nPlotRow)
199 layout( matrix(1:(nPlotRow*nPlotCol),nrow=nPlotRow,n      col=nPlotCol,byrow=T) )
200 par( mar=c(4,3,2.5,0) , mgp=c(2,0.7,0) )
201 histInfo = plotPost( sigmaSamp , xlab="Sigma Value" , compV      al=NULL ,
202                      breaks=30 , main=bquote(sigma[y]) ,
203                      cex.main=1.67 , cex.lab=1.33 )
204 histInfo = plotPost( b0Samp , xlab="Intercept Value" , comp      Val=NULL ,
205                      breaks=30 , main=bquote(.predictedName) ** at ** x==0) ,
206                      cex.main=1.67 , cex.lab=1.33 )
207 for ( sIdx in 1:nPredictors ) {
208   histInfo = plotPost( bSamp[,sIdx] , xlab="Slope Value" , co      mpVal=0.0 ,
209                      breaks=30 ,

```

```

210           main=bquote( Delta * .(predictedName) /
211                           Delta * .(predictorNames[sIdx]) ) ,
212                           cex.main=1.67 , cex.lab=1.33 )
213     }
214 dev.copy2eps(file=paste(fname,"PostHist.eps",sep=""))
215
216 # Posterior prediction:
217 # Specify x values for which predicted y's are needed.
218 # xPostPred is a matrix such that ncol=nPredictors and nrow= nPostPredPts.
219 xPostPred = rbind(
220   apply(x,2,mean)-3*apply(x,2,sd) , # mean of data x minus th rice SD of data x
221   apply(x,2,mean) , # mean of data x
222   apply(x,2,mean)+3*apply(x,2,sd) # mean of data x plus thri ce SD of data x
223 )
224 # Define matrix for recording posterior predicted y values f or each xPostPred.
225 # One row per xPostPred value, with each row holding random pr edicted y values.
226 postSampSize = chainLength
227 yPostPred = matrix( 0 , nrow=NROW(xPostPred) , ncol=postSa mpSize )
228 # Define matrix for recording HDI limits of posterior predic ted y values:
229 yHDIlim = matrix( 0 , nrow=NROW(xPostPred) , ncol=2 )
230 # Generate posterior predicted y values.
231 # This gets only one y value, at each x, for each step in the chai n.
232 for ( chainIdx in 1:chainLength ) {
233   yPostPred[,chainIdx] = rnorm( NROW(xPostPred) ,
234                               mean = b0Samp[chainIdx]
235                               + xPostPred %*% cbind(bSamp[chainIdx],) ,
236                               sd = rep( sigmaSamp[chainIdx] , NROW(xPostPred) ) )
237 }
238 source("HDIofMCMC.R")
239 for ( xIdx in 1:NROW(xPostPred) ) {
240   yHDIlim[xIdx,] = HDIofMCMC( yPostPred[xIdx,] )
241 }
242 cat( "\nPosterior predicted y for selected x:\n" )
243 show( cbind( xPostPred , yPostPredMean=rowMeans(yPostPr ed) , yHDIlim ) )
244
245 #-----
```

17.5.2 Multiple linear regression with hyperprior on coefficients

This program was used to create Figure 17.7, among others. See comments regarding `inprod()` before the previous program.

(MultiLinRegressHyperBrugs.R)

```

1 graphics.off()
2 rm(list=ls(all=TRUE))
3 fname = "MultiLinRegressHyper"
4 library(BRugs)      # Kruschke, J. K. (2010). Doing Bayesian dat a analysis:
5                      # A Tutorial with R and BUGS. Academic Press / Elsevier.
6 #-----
```

```

7 # THE MODEL.
8
9 modelstring = "
10 # BUGS model specification begins here...
11 model {
12   for( i in 1 : nData ) {
13     y[i] ~ dnorm( mu[i] , tau )
14     mu[i] <- b0 + inprod( b[] , x[i,] )
15 }
```

```

16     tau ~ dgamma(.01,.01)
17     b0 ~ dnorm(0,1.0E-12)
18     for ( j in 1:nPredictors ) {
19         b[j] ~ dt( muB , tauB , tdfB )
20     }
21     muB ~ dnorm( 0 , .100 )
22     tauB ~ dgamma(.01,.01)
23     udfB ~ dunif(0,1)
24     tdfB <- 1 + tdfBgain * ( -log( 1 - udfB ) )
25 }
26 # ... end BUGS model specification
27 " # close quote for modelstring
28 writeLines(modelstring,con="model.txt")
29 modelCheck( "model.txt" )
30
31 #-----#
32 # THE DATA.
33
34 tdfBgain = 1
35
36 dataSource = c("Guber1999","McIntyre1994","random")[3      ]
37
38 if ( dataSource=="Guber1999" ) {
39     fname = paste("Guber1999Brugs","tdf",tdfBgain,sep="")
40     dataMat = read.table( file="Guber1999data.txt" ,
41                           col.names = c( "State","Spend","StuTchRat","Salary",
42                                         "PrcntTake","SATV","SATM","SATT" ) )
43     # Specify variables to be used in BUGS analysis:
44     predictedName = "SATT"
45     predictorNames = c( "Spend" , "PrcntTake" )
46     #predictorNames = c( "Spend" , "PrcntTake" , "Salary" , "Stu      TchRat" )
47     nData = NROW( dataMat )
48     y = as.matrix( dataMat[,predictedName] )
49     x = as.matrix( dataMat[,predictorNames] )
50     nPredictors = NCOL( x )
51 }
52
53 if ( dataSource=="McIntyre1994Hyper" ) {
54     fname = paste("McIntyre1994Brugs","tdf",tdfBgain,sep=      "")
55     dataMat = read.csv(file="McIntyre1994data.csv")
56     predictedName = "CO"
57     predictorNames = c("Tar","Nic","Wt")
58     nData = NROW( dataMat )
59     y = as.matrix( dataMat[,predictedName] )
60     x = as.matrix( dataMat[,predictorNames] )
61     nPredictors = NCOL( x )
62 }
63
64 if ( dataSource=="random" ) {
65     fname = paste("RandomHyper","tdf",tdfBgain,sep="")
66     # Generate random data.
67     # True parameter values:
68     betaTrue = c( 100 , 1 , 2 , rep(0,21) ) # beta0 is first component
69     nPredictors = length( betaTrue ) - 1
70     sdTrue = 2
71     tauTrue = 1/sdTrue^2
72     # Random X values:
73     set.seed(47405)
74     xM = 5 ; xSD = 2

```

```

75  nData = 100
76  x = matrix( rnorm( nPredictors*nData , xM , xSD ) , nrow=nData      )
77  predictorNames = colnames(x) = paste("X",1:nPredictors,      sep="")
78  # Random Y values generated from linear model with true parameter values:
79  y = x %*% matrix(betaTrue[-1],ncol=1) + betaTrue[1] + rnorm      (nData,0,sdTrue)
80  predictedName = "Y"
81  # Select which predictors to include
82  includeOnly = 1:nPredictors # default is to include all
83  #includeOnly = 1:6 # subset of predictors overwrites default
84  x = x[,includeOnly]
85  predictorNames = predictorNames[includeOnly]
86  nPredictors = NCOL(x)
87 }
88
89 # Prepare data for BUGS:
90 # Re-center data at mean, to reduce autocorrelation in MCMC sampling.
91 # Standardize (divide by SD) to make initialization easier.
92 standardizeCols = function( dataMat ) {
93   zDataMat = dataMat
94   for ( colIdx in 1:NCOL( dataMat ) ) {
95     mCol = mean( dataMat[,colIdx] )
96     sdCol = sd( dataMat[,colIdx] )
97     zDataMat[,colIdx] = ( dataMat[,colIdx] - mCol ) / sdCol
98   }
99   return( zDataMat )
100 }
101 zx = standardizeCols( x )
102 zy = standardizeCols( y )
103
104 # Get the data into BUGS:
105 dataList = list(
106   x = zx ,
107   y = as.vector( zy ) , # BUGS does not treat 1-column mat as vector
108   nPredictors = nPredictors ,
109   nData = nData ,
110   tdfBgain = tdfBgain
111 )
112 modelData( bugsData( dataList ) )
113
114 #-----#
115 # INITIALIZE THE CHAINS.
116
117 nChain = 3
118 modelCompile( numChains = nChain )
119
120 genInitList <- function(nPred=nPredictors) {
121   lmInfo = lm( y ~ x ) # R function returns least-squares (normal      MLE) fit.
122   bInit = lmInfo$coeff[-1] * apply(x,2,sd) / sd(y)
123   tauInit = (length(y)*sd(y)^2)/sum(lmInfo$res^2)
124   list(
125     b0 = 0 ,
126     b = bInit ,
127     tau = tauInit ,
128     muB = mean( bInit ) ,
129     tauB = 1 / sd( bInit )^2 ,
130     udfB = 0.95 # tdfB = 4
131   )
132 }
133 for ( chainIdx in 1 : nChain ) {

```

```

134     modellInits( bugsInits( genInitList ) )
135 }
136
137 #-----
138 # RUN THE CHAINS
139
140 # burn in
141 BurnInSteps = 100
142 modelUpdate( BurnInSteps )
143 # actual samples
144 samplesSet( c( "b0" , "b" , "tau" , "muB" , "tauB" , "tdfB" ) )
145 stepsPerChain = ceiling(10000/nChain)
146 thinStep = 2
147 modelUpdate( stepsPerChain , thin=thinStep )
148
149 #-----
150 # EXAMINE THE RESULTS
151
152 source("plotChains.R")
153 source("plotPost.R")
154
155 checkConvergence = F
156 if ( checkConvergence ) {
157   b0Sum = plotChains( "b0" , saveplots=F , filenameroot=fnam e )
158   bSum = plotChains( "b" , saveplots=F , filenameroot=fnme )
159   tauSum = plotChains( "tau" , saveplots=F , filenameroot=fna me )
160   muBSum = plotChains( "muB" , saveplots=F , filenameroot=fna me )
161   tauBSum = plotChains( "tauB" , saveplots=F , filenameroot=fna me )
162   tdfBSum = plotChains( "tdfB" , saveplots=F , filenameroot=fna me )
163 }
164
165 # Extract chain values:
166 zb0Samp = matrix( samplesSample( "b0" ) )
167 zbSamp = NULL
168 for ( j in 1:nPredictors ) {
169   zbSamp = cbind( zbSamp , samplesSample( paste("b[,j]", sep="" ) ) )
170 }
171 zTauSamp = matrix( samplesSample( "tau" ) )
172 zSigmaSamp = 1 / sqrt( zTauSamp ) # Convert precision to SD
173 chainLength = length(zTauSamp)
174
175 # Convert to original scale:
176 bSamp = zbSamp * matrix( sd(y)/apply(x,2,sd) , byrow=TRUE ,
177                           ncol=nPredictors , nrow=NROW(zbSamp) )
178 b0Samp = ( zb0Samp * sd(y)
179             + mean(y)
180             - rowSums( zbSamp
181                         * matrix( sd(y)/apply(x,2,sd) , byrow=TRUE ,
182                                       ncol=nPredictors , nrow=NROW(zbSamp) )
183                         * matrix( apply(x,2,mean) , byrow=TRUE ,
184                                       ncol=nPredictors , nrow=NROW(zbSamp) ) ) )
185 sigmaSamp = zSigmaSamp * sd(y)
186
187 # Scatter plots of parameter values, pairwise:
188 if ( nPredictors <= 6 ) { # don't display if too many predictors
189   windows()
190   thinIdx = round(seq(1,length(zb0Samp),length=200))
191   pairs( cbind( zSigmaSamp[thinIdx] , zb0Samp[thinIdx] , zb      Samp[thinIdx,] ) ,
192         labels=c("Sigma zy","zIntercept",paste("zSlope",pred      ictorNames,sep="")))

```

```

193 windows()
194 thinIdx = seq(1,length(b0Samp),length=700)
195 pairs( cbind( sigmaSamp[thinIdx] , b0Samp[thinIdx] , bSam
196   labels=c( "Sigma y" , "Intercept", paste("Slope",predict
197 dev.copy2eps(file=paste(fname,"PostPairs.eps",sep=""))
198 }
199 # Show correlation matrix on console:
200 cat("\nCorrelations of posterior sigma, b0, and bs:\n")
201 show( cor( cbind( sigmaSamp , b0Samp , bSam ) ) )
202
203 # Display the posterior:
204 nPlotPerRow = 5
205 nPlotRow = ceiling((2+nPredictors)/nPlotPerRow)
206 nPlotCol = ceiling((2+nPredictors)/nPlotRow)
207 windows(3.5*nPlotCol,2.25*nPlotRow)
208 layout( matrix(1:(nPlotRow*nPlotCol),nrow=nPlotRow,n      col=nPlotCol,byrow=T) )
209 par( mar=c(4,3,2.5,0) , mgp=c(2,0.7,0) )
210 histInfo = plotPost( sigmaSamp , xlab="Sigma Value" , compV    al=NULL ,
211   breaks=30 , main=bquote(sigma[y]) ,
212   cex.main=1.67 , cex.lab=1.33 )
213 histInfo = plotPost( b0Samp , xlab="Intercept Value" , comp      Val=NULL ,
214   breaks=30 , main=bquote(.predictedName) ** at "* x==0) ,
215   cex.main=1.67 , cex.lab=1.33 )
216 for ( sIdx in 1:nPredictors ) {
217 histInfo = plotPost( bSam[sIdx] , xlab="Slope Value" , co      mpVal=0.0 ,
218   breaks=30 ,
219   main=bquote( Delta * .predictedName) /
220   Delta * .(predictorNames[sIdx]) ) ,
221   cex.main=1.67 , cex.lab=1.33 )
222 }
223 dev.copy2eps(file=paste(fname,"PostHist.eps",sep=""))
224
225 # Posterior prediction:
226 # Specify x values for which predicted y's are needed.
227 # xPostPred is a matrix such that ncol=nPredictors and nrow=      nPostPredPts.
228 xPostPred = rbind(
229   apply(x,2,mean)-3*apply(x,2,sd) , # mean of data x minus th      rice SD of data x
230   apply(x,2,mean)                  , # mean of data x
231   apply(x,2,mean)+3*apply(x,2,sd)   # mean of data x plus thri      ce SD of data x
232 )
233 # Define matrix for recording posterior predicted y values f      or each xPostPred.
234 # One row per xPostPred value, with each row holding random pr edicted y values.
235 postSampSize = chainLength
236 yPostPred = matrix( 0 , nrow=NROW(xPostPred) , ncol=postSa mpSize )
237 # Define matrix for recording HDI limits of posterior predic      ted y values:
238 yHDIlim = matrix( 0 , nrow=NROW(xPostPred) , ncol=2 )
239 # Generate posterior predicted y values.
240 # This gets only one y value, at each x, for each step in the chai      n.
241 for ( chainIdx in 1:chainLength ) {
242   yPostPred[,chainIdx] = rnorm( NROW(xPostPred) ,
243     mean = b0Samp[chainIdx]
244     + xPostPred %*% cbind(bSam[chainIdx,]) ,
245     sd = rep( sigmaSamp[chainIdx] , NROW(xPostPred) ) )
246 }
247 source("HDIofMCMC.R")
248 for ( xIdx in 1:NROW(xPostPred) ) {
249   yHDIlim[xIdx,] = HDIofMCMC( yPostPred[xIdx,] )
250 }
251 cat( "\nPosterior predicted y for selected x:\n" )

```

```

252 show( cbind( xPostPred , yPostPredMean=rowMeans(yPostPred) , yHDIlim ) )
253 #-----
254 -----
```

17.6 Exercises

Exercise 17.1. [Purpose: View the prior on the regression coefficients, when there is a hyperprior.] The hyperprior on regression coefficients in Figure 17.6 may be difficult to intuit. Therefore, generate graphs of prior distribution on the regression coefficients for the program in Section 17.5.2 (MultiLinRegressHyperBrugs.R), when `tdf` is set at different values. To accomplish this, do the following:

1. Because BUGS balks at extreme values of tau, and extremes can be sampled under the vague prior, the gamma distributions for tau need to be censored. Therefore, change the tau specifications to

```
tau    dgamma(.01,.01)(0.0001,10000)
tauB   dgamma(.01,.01)(0.0001,10000)
```

2. In the `data` list, comment out only the single line that specifies `the` values. The other lines, that specify `the` values and the number of predictors, etc., must remain, ~~because~~ they specify the structure of the model. The model only predicts a function of `x`; the model does not predict `y`.

3. In the initialization of the chains, comment out the `data` initialization, and instead use `modelGenInits()`.

4. Do not use the `plotChains()` command, because some of the MCMC values may be too extreme for some of the BUGS graphics routines that already `plotChains()`, and BUGS will crash.

5. When plotting the slope values in the last part of the code, manually specify the axis limits so that you can see the central parts of the distribution, perhaps like this:

```
histInfo = plotPost( bSamp[,sIdx] , xlab="Slope Value" , co      mpVal=0.0 ,
                     breaks=c(-1000000,seq(-400,400,21),1000000) , xlim=c(      -400,400) , ...)
```

Run the program and display the prior for `tdfBgain=1` and for `tdfBgain=100`. Point out and discuss any differences in the priors for those different hyperprior constants.

Exercise 17.2. [Purpose: Power analysis for multiple regression.] Consider the SAT data shown in Figure 17.3, p. 374, and the posterior for a linear regression model (with no interaction), shown in Figure 17.5, p. 377. The marginal posterior distribution for the slope on spending-per-pupil has a 95% HDI that excludes zero, and that has a width nearly 17. Consider two different goals for the research: One goal is to show that the 95% HDI spending-per-pupil excludes a ROPE of $[0; +1]$, and the second goal is to show that the width of the 95% HDI is less than 10. How can we assess the probability of meeting these goals?

To address the question, we think of `the` values as being randomly generated according to the linear regression model, at `the` values that are provided by the world. The model only describes the dependency of `y` on `x`; the model does not describe the distribution of the `x` values. In different applications `the` values have different actual generators. In experiments, `the` values are selected by the experimenter, and power analysis explicitly manipulate the values of `the` and the frequencies of each. In observational studies, `the` values are generated by the world, not by the experimenter. In observational studies, we can think of `the` values as randomly drawn from some underlying distribution. `x` could be a person's height, and we can easily randomly sample another person to get

another x value. In other observational studies, it is a conceptualist to think of x as being a random value that is sampled from some underlying population; e.g., spending-per-pupil in each of the 50 states. But even in that last case, we think of each state as being representative of some universe of possible states, from which 50 were drawn.

In the present non-experimental, observational study, (the SAT data for 50 states), the x values are not selected by researcher. Moreover, we have no model of the x values. Therefore, we use the actual values themselves as the best available description of the underlying distribution of x values. When we simulate a new datum, we first randomly select one of the points from the actual values, and then we randomly generate a value according the credible parameter values in the model.

In our specific situation there are only 50 states altogether therefore we can do a retrospective power analysis using a sample size of 50. But also imagine increasing the sample size if, instead of using state-average data, imagine that the data are from individual school districts within states. Under this interpretation, the posterior from the state-average data, and the range of values, might not be fully representative of district data, but at least it's better than nothing.

Answer these questions: What is the retrospective power of the two goals? (Hint: It's about .67 and .00, when the points are randomly sampled with replacement from the 50 actual points.) What is the power for each of the goals when $N = 130$? (Hint: It's about .85 and .30.) What's the minimal N needed to achieve a power of .80 for the second goal?

Programming hints: Use the power-analysis program of Section 13.6.2 (`FilconBrugsPower.R`) as a template for adapting the linear regression program Section 17.5.1 (`MultipleLinearRegressionBrugs.R`). You'll want to pass the raw data into the `GoalAchievedForSample()` function because that way the data can be standardized for SBU but the resulting parameter values can be converted back to original scale (using information from the original data).

Chapter 18

Metric Predicted Variable with One Nominal Predictor

Contents

18.1 Bayesian oneway ANOVA	402
18.1.1 The hierarchical prior	340
18.1.1.1 Homogeneity of variance	404
18.1.2 Doing it with R and BUGS	404
18.1.3 A worked example	406
18.1.3.1 Contrasts and complex comparisons	407
18.1.3.2 Is there a difference?	408
18.2 Multiple comparisons	409
18.3 Two group Bayesian ANOVA and the NHST test	412
18.4 R code	413
18.4.1 Bayesian oneway ANOVA	413
18.5 Exercises	417

Familywise error rates breed rumors of incest,
Hounding for quarry in multiple t tests.
Barking at research, poor dog got run over;
Should have done Bayesian oneway ANOVA.

In this chapter we consider situations with a metric predicted variable and a nominally-scaled predictor variable. There are many cases of these situations in real-world research. For example, we may want to predict weight loss (a metric variable) as a function of which diet the person follows (e.g., low-carb, vegetarian, or-fatty). As another example, we may want to predict severity of psychosis (measured on a metric scale) as a function of which anti-psychotic drug the person takes. Or we may want to predict income as a function of political party affiliation. This combination of predicted and predictor scales occurs in the first row, fourth cell, of Table 14.1 (p. 312).

In traditional NHST analyses, these situations are addressed by “oneway analysis of variance” (ANOVA). The term “oneway” refers to the fact that a single nominal variable is being used as the predictor. The phrase “analysis of variance” refers to the fact that the overall variance among all the values is decomposed, analyzed, into two parts: variance within the levels of the nominal predictors, and variance between the levels of the nominal

predictors. The variance within levels of the nominal predictor is called noise or error, i.e., variability that cannot be predicted by the predictor. The complementary variance between the levels of the nominal predictor is called the effect of the predictor. Usually we do the research with the goal of detecting an effect, which means that we would like the magnitude of the variance between levels to be large compared to the variance within levels. The ratio, of variance between to variance within, is called the F-ratio. In the Bayesian approach, we rarely if ever refer to the F-ratio. But because the model we use is based on the model of traditional ANOVA, we will refer to our analysis as Bayesian ANOVA, or sometimes BANOVA.

18.1 Bayesian oneway ANOVA

The basic idea of oneway ANOVA was introduced in Section 6.4.1 p. 300. The predictor is a variable measured on a nominal scale. For example, affiliation is predicted as a function of political party affiliation, notice that the predictor has nominal levels such as libertarian, green, democratic, republican, and so on. We denote the predictable axis, which is a vector with one component per nominal level. For example, suppose that the predictor is political party affiliation, with Green as level 1, Democrat as level 2, Republican as level 3, Libertarian as level 4, and Other as level 5. Then Democrat is represented as $x = [0; 1; 0; 0; 0]$, and Libertarian is represented as $x = [0; 0; 0; 1; 0]$. Political party affiliation is being treated here as a categorical label only, no ordering along a liberal-conservative scale.

The formal model indicates how to derive the predicted value of the predictor. The idea is that there is a baseline quantity of the predicted value, and each level of the predictor indicates a deviation above or below that baseline. We will denote the baseline value of the prediction as α_0 . The deviation for the j^{th} level of the predictor is denoted β_j .

When the predictor has value $x = [\alpha_0; x_1; \dots; x_j; \dots; x_n]$, then the predicted value is

$$\begin{aligned} X &= \alpha_0 + \sum_j \beta_j x_{ji} \\ &= \alpha_0 + \beta_j x_i \end{aligned} \quad (18.1)$$

where the notation $\beta_j x_i$ denotes the “dot product” of the vectors. In Equation 18.1, the coefficient β_j indicates how much changes when x_i changes from neutral to level j . In other words, β_j indicates how much changes when x_i changes from $alk_j = 0$ to $x_i = 1$. The baseline is constrained such that the deviations sum to zero across the levels of:

$$\sum_{j=1}^J \beta_j = 0 \quad (18.2)$$

The expression of the model in Equation 18.1 is not complete without the constraint in 18.2. Examples were shown in Figure 14.4, p. 301, and it is time to go now to that Figure for a quick review.

The predicted value \hat{y}_i in Equation 18.1 is of the central tendency in the data. The data themselves are assumed to be randomly generated around that central tendency. As usual, we will assume a normal distribution $y_i \sim N(\hat{y}_i; \sigma^2)$, where σ^2 is the precision of the normal distribution. As discussed in previous chapters, if the data have outliers, a t-distribution may be used instead.

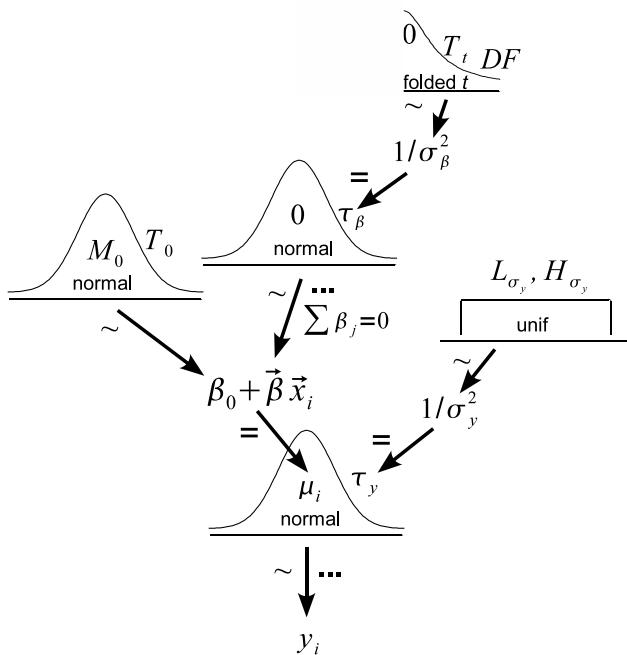


Figure 18.1: Hierarchical dependencies for model of one Bayesian ANOVA. The baseline is β_0 , and the deviation away from that baseline for the i^{th} level of x is β_i . The standard deviation of the β 's has a folded-prior. The variance within levels of x is estimated by the precision τ_y , which is here assumed to be homogeneous across groups, although it need not be in general.

18.1.1 The hierarchical prior

Our primary interest is in estimating the deviation parameters β_j for each level of x . We could just put a separate prior on each parameter and estimate them separately from each other. It is typical, however, that the levels of x are not utterly unrelated to each other, and therefore data from one level may inform estimates in other levels. For example, the deviations for republicans, libertarians, and greens ~~form~~ an estimate of the deviation for democrats. Thus, if the deviation for libertarians is 0, for republicans is 0.5, and for greens is 1.0, then the deviation for democrats should be somewhere in the general range, and not out at, say, 12.0. At least, we might have prior beliefs that the deviations for most levels of x may be small, with only a few deviations being large, and ~~there~~ we can let the various levels mutually inform each other's estimates based on this structural assumption.

The form of the hierarchical model for oneway BANOVA is displayed in Figure 18.1 (Gelman, 2005, 2006). In the upper middle of the diagram is a normal distribution that describes the distribution of deviations β_j across levels of x . This normal distribution has a mean at zero, reflecting the fact that the deviations are constrained to fall both above and below the baseline, because they must sum to zero. ~~and~~ the precision of this normal distribution, τ_β , is estimated, not pre-set at a constant. Thus, if many levels of x have a small deviation in the data, then the precision is estimated to be high, and this in turn shrinks the estimates of other

The prior for τ_β derives from the recommendation of Gelman (2006). First, the preci-

sion is converted to standard deviation:= $1 = \sigma^2$. Then a folded-t-distribution is used as a prior on σ . The folded-t is just the positive side of the usual distribution. Notice that it is defined only over non-zero values, as is required for σ and it extends to positive infinity. Unlike the gamma() distribution that is often used for precisions, however, folded-t does not have infinite density near zero. Because noisy data never rule out decisions of all zero, there can cause unintended distortions in the posterior if the prior places extreme densities at either end of the scale (see Gelman, 2006, [false](#)).

A folded-t prior could also be used for the noise, but we will use a uniform, again as recommended by Gelman (2006). One motivation is that σ may have a more intuitive form than a folded-t when expressing a prior belief. A second reason is that the infelicities of estimation that affect σ are not present so prominently at this level in the model, because the within-level noise is typically not zero, and there are enough data points to overwhelm any mildly informed prior.

18.1.1.1 Homogeneity of variance

The model described here assumes equal variances across levels of x . As a concrete example, the model assumes that the variance of income is the same for republicans as for democrats as for libertarians as for greens. This assumption of homogeneous variances is vestigial from two precursors. First, the analogous assumption is made in linear regression, and ANOVA may be construed, mathematically, as a special case of linear regression. Second, homogeneity of variance is assumed in NHST ANOVA to simplify derivation of F distributions. Neither of these precursors actually demands that we make the assumption of equal variances in BANOVA.

The model described here assumes homogeneity of variance ~~for simplicity in presentation~~. By assuming equal variances for all levels, the focus could be on the estimation of the group decision parameters. Also, by assuming equal variances, the results of BANOVA can be more directly compared to the results of NHST ANOVA, if such a comparison desired.

In principle, the BANOVA model can (and should) estimate variance components for each level of x . The model in Figure 18.1 can be expanded analogously to ~~the~~ in Figure 16.11, p. 356. Instead of a single precision used for all levels of x , a separate precision σ_j^2 is estimated for each level of x , as in the lower-right of Figure 16.11. A higher-level distribution describes the spread of the across levels of x . This structure provides shrinkage of the estimates of σ_j^2 , to the extent that the data suggest homogeneity of variance. Exercise 18.3 has you give this scheme a test drive.

18.1.2 Doing it with R and BUGS

As usual, every arrow in the hierarchical diagram of Figure 18.1 has a corresponding line in the BUGS model specification. The parameters that appear ~~as~~ in Figure 18.1 are denoted by “ α_{ij} ” in the model specification.

To understand the way that the model is specified in the BUGS code, it is important to understand how the data are formatted. The values in the program are coded as integer indices 1, 2, 3, ..., and as vectors $1; 0; 0; : : i, h_1; 1; 0; : : i, h_0; 0; 1; : : i, \dots$. By coding x as integers, then nested indexing can be used instead of ~~lists~~ of vectors. Thus, x becomes coded as $[x]$, not $\text{inprod}(a[],x[])$. For the i^{th} observation, the value of

is coded as $x[i]$. Thus, $x[i] \sim f_1; 2; 3; \dots; NxLvl$ for $i \in 1:Ntotal$ where $NxLvl$ is the number of levels of x and $Ntotal$ is the total number of observations.

Here is the BUGS model specification: (ANOVAonewayBRugs.R)

```

11 model {
12   for ( i in 1:Ntotal ) {
13     y[i] ~ dnorm( mu[i] , tau )
14     mu[i] <- a0 + a[x[i]]
15   }
16   #
17   tau <- pow( sigma , -2 )
18   sigma ~ dunif(0,10) # y values are assumed to be standardized
19   #
20   a0 ~ dnorm(0,0.001) # y values are assumed to be standardized
21   #
22   for ( j in 1:NxLvl ) { a[j] ~ dnorm( 0.0 , atau ) }
23   atau <- 1 / pow( aSD , 2 )
24   aSD <- abs( aSDunabs ) + .1
25   aSDunabs ~ dt( 0 , 0.001 , 2 )
26 }
```

The constraint, that the deviations sum to zero, does not appear in the model specification. The BUGS code estimates the baseline and deviation from the constraint, but the MCMC estimates are re-centered at zero by subsequent R code. The non-centered baseline is denoted in the BUGS model as, and the non-centered deviations are denoted. Those non-centered estimates are transformed to respect the sum-to-zero constraint merely by subtracting the mean of the 's from each $a[j]$, and adding the mean to the baseline. Thus, $b[j] = a[j] - \text{mean}(a)$ and $b0 = a0 + \text{mean}(a)$

The constants for the top-level priors are set with the assumption that the data values, y , have been standardized. according to Equation 16.1, p. 318 course, the values cannot be standardized because they are nominal.) This standardization makes it easier to establish reasonable default priors for a range of applications without having to change the prior constants when the application changes, for example from income, on the order of 10^5 dollars, to width of hairs, on the order of 10 millimeters. Nevertheless, when there is strong prior information, it should be incorporated. Exercise 18.2 has you explore robustness of the results when you use different priors.

There is one other trick in the BUGS model specification that is not in the hierarchical diagram of Figure 18.1. One line of the BUGS model specification is the standard deviation of the group effects, denoted $aSDunabs$, comes from a distribution: $aSDunabs \sim dt(0,0.001,2)$. Another line takes the absolute value to "fold" the distribution onto the non-negative numbers, yielding the value $aSD <- abs(aSDunabs) + .1$. But that line also mysteriously adds a small constant, namely 0.1. This constant keeps aSD from venturing extremely close to zero. The reason for keeping aSD away from zero is that shrinkage can become overwhelmingly strong when there are many groups with few data points per group. This becomes especially problematic in the next chapter when we begin interaction of factors.

It turns out that MCMC sampling for this model can be extremely slow. One important way to reduce burn-in time is to start the chains at reasonable positions. We start the overall baseline at the grand mean of the data, and the deviations at the level means minus the grand mean. The variances are also initialized to the corresponding data variances. The full code, including initialization of $chains$, is presented in Section 18.4.1 (ANOVAonewayBRugs.R)

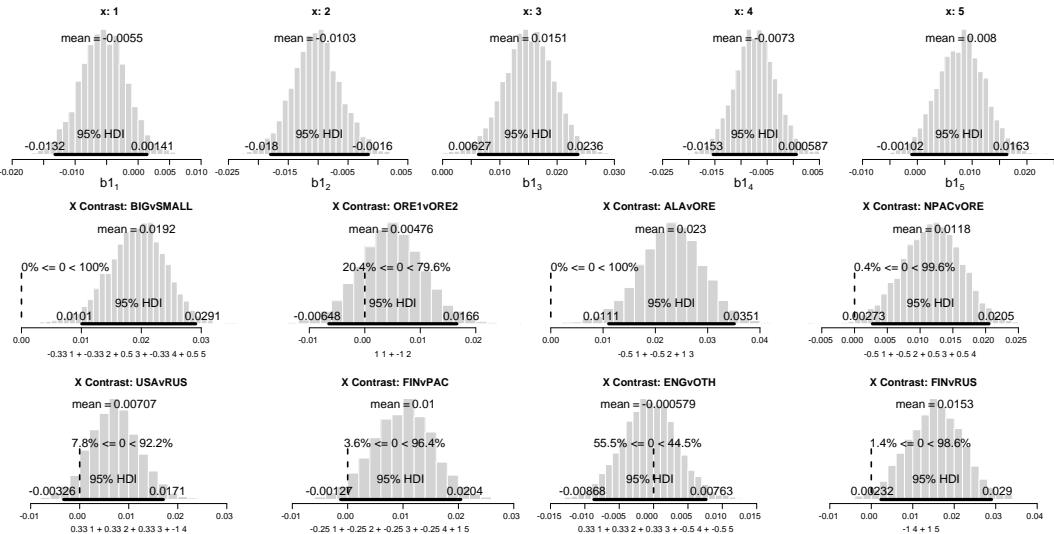


Figure 18.2: Upper row: Posterior estimates ϕ values for data from McDonald (2009; McDonald et al., 1991). The values indicate deviations by each group away from the overall central tendency. Middle and bottom rows: Various complex comparisons of ϕ values. For example, the bottom row, third panel, compares the three English-speaking sites against the two non-English speaking sites (where they might say “midiya myshtsy” or “simpukka lihas” instead of “mussel muscle”).

Because the chains can be highly autocorrelated, extensive thinning is needed, keeping a step only once out of several hundred. Running such long chains can take a long time, and become boring for your computer, which would rather be checking the web for exciting software updates. In the examples presented here, we tolerate the modest waiting times. But there are various methods for re-parameterizing models so that the chains are sampled with less auto-correlation (e.g., Gelman, 2008; Gelman & Hill, 2007, Ch. 19).

18.1.3 A worked example

With all the emphasis these days on physical fitness and muscle building, it's only appropriate to consider an example about muscles. In particular, like to know if geographical location influences muscle size, which might be affected by the weather or amount of daylight. Consider some data regarding muscles from five geographical locations: (1) Tillamook, Oregon; (2) Newport, Oregon; (3) Petersburg, Alaska; (4) Vladivostok, Russia (Pacific coast); and (5) Tvarminne, Finland. The values in the data set are ratios of the anterior adductor muscle scar length to total muscle length, in the mussel species *Mytilus trossulus*. These ratios of scar length to total length tend to be between 5% and 15% (McDonald, 2009; McDonald, Seed, & Koehn, 1991).

Results of the BRugs program listed in Section 18.4 (ANOVA one-way BRugs) are shown in Figure 18.2. The histograms in the upper row show the (main) posterior distributions of the ϕ values for the five geographical locations. These values are deviations away from the baseline ϕ_0 , which is not shown. Some things to keep in mind when interpreting the results: First, the estimates of deviation are subject to “shrinkage”, because the model incorporates the prior structural assumption that all the locations come from the same

overarching distribution. The mean deviations shown in Figure 18.2 are, in fact, a little smaller than the deviations of the actual sample means. ~~Second~~ The model assumes that the precision is the same for all groups; i.e., there is homogeneity of variance. The posterior β_j values are the ones that are believable when also assuming homogeneous variances. If the groups actually have wildly different variances, then the estimates from may be distorted. Third, the marginal distributions on the β_j cannot be used to directly infer differences between groups, because the parameters might be correlated; the deviations tend to be negatively correlated, because increasing the estimate for one group suggests decreasing the estimated deviation for another, if they remain symmetric around the baseline. To judge differences between groups, the differences must be computed directly.

18.1.3.1 Contrasts and complex comparisons

The middle and bottom rows of Figure 18.2 shows several ~~comparisons~~ for the mussel muscle results. A comparison amounts to a difference between an average of some groups and an average of other groups. For example, to compare ~~the Pacific Ocean mussels~~ against the one non-Pacific (Baltic Sea) mussel, we multiply deviations (β_j 's) of the first four groups by $\frac{1}{4}$ to get their average, and subtract it from the deviation of the fifth group to get the difference. The difference is called a contrast and when the comparison involves a contrast of averages, instead of ~~a test~~ of two specific groups, it is sometimes called a complex comparison. The contrast is fully specified by the coefficients on the groups, which can be placed into a vector of contrast coefficients. For example, the contrast coefficients for comparing Pacific Ocean mussels against Baltic mussels are $1/4; 1/4; 1/4; 1/4; +1$. Notice that the coefficients sum to zero. We compute the difference at every step in the MCMC chain, and examine the ~~result~~ distribution of believable differences. The distribution for this particular example is shown in the bottom row, second panel, of Figure 18.2, where it can be seen that over 96% of the believable differences lie on one side of zero, and the 95% HDI just spans zero. From these results we may not want to declare categorically that there is a ~~big~~ difference between Finland and the other sites; the decision depends on how you set your ROPE. Regardless of your decision rule, the posterior does tell us the most believable difference and the uncertainty in that difference.

Figure 18.2 shows a variety of comparisons that might be of interest. For example, the first panel of the middle row compares the two sites with the ~~giant~~ muscles against the three other sites. This sort of comparison would be labeled “hoc” by traditional analyses, because we might not have specified which sites would ~~exist~~ before collecting the data. The second panel in the middle row contrasts the ~~two sites~~ Oregon. The third panel in the middle row compares the Alaska site against the ~~average~~ of the two Oregon sites.

We can make all the comparisons shown in Figure 18.2, and ~~any others~~ as we like, without worrying about inflated false alarm rates, ~~because~~ the posterior distribution does not change when we consider additional comparisons. ~~The~~ Posterior distribution is the best inference we can make based on the data we have ~~and the beliefs~~ we started with. It is possible that the random data in our sample are ~~especially~~ unrepresentative of the underlying population, but we cannot know. Fortunately, because of the incorporation of our prior knowledge about how estimates in different locations can mutually inform each other, the estimates undergo shrinkage, which helps mitigate the effect of rogue data. In many applications, the shrinkage yields decisions similar to those that would result from NHST “corrections” for multiple comparisons. But unlike NHST corrections, the

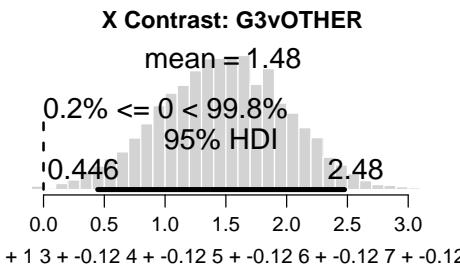


Figure 18.3: A comparison of Group 3 versus the average of other groups, for the data in Solari et al. (2008, Table 3, p. 495). (The specification of contrast coefficients on the x-axis overflows the margins of the figure because there are so many groups. The contrast coefficients on the nine groups are 1=8; 1=8;+1; 1=8;:::, which, when rounded to two decimal places, appear as 0:12; 0:12;+1; 0:12;:::.)

shrinkage in the Bayesian approach is based on explicit prior knowledge, and is not affected by which or how many comparisons are intended. (For a discussion of these issues, see Section 17.2, regarding decisions about regression coefficients, and Section 11.4, regarding multiple comparisons of groups)

18.1.3.2 Is there a difference?

The contrasts and complex comparisons in Figure 18.2 were judged to be credibly non-zero if the 95% HDI excluded (a ROPE around) zero. Evidence would be deemed to be practically equivalent to zero if its HDI fell entirely within a ROPE. This decision procedure is attractive because all the groups are simultaneously estimated, with mutually informed shrinkage, and from priors that are also appropriately informed (which entails also being agreeable to a skeptical audience).

Some researchers prefer to pose the question, “Is there evidence?”, as a model comparison on two priors. One prior expresses the null hypothesis that the contrast has zero magnitude, while the other prior expresses a complementary hypothesis that any magnitude contrast is possible. This approach was discussed extensively in Section 12.2, beginning on p. 245.

There are two attractions to the two-prior, model-comparison approach. One attraction is that the model comparison can yield posterior odds in favor of the null, unlike NHST, which can only reject a null hypothesis but never accept it. Another attraction is that the complementary hypothesis is usually intended to be an “automatic” uninformed prior that is chosen for mathematical felicity. The hope is that an automatic prior obviates debate about how prior information should be expressed.

As was argued in Section 12.2, the two-prior approach should be applied cautiously. First, it is important to emphasize that the two-prior approach only indicates which prior is relatively less unbelievable. If either prior is theoretically untenable in the first place, then the “automatic” model comparison is automatically uninformative. Thus, the two-prior approach should only be applied to situations in which (i) it is theoretically appropriate to posit that a particular contrast really can be exactly zero, and (ii) the alternative prior incorporates prior knowledge about the plausible magnitude of the difference.

As an example, consider a situation presented by Solari and Sun (2008, Table 3, p. 495). There were nine groups, with a metric dependent variable. The dependent variable was the acetic acid content of tomatoes, and the nine groups were different types of manuring during growth of the tomatoes. The mean of Group 3 was found to be different than other groups. To test whether Group 3 was different, the authors conducted a Bayesian model comparison of two priors: The null-hypothesis prior had all nine groups with iden-

tical means. The alternative prior had Group 3 with a separate estimated mean while the other eight groups had identical means. The resulting Bayes factor (BF) strongly favored the alternative prior. Does this result suggest that the alternative prior is what we should believe? Unfortunately, no. The BF tells us that the prior ~~with~~ weight equal means and one different mean for Group 3 is more believable than the prior ~~with~~ equal means (assuming that the priors on the two hypotheses were 50-50). But the prior with eight equal means, on groups other than Group 3, is already untenable because we ~~not~~ believe that the eight groups have identical means. Moreover, the estimate of difference, between Group 3 and the other groups, is not what we want, because the estimate ~~do not~~ take into account variation among the eight other groups.

When instead we conduct a Bayesian analysis using the BANOVA model, we obtain a posterior that simultaneously estimates all the separate group effects, with shrinkage, from a plausibly informed prior. The complex comparison of Group 3 against the other eight groups is shown in Figure 18.3, where it can be seen that the magnitude of the contrast is credibly greater than zero. In this application, there is need to pursue a BF approach to group comparisons.

It is also worth reiterating that the two-prior, model-comparison approach can arrive at a conclusion opposite that of the one-prior, estimation approach. Recall Figure 12.5, p. 249, which showed that a model comparison preferred the hypothesis of identical groups to the alternative hypothesis of all different groups, even though an estimation of effects in the alternative hypothesis showed a credible difference among groups. The point in that case was that the null model, even though it was a poor fit, was less bad than the alternative model. Follow-up model comparisons would be required to narrow down which combination of group equivalences was least implausible. And even after that, we would not necessarily want to believe that any of the groups truly equivalent, because we know in advance that they were treated differently. Instead, we desire an estimate of the differences and the precision of the estimate. That situation involved a dichotomous dependent variable, but the analogous situation can arise with metric dependent variables.

The two-prior, model-comparison approach can be appropriate in situations where actual equivalence is tenable and the goal is to identify which conditions are plausibly equivalent, or situations in which zero-magnitude effects are tenable and the goal is to identify which conditions have zero effect. In those situations, it behoves the researcher to pursue the model-comparison or related approaches (see, e.g., Hochberg, 1999; Gopalan & Berry, 1998; Mueller et al., 2007; Scott & Berger, 2006). Moreover, Bayesian model comparison is highly advisable when the two models are genuinely competitive that express different explanations of the data. In these situations, it is important that the priors in the two models are equivalently informed, so that one model is not at a disadvantage because of an infelicity in an arbitrary, automatic prior.

18.2 Multiple comparisons

In twentieth century null-hypothesis significance testing (NHST), there is an immense literature regarding how to compute the “true” significance (probability of false alarm) of an apparent difference between groups, when the analyst is conducting comparisons of multiple groups. The problem is that when more comparisons are conducted, there are more opportunities for a spuriously large difference to appear by accident. In other words, there are more opportunities for false alarms. Notice that this problem of inflated false alarm

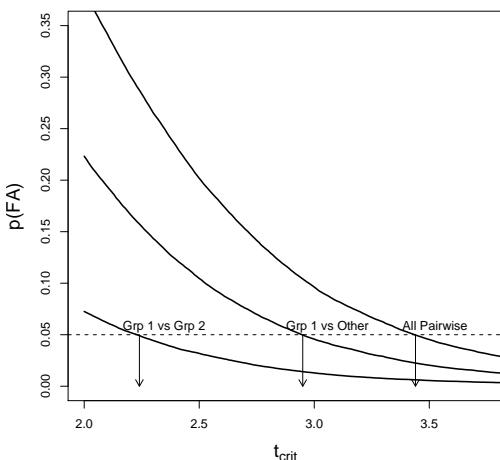


Figure 18.4: The probability of false alarm as a function of critical value, with separate curves for different sets of comparisons. All groups have $N = 6$ fixed by intention. The curve labeled “Grp 1 vs Grp 2” is for a single comparison of two groups, and corresponds with the usual two-group distribution. The curve labeled “Grp 1 vs Other” refers to four paired comparisons, of Group 1 versus each of the other four groups. The last curve is for the set of all 10 pairwise comparisons.

rates arises because NHST is based on the intentions of the analyst. If the analyst intends to make lots of comparisons between various combinations of groups, then there is greater opportunity for false alarms. If the analyst intends to make a few comparisons between groups, then there is less opportunity for false alarms.

For example, consider again the sea mussel data. Group 4 (Russia) and Group 5 (Finland) seem to be different, and it is meaningful to plan a comparison between them because of their geographical difference. If we run a two-group test, we get $t = 2.53$; $p = 0.028$, which denotes a significant difference. On the other hand, if we run a post-hoc test of all pairwise comparisons, using Tukey's Honest Significant Difference correction, then we find that $p = 0.093$, and the difference is not significant. So, do Russia and Finland really differ or not? According to NHST, the answer depends on your intentions: If you intended to compare only those two, then they are significantly different, but if you intended to make all pairwise comparisons, then they are not significantly different.

Section 11.4, p. 227, discussed multiple comparisons in the context of a dichotomous dependent variable. Here we reiterate those in the context of a metric dependent variable.

Suppose that we have two groups: One group is patients treated with a placebo and a second group is patients treated with a totally ineffective drug. We measure a metric variable, e.g., body temperature. Because there is no actual difference between the treatments, the underlying distributions of body temperatures are identical for the two groups; we will suppose that they are normally distributed with equal means and equal variances. When we run an experiment, we are collecting a random sample of 6 from each of the groups. The random samples might show a spuriously large difference between their means, just by chance, despite the fact that on average, in the long run, the true differences are identical.

To determine how often the spuriously large differences occur, we can simulate conducting the experiment over and over. For every simulated experiment, we compute the difference of means between the samples from the groups. The difference of sample means is in units of the original measurement scale, e.g., degrees Fahrenheit or degrees Celsius. To get rid of the arbitrary nature of the measurement scale, standardize the difference of means and call the result the t -statistic. Because the true difference between groups is zero, the t value typically will be near zero. Occasionally, by chance, the t value will be

far above or far below zero. The lowest curve in Figure 18.4 is the probability that the sampled value falls above the critical value on the abscissa. For example, the probability that the sample value falls above $t_{crit} = 2.23$ is $p(FA) = 0.05$; this is marked by an arrow. In NHST, the decision rule is to reject the null hypothesis if the sample exceeds a critical value that is selected to keep false alarms to only 5%. Then comparing Group 1 with Group 2, we would reject the null hypothesis if $t > 2.23$, because that would happen only 5% of the time by chance alone.

Now consider an expanded experiment, in which there is a placebo treatment and four distinct drugs, for a total of five treatment groups. According to the null hypothesis, the five treatment groups have identical distributions of body temperatures (normally distributed with equal means and variances). However, because of random sampling in any particular experiment, some treatment samples will have higher or lower mean temperatures than other treatment samples. Suppose that before we collect any data, we plan to compare the placebo group (Group 1) with each of the four drug groups, we plan four pairwise comparisons. Each of these comparisons might yield a large difference merely by chance, even when there is truly no difference in the underlying distributions. We can determine how often these chance extremes happen by running Monte Carlo simulation. For a simulated experiment, we randomly sample 6 scores from each of the 5 groups, and compute the values of each of the 4 comparisons. The simulated experiments are repeated many times. For each candidate, we see what proportion of simulated experiments had a comparison that exceeded that critical value. The middle panel of Figure 18.4 shows the result. Notice that at any given value of t_{crit} , there is now a much higher probability that the simulated experiment will have at least one comparison larger. In particular, to keep the false alarm rate down to 5% t_{crit} must be about 2.95 instead of 2.33.

If we did not plan only four tests, but instead decided to compare every group with every other group, then we would have even more opportunity for false alarms. With 5 groups, there are 10 pairwise comparisons. If we simulate experiments from equal distributions as before, but this time consider all 10 values, the probability of false alarm is higher yet, as shown in the right curve of Figure 18.4. The critical value has risen even higher, to approximately 3.43.

Now, suppose we actually run the experiment. We randomly assign 30 people to the 5 groups, 6 people per group. The first group gets the placebo, the other four groups get the corresponding four drugs. We are careful to make this a double-blind experiment: Neither the subjects nor experimenters know who is getting which treatment. Moreover, no one knows whether any other person is even in the experiment. We collect the data. Our first question is to compare the placebo and the first drug, group 1 versus group 2. We compute the statistic for the data from the two groups and find that $t = 2.95$. Do we decide that the two treatments had significantly different effects?

The answer, bizarrely, depends on the intentions of the experimenter. Suppose, for instance, that we handed the data from the first two groups to a research assistant, who is asked to test for a difference between groups. The assistant runs a t -test and finds $t = 2.95$, declaring it to be highly significant because it greatly exceeds the critical value of 2.23 for a two-group-test. Suppose, on the other hand, that we handed the data from all five groups to a different research assistant, who is asked to compare the placebo against each of the other four. This assistant runs a t -test of group 1 versus group 2 and finds $t = 2.95$,

¹For a discussion of various correction procedures and where they come from, see Figure 5.1 of Maxwell and Delaney (2004). If you must learn NHST methods, this is a good resource.

declaring it to be marginally significant because it just squeezes past the critical value of 2.95 for these four planned comparisons. Suppose, on the other hand, that we handed the data from all five groups to a different research assistant, who is told to conduct all pairwise comparisons post-hoc because we have no strong hypotheses about which treatments will have beneficial or detrimental or neutral effects. This assistant runs a t -test of group 1 versus group 2 and $t_{ndf} = 2.95$, declaring it to be not significant because it fails to exceed the critical value of 3.43 that is used for post-hoc pairwise comparisons. Notice that regardless of which assistant analyzed the data, the t -value for the two groups stayed the same because the data of the two groups stayed the same. Indeed, the t -value completely uninfluenced by the intentions of the analyst. So why should the interpretation of the data be influenced by the intentions of the analyst? It shouldn't.

And if you believe that the interpretation should be informed by the intention of the analyst, how do you determine the intention of the analyst? Did the analyst truly plan only those particular comparisons, or did the analyst really consider but jettison them once the data were in? Or, did the analyst actually plan fewer comparisons, but realize later that additional comparisons should be made to address theoretical issues? Or, did the analyst actually plan to include two other treatment groups in the study, but then not actually include those groups in the analysis because of administrative constraints committed during the data collection? Or what if the experiment was planned by a team of people, some of whom planned some comparisons, and others of whom planned other comparisons? Conclusion: Establishing the true intention of the analyst is not only impossible, it is also impossible.

Multiple comparisons are not a problem in a Bayesian analysis (Gelman et al., 2009). The posterior distribution is a fixed entity in high-dimensional parameter space, and making comparisons between groups is simply examining the posterior distribution from different perspectives or margins. The posterior does not change when new comparisons come to mind.

The posterior is not immune to spurious coincidences of data, of course. False alarms are mitigated, however, by incorporating prior knowledge into the structure of the model. The estimates of the groups are mutually informative, estimation of higher-level structure, and shrinkage of estimates across groups attenuates false alarms. The attenuation of false alarms is governed by the data, not by unknowable truths.

18.3 Two group Bayesian ANOVA and the NHST t test

The idea behind a NHST t test is simple: We have two groups, each with a mean. We compute the difference of the means, and standardize that difference relative to the standard deviation of the scores within the groups. The resulting standardized difference is called the t value. We want to know whether the observed t value is significantly different from zero, so we compare the t value to a sampling distribution of t values (Gosset, 1908). The sampling distribution assumes that the intention of the experimenter was to stop when there were exactly N_1 values observed for the first group, and exactly N_2 values observed for the second group.

The t test is a special case of NHST ANOVA when there are only two groups. More specifically, when the two groups are assumed to have equal variances in the underlying population, then the t value squared equals the F value in two-group ANOVA. (And what's an F value, you may ask? The F value is the summary statistic used in NHST ANOVA to express how much the groups differ from each other. It's the ratio of the variance between

group means, relative to the variance within groups.)

In typical applications of BANOVA, the prior on the betweengroup variance is only mildly informed. In this case, a BANOVA on two groups imposes shrinkage on the group estimates because there are so few groups. It is only when several groups “gang up” that they strongly inform the estimate of the variation between groups, and therefore constrain the estimates of other groups. When the prior on variance within groups is also vague, the results of a two-group BANOVA closely agree with the results of an NHST t-test. Exercise 18.1 has you explore this correspondence.

18.4 R code

18.4.1 Bayesian oneway ANOVA

(ANOVAonewayBRugs.R)

```

1 graphics.off()
2 rm(list=ls(all=TRUE))
3 fnroot = "ANOVAonewayBrugs"
4 library(BRugs)      # Kruschke, J. K. (2010). Doing Bayesian data analysis:
5                      # A Tutorial with R and BUGS. Academic Press / Elsevier.
6 #-----
7 # THE MODEL.
8
9 modelstring = "
10 # BUGS model specification begins here...
11 model {
12   for ( i in 1:Ntotal ) {
13     y[i] ~ dnorm( mu[i] , tau )
14     mu[i] <- a0 + a[x[i]]
15   }
16   #
17   tau <- pow( sigma , -2 )
18   sigma ~ dunif(0,10) # y values are assumed to be standardized
19   #
20   a0 ~ dnorm(0,0.001) # y values are assumed to be standardized
21   #
22   for ( j in 1:NxLvl ) { a[j] ~ dnorm( 0.0 , atau ) }
23   atau <- 1 / pow( aSD , 2 )
24   aSD <- abs( aSDunabs ) + .1
25   aSDunabs ~ dt( 0 , 0.001 , 2 )
26 }
27 # ... end BUGS model specification
28 " # close quote for modelstring
29 # Write model to a file, and send to BUGS:
30 writeLines(modelstring,con="model.txt")
31 modelCheck( "model.txt" )
32
33 #-----
34 # THE DATA.
35
36 # Specify data source:
37 dataSource = c( "McDonaldSK1991" , "SolariLS2008" , "Rando m" )[1]
38 # Load the data:
39
40 if ( dataSource == "McDonaldSK1991" ) {
41   fnroot = paste( fnroot , dataSource , sep="" )

```

```

42 datarecord = read.table( "McDonaldSK1991data.txt", head      =T ,
43                         colClasses=c("factor","numeric") )
44 y = as.numeric(datarecord$Size)
45 Ntotal = length(datarecord$Size)
46 x = as.numeric(datarecord$Group)
47 xnames = levels(datarecord$Group)
48 NxLvl = length(unique(datarecord$Group))
49 contrastList = list( BIGvSMALL = c(-1/3,-1/3,1/2,-1/3,1/2) ,
50                      ORE1vORE2 = c(1,-1,0,0,0) ,
51                      ALAvORE = c(-1/2,-1/2,1,0,0) ,
52                      NPACvORE = c(-1/2,-1/2,1/2,1/2,0) ,
53                      USAvRUS = c(1/3,1/3,1/3,-1,0) ,
54                      FINvPAC = c(-1/4,-1/4,-1/4,-1/4,1) ,
55                      ENGvOTH = c(1/3,1/3,1/3,-1/2,-1/2) ,
56                      FINvRUS = c(0,0,0,-1,1) )
57 }
58
59 if ( dataSource == "SolariLS2008" ) {
60   fnroot = paste( fnroot , dataSource , sep="" )
61   datarecord = read.table("SolariLS2008data.txt", header      =T ,
62                         colClasses=c("factor","numeric") )
63   y = as.numeric(datarecord$Acid)
64   Ntotal = length(datarecord$Acid)
65   x = as.numeric(datarecord$type)
66   xnames = levels(datarecord$type)
67   NxLvl = length(unique(datarecord$type))
68   contrastList = list( G3vOTHER = c(-1/8,-1/8,1,-1/8,-1/8,
69                         -1/8,-1/8,-1/8,-1/8) )
70 }
71 if ( dataSource == "Random" ) {
72   fnroot = paste( fnroot , dataSource , sep="" )
73   #set.seed(47405)
74   ysdtrue = 4.0
75   a0true = 100
76   atrue = c( 2 , -2 ) # sum to zero
77   npercell = 8
78   datarecord = matrix( 0, ncol=2 , nrow=length(atrue)*nperc      = ell )
79   colnames(datarecord) = c("y","x")
80   rowidx = 0
81   for ( xidx in 1:length(atrue) ) {
82     for ( subjidx in 1:npercell ) {
83       rowidx = rowidx + 1
84       datarecord[rowidx,"x"] = xidx
85       datarecord[rowidx,"y"] = ( a0true + atrue[xidx] + rnorm(1,      0,ysdtrue) )
86     }
87   }
88   datarecord = data.frame( y=datarecord[, "y"] , x=as.factor      = r(datarecord[, "x"]) )
89   y = as.numeric(datarecord$y)
90   Ntotal = length(y)
91   x = as.numeric(datarecord$x)
92   xnames = levels(datarecord$x)
93   NxLvl = length(unique(x))
94   # Construct list of all pairwise comparisons, to compare with      h NHST TukeyHSD:
95   contrastList = NULL
96   for ( g1idx in 1:(NxLvl-1) ) {
97     for ( g2idx in (g1idx+1):NxLvl ) {
98       cmpVec = rep(0,NxLvl)
99       cmpVec[g1idx] = -1
100      cmpVec[g2idx] = 1

```

```

101     contrastList = c( contrastList , list( cmpVec ) )
102   }
103 }
104 }
105
106 # Specify the data in a form that is compatible with BRugs mode I, as a list:
107 ySDorig = sd(y)
108 yMorig = mean(y)
109 z = ( y - yMorig ) / ySDorig
110 dataList = list(
111   y = z ,
112   x = x ,
113   Ntotal = Ntotal ,
114   NxLvl = NxLvl
115 )
116 # Get the data into BRugs:
117 modelData( bugsData( dataList ) )
118
119 #-----
120 # INTIALIZE THE CHAINS.
121
122 # Autocorrelation within chains is large, so use several cha      ins to reduce
123 # degree of thinning. But we still have to burn-in all the chai      ns, which takes
124 # more time with more chains (on serial CPUs).
125 nchain = 5
126 modelCompile( numChains = nchain )
127
128 if ( F ) {
129   modelGenInits() # often won't work for diffuse prior
130 } else {
131   # initialization based on data
132   theData = data.frame( y=datalist$y , x=factor(x,labels=x           names) )
133   a0 = mean( theData$y )
134   a = aggregate( theData$y , list( theData$x ) , mean )[,2] - a0
135   ssw = aggregate( theData$y , list( theData$x ) ,
136                     function(x){var(x)*(length(x)-1)} )[,2]
137   sp = sqrt( sum( ssw ) / length( theData$y ) )
138   genInitList <- function() {
139     return(
140       list(
141         a0 = a0 ,
142         a = a ,
143         sigma = sp ,
144         aSDunabs = sd(a)
145       )
146     )
147   }
148   for ( chainIdx in 1 : nchain ) {
149     modellnits( bugsInits( genInitList ) )
150   }
151 }
152
153 #-----
154 # RUN THE CHAINS
155
156 # burn in
157 BurnInSteps = 10000
158 modelUpdate( BurnInSteps )
159 # actual samples

```

```

160 samplesSet( c( "a0" , "a" , "sigma" , "aSD" ) )
161 stepsPerChain = ceiling(5000/nchain)
162 thinStep = 750
163 modelUpdate( stepsPerChain , thin=thinStep )
164
165 #-----
166 # EXAMINE THE RESULTS
167
168 source("plotChains.R")
169 source("plotPost.R")
170
171 checkConvergence = T
172 if ( checkConvergence ) {
173   sumInfo = plotChains( "a0" , saveplots=T , filenameroot=fn
174   sumInfo = plotChains( "a" , saveplots=T , filenameroot=fnr
175   sumInfo = plotChains( "sigma" , saveplots=T , filenameroot
176   sumInfo = plotChains( "aSD" , saveplots=T , filenameroot=f
177 }
178
179 # Extract and plot the SDs:
180 sigmaSample = samplesSample("sigma")
181 aSDSample = samplesSample("aSD")
182 windows()
183 layout( matrix(1:2,nrow=2) )
184 par( mar=c(3,1,2.5,0) , mgp=c(2,0.7,0) )
185 plotPost( sigmaSample , xlab="sigma" , main="Cell SD" , bre aks=30 )
186 plotPost( aSDSample , xlab="aSD" , main="a SD" , breaks=30 )
187 dev.copy2eps(file=paste(fnroot,"SD.eps",sep=""))
188
189 # Extract a values:
190 a0Sample = samplesSample( "a0" )
191 chainLength = length(a0Sample)
192 aSample = array( 0 , dim=c( datalist$NxLvl , chainLength ) )
193 for ( xidx in 1:datalist$NxLvl ) {
194   aSample[xidx,] = samplesSample( paste("a[",xidx,"]",se p="" ) )
195 }
196
197 # Convert to zero-centered b values:
198 mSample = array( 0 , dim=c( datalist$NxLvl , chainLength ) )
199 for ( stepIdx in 1:chainLength ) {
200   mSample[stepIdx ] = ( a0Sample[stepIdx] + aSample[stepIdx] )
201 }
202 b0Sample = apply( mSample , 2 , mean )
203 bSample = mSample - matrix(rep( b0Sample ,NxLvl),nrow=NxLvl,byrow=T)
204 # Convert from standardized b values to original scale b values:
205 b0Sample = b0Sample * ySDorig + yMorig
206 bSample = bSample * ySDorig
207
208 # Plot b values:
209 windows(datalist$NxLvl*2.75,2.5)
210 layout( matrix( 1:datalist$NxLvl , nrow=1 ) )
211 par( mar=c(3,1,2.5,0) , mgp=c(2,0.7,0) )
212 for ( xidx in 1:datalist$NxLvl ) {
213   plotPost( bSample[xidx,] , breaks=30 ,
214             xlab=bquote(beta*1[.(xidx)]) ,
215             main=paste("x:",xnames[xidx]) )
216 }
217 dev.copy2eps(file=paste(fnroot,"b.eps",sep=""))
218

```

```

219 # Display contrast analyses
220 nContrasts = length( contrastList )
221 if ( nContrasts > 0 ) {
222   nPlotPerRow = 5
223   nPlotRow = ceiling(nContrasts/nPlotPerRow)
224   nPlotCol = ceiling(nContrasts/nPlotRow)
225   windows(3.75*nPlotCol,2.5*nPlotRow)
226   layout( matrix(1:(nPlotRow*nPlotCol),nrow=nPlotRow,n      col=nPlotCol,byrow=T) )
227   par( mar=c(4,0.5,2.5,0.5) , mgp=c(2,0.7,0) )
228   for ( cIdx in 1:nContrasts ) {
229     contrast = matrix( contrastList[[cIdx]],nrow=1) # make it           a row matrix
230     inclIdx = contrast!=0
231     histInfo = plotPost( contrast %*% bSample , compVal=0 , brea    ks=30 ,
232                           xlab=paste( round(contrast[inclIdx],2) , xnames[inclIdx]           ,
233                                         c(rep("+",sum(inclIdx)-1),"") , collapse=" " ) ,
234                           cex.lab = 1.0 ,
235                           main=paste( "X Contrast:", names(contrastList)[cIdx] ) )
236   }
237   dev.copy2eps(file=paste(fnroot,"xContrasts.eps",sep      =""))
238 }
239
240 #####  

241 # Do NHST ANOVA and t tests:  

242
243 theData = data.frame( y=y , x=factor(x,labels=xnames) )
244 aovresult = aov( y ~ x , data = theData ) # NHST ANOVA
245 cat("\n-----\n")
246 print( summary( aovresult ) )
247 cat("\n-----\n")
248 print( model.tables( aovresult , "means" ) , digits=4 )
249 windows()
250 boxplot( y ~ x , data = theData )
251 cat("\n-----\n")
252 print( TukeyHSD( aovresult , "x" , ordered = FALSE ) )
253 windows()
254 plot( TukeyHSD( aovresult , "x" ) )
255 if ( T ) {
256   for ( xIdx1 in 1:(NxLvl-1) ) {
257     for ( xIdx2 in (xIdx1+1):NxLvl ) {
258       cat("\n-----\n")
259       cat( "xIdx1 = " , xIdx1 , " , xIdx2 = " , xIdx2 ,
260             " , M2-M1 = " , mean(y[x==xIdx2])-mean(y[x==xIdx1]) , "\n" )
261       print( t.test( y[x==xIdx2] , y[x==xIdx1] , var.equal=T ) ) #      t test
262     }
263   }
264   cat("\n-----\n")
265
266 #####  


```

18.5 Exercises

Exercise 18.1 [Purpose: Notice that Bayesian ANOVA with two groups tends to agree with an NHST t test.] The BRugs program of Section 18.4 ANOVAonewayBRugs.R allows you to specify random data. It executes a Bayesian ANOVA, and at the end of the program it also conducts an NHST ANOVA and t tests (using R's `aov` and `t.test` functions). Run the program

ten times with different random data by commenting out the `seed` command. Specify `ysdtrue = 4.0`, `atru = c(2,-2)` (which implies two groups because there are two dections), and `percell = 8`. For each run, record, by hand, (i) how much of the posterior difference between means falls on one side of zero (see the histogram with the main title "X Contrast" and x-axis labeled "1 1 + 1 2"), (ii) whether the 95% HDI excludes zero, and (iii) the confidence interval and p-value of the NHST t-test. Do the t-test and the BANOVA usually agree in their decisions about whether the group means are different?

Exercise 18.2. [Purpose: Influence of the prior in Bayesian ANOVA.] In the model section of the BRugs program of Section 18.4.1~~(ANOVAonewayBRugs.R)~~, and correspondingly in the diagram of Figure 18.1, there are several constants that determine the prior. These constants include the mean value of the baseline (M_0 in the diagram), the precision on the baseline (H_y in the diagram), the precision of the foldedt distribution (T_t in the diagram), and the upper value of the uniform distribution on y ($H_{\bar{y}}$ in the diagram). Because the data are standardized, M_0 should be set at zero, and H_y can be modest (not terribly small). $H_{\bar{y}}$ also can be set to a modest value because the data are standardized. But what is the precision of the foldedt distribution, T_t ? This constant modulates the degree of shrinkage: A large T_t indicates prior knowledge that the groups do not differ much, and imposes a high degree of shrinkage that must be overcome by the data.

Run the program on the mussel data using a small value of T_t , such as 1.0E-6, and a large value of T_t , such as 1000. Are the results very different? Discuss which prior value might be appropriate.

Exercise 18.3. [Purpose: Bayesian ANOVA without assuming equal variances.] Modify the program in Section 18.4.1~~(ANOVAonewayBRugs.R)~~ so that it allows a different variance for each group, with the different variances coming from a hyperdistribution that has precision informed by the data. In other words, instead of assuming the same σ^2_y ($= 1 = \frac{2}{N}$) for all the levels of x , we allow each group to have its own variance. Denote the variance of j^{th} group as σ^2_j , analogous to the deviation j . Just as the group deviations are assumed to come from a higher-level distribution, we will assume that group SD's come from a higher-level distribution. Because SD's must be non-negative, use a gamma density for the higher-level distribution. The gamma distribution has parameters for which you need to establish a prior. See the right side of Figure 16.11, p. 356, for guidance. Corresponding code is provided in a hint, below. Run the program on the mussel muscle data. Are the conclusions about the group means any different than when assuming equal variances across groups?

Hint regarding the conclusion: The posteriors on the group means are only a little different in this case, because the group variances are roughly the same. But, because the group variances are less constrained when they are allowed to be different, they are less certain. Therefore the group means are a little less certain, and therefore the differences of means are a little less certain.

Programming hints: Here are some code snippets, showing model specification and chain initialization. (ANOVAonewayNonhomogvarBRugs.R)

```

11 model {
12   for ( i in 1:Ntotal ) {
13     y[i] ~ dnorm( mu[i] , tau[x[i]] )
14     mu[i] <- a0 + a[x[i]]
15   }
16   a0 ~ dnorm(0,0.001)
17   for ( j in 1:NxLvl ) {
```

```

18      a[j] ~ dnorm( 0.0 , atau )
19      tau[j] ~ dgamma( sG , rG )
20  }
21  sG <- pow(m,2)/pow(d,2)
22  rG <- m/pow(d,2)
23  m ~ dgamma(1,1)
24  d ~ dgamma(1,1)
25  atau <- 1 / pow( aSD , 2 )
26  aSD <- abs( aSDunabs ) + .1
27  aSDunabs ~ dt( 0 , 0.001 , 2 )
28 }

```

(ANOVAonewayNonhomogvarBrugs.R)

```

133 # initialization based on data
134 theData = data.frame( y=datalist$y , x=factor(x,labels=x           names) )
135 a0 = mean( theData$y )
136 a = aggregate( theData$y , list( theData$x ) , mean )[,2] - a0
137 tau = 1/(aggregate( theData$y , list( theData$x ) , sd )[,2])      ^2
138 genInitList <- function() {
139   return(
140     list(
141       a0 = a0 ,
142       a = a ,
143       tau = tau ,
144       m = mean( tau ) ,
145       d = sd( tau ) ,
146       aSDunabs = sd(a)
147     )
148   )
149 }
150 for ( chainIdx in 1 : nchain ) {
151   modellnits( bugsInits( genInitList ) )
152 }

```


Chapter 19

Metric Predicted Variable with Multiple Nominal Predictors

Contents

19.1	Bayesian multi-factor ANOVA	22
19.1.1	Interaction of nominal predictors	422
19.1.2	The hierarchical prior	442
19.1.3	An example in R and BUGS	425
19.1.4	Interpreting the posterior	428
19.1.4.1	Metric predictors and ANCOVA	428
19.1.4.2	Interaction contrasts	429
19.1.5	Non-crossover interactions, rescaling, and homogenes variances	430
19.2	Repeated measures, a.k.a. within-subject designs	432
19.2.1	Why use a within-subject design? And why not?	434
19.3	R code	435
19.3.1	Bayesian two-factor ANOVA	435
19.4	Exercises	444

Sometimes I wonder just how it could be, that
Factors aligned so you'd end up with me.
All of the priors made everyone think, that
Our interaction was destined to shrink.

In this chapter we consider situations with a metric predicted variable and multiple nominal predictor variables. For example, we might want to predict income (a metric variable) on the basis of political party affiliation (a nominal variable) and ethnicity (another nominal variable). Or, we may want to predict response time (a metric variable) on the basis of hand used for the response (a nominal value: dominant hand or non-dominant hand) and modality of stimulus (another nominal value: visual, auditory or tactile). These situations are modeled by the cell in the first row and last column of Table 1, p. 312.

In traditional NHST, this situation is known as multifactor ANOVA. We use the same underlying model, but without reference to sampling distributions, and instead with hierarchical priors that provide additional structural constraints. Multifactor ANOVA is a straightforward extension of the model in the previous chapter including the a new

concept of interaction between nominal variables. Just as simple regression considered interaction of metric predictors, multifactor ANOVA considers interaction of nominal predictors.

19.1 Bayesian multi-factor ANOVA

Recall from the previous chapter that in oneway ANOVA we describe the effect of each level of the predictor as a deviation away from an overall mean, where the baseline is the central tendency across all levels of the predictor. In multifactor ANOVA, the same idea applies to two or more predictors, and the deviations are to each predictor added. We'll use notation analogous to the previous chapter, but with extra subscripts to indicate the different predictors, just as we used in multiple regression on continuous predictors.

The mathematical notation was introduced as a case of the general linear model in Section 14.1.6.2, p. 302. Suppose we have two nominal predictors denoted x_1 and x_2 . These predictor vectors can only take on values 0, 1, ..., 0, ..., 1, ..., and so on, with the j^{th} component having the value 1 when the predictor has its nominal level.

When the effects of the two predictors are additive, the predicted values are as follows:

$$\begin{aligned} y &= \beta_0 + \beta_1 x_{1j} + \beta_2 x_{2k} \\ &= \beta_0 + \sum_{j=1}^{J-1} \beta_{1;j} x_{1;j} + \sum_{k=1}^{K-1} \beta_{2;k} x_{2;k} \end{aligned}$$

To make the parameter values unique, we include the constraints

$$\sum_{j=1}^{J-1} \beta_{1;j} = 0 \quad \text{and} \quad \sum_{k=1}^{K-1} \beta_{2;k} = 0$$

Those equations repeat Equations 14.7 and 14.8. In words, β_0 establishes the overall baseline from which the predictors indicate deviations. When predictor x_1 has value $x_{1;j}$, a deviation of $\beta_{1;j}$ is added to the baseline, and when predictor x_2 has value $x_{2;k}$, a deviation of $\beta_{2;k}$ is also added to the baseline. The deviations may be negative, and across all levels of the predictors, the constraints demand much negative deviation as positive deviation, so that the deviations sum to zero for each predictor.

19.1.1 Interaction of nominal predictors

The effect of two predictors may be non-additive, in which case we say that there is an "interaction" of the predictors. For example, if a flame is under a hot-air balloon, its levity will increase. And if hydrogen is added to a balloon, its levity will increase. But if hydrogen and flame are added to a balloon, there is a non-additive interaction, such that levity is not increased.

Figure 19.1 displays a simple interaction. Both predictors have only two levels. The abscissa groups the two levels of predictor, and the shading of the bars indicates the two levels of predictor x_2 . All three panels of Figure 19.1 show the same data, but the nature of the interaction is highlighted differently in each panel.

In the left panel of Figure 19.1, the dashed parallelogram indicates the best additive model for the data. The dashed lines indicate the average when the levels of the

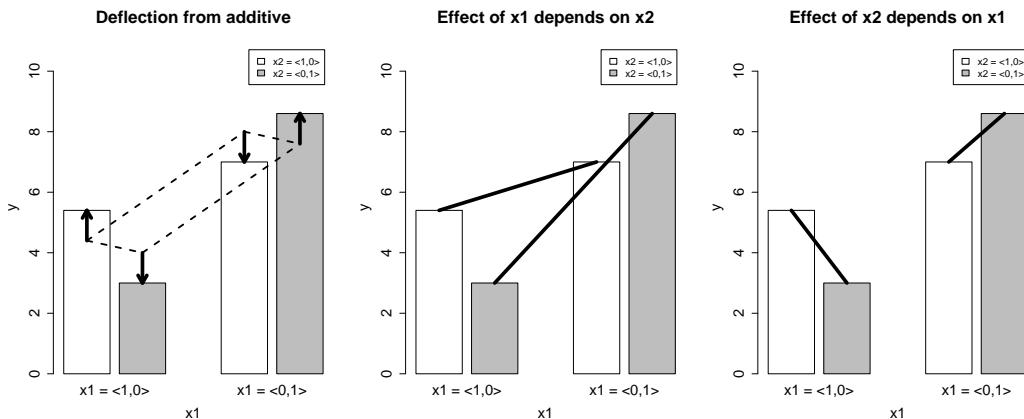


Figure 19.1: An interaction of nominal variables x_1 and x_2 , parsed three ways. The left panel emphasizes that the interaction involves a non-additive, torsion-like deflection away from the additive model, as indicated by the arrows. The middle panel shows the same data, with lines that emphasize that the effect of x_1 depends on the value of x_2 . The right panel again shows the same data, but with lines that emphasize that the effect of x_2 depends on the value of x_1 .

The vertical arrows highlight the non-additive deflections away from the additive average, that constitute the interaction. Notice that the arrows sum to zero across each edge of the parallelogram. Thus, the interaction components do not change the average deflections of each predictor.

The middle and right panels of Figure 19.1 highlight different interpretations of the interaction. The middle panel shows that the effect of x_1 , i.e., the amount that changes when x_1 changes, depends on the level of x_2 : When $x_2 = h1; 0i$, there is only a small change in y when x_1 changes, but when $x_2 = h0; 1i$, there is only a larger change in y when x_1 changes. The right panel makes the same point but with the roles of the predictors reversed: When $x_1 = h1; 0i$, the effect of x_2 is to decrease y , but when $x_1 = h0; 1i$, the effect of x_2 is to increase y .

The average deflection from baseline due to a predictor is called the main effect of the predictor. The main effects of the predictors correspond to the dashed lines in the panel of Figure 19.1. When there is non-additive interaction between predictors, the effect of one predictor depends on the level of the other predictor. The deflection from baseline for a predictor, at a fixed level of the other predictor, is called the simple effect of the predictor at the level of the other predictor. When there is interaction, the simple effects do not equal the main effect.

It may be edifying to compare Figure 19.1, which shows interaction of nominal predictors, with Figure 17.8, p. 384, which shows interaction of metric predictors. The essential notion of interaction is the same in both cases: Interaction means the non-additive portion of the prediction, and interaction means that the effect of one predictor depends on the level of the other predictor.

The mathematical formalism for non-additive interaction was introduced in Section 14.1.6.3, p. 303, as is repeated here. The non-additive components, indicated by the vertical arrows in Figure 19.1, are denoted $\alpha_{2;j;k}$, which means the interaction of predictors 1 and 2 (denoted 12) at level j of predictor 1 and level k of predictor 2. The formal

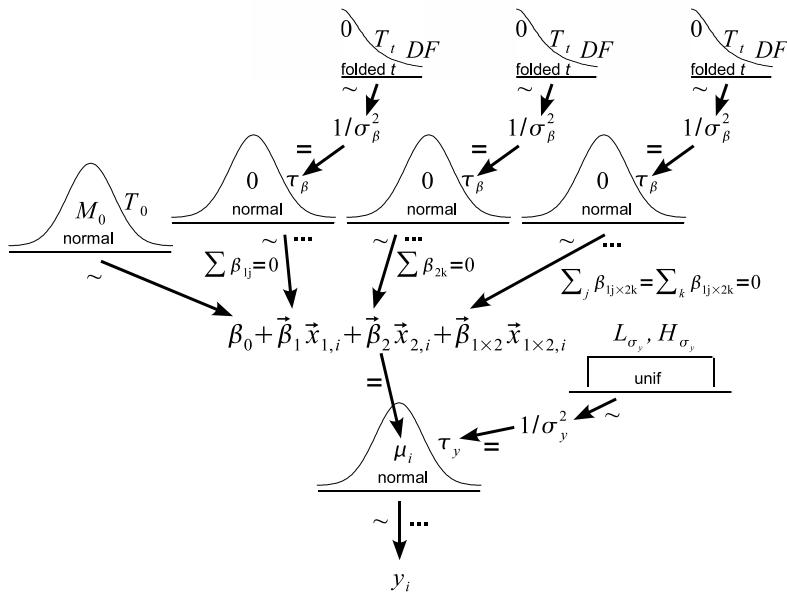


Figure 19.2: Hierarchical dependencies for model of two-way Bayesian ANOVA. Compare with Figure 18.1.

expression merely expands the additive model with by including the interaction. Recall from Equations 14.9 and 14.10 that the model with interaction can be written as

$$\begin{aligned} y &= \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_{12} x_{12} \\ &= \beta_0 + \sum_{j=1}^{1;j} \beta_{1;j} x_{1;j} + \sum_{k=1}^{2;k} \beta_{2;k} x_{2;k} + \sum_{j=1}^{1;2;j;k} \beta_{12;j;k} x_{12;j;k} \end{aligned}$$

with the constraints

$$\sum_{j=1}^{1;j} \beta_{1;j} = 0 \quad \text{and} \quad \sum_{k=1}^{2;k} \beta_{2;k} = 0 \quad \text{and} \quad \sum_{j=1}^{1;2;j;k} \beta_{12;j;k} = 0 \quad \sum_{k=1}^{2;k} \beta_{12;j;k} = 0$$

In those last equations, the symbol “ \forall ” means “for all”. In words, the last two equations simply mean that the interaction de ections sum to zero at every level of the two predictors. A graphic example of this was shown in the left part of Figure 19.1, where it can be seen that the heights of the arrows sum to zero along edges of the parallelogram.

Our goal is to estimate the additive and interactive de ects, based on the observed data. It is important to understand that the observed data are not the bars in Figure 19.1; instead, the data are swarms of points at various heights the heights of the bars. The bars represent the central tendency of the data at each combination of predictors. Thus, what the equations above actually predict is the central tendency at each combination of predictors, and the data are typically modeled as being uniformly distributed around.

19.1.2 The hierarchical prior

The complete generative model of the data is shown in Figure 19.2. It might look daunting, but it really is merely the diagram for oneway ANOVA, in Figure 18.1, with the hyperprior replicated for each predictor and interaction.

The lowest level of Figure 19.2 indicates that the observed points, y_i , are distributed normally around the predicted value μ_i . Moving upward in the diagram, the arrow impinging upon μ_i indicates that the predicted value is baseline plus additive effect due to each predictor plus interactive de effect due to the combination of predictors. The upper levels of the diagram indicate prior structural assumptions about de effects. We assume that the de effects produced by a predictor are centered at zero, we allow the variance (i.e., precision) of the de effects to be estimated from the data, if most of the de effects are small, the estimated variance is small, and the hyperdistribution creates shrinkage in the estimates of other de effects.

A key conceptual aspect of the hyperdistributions is that they apply separately to the different predictors and interactions. In other words, there is just one hyperdistribution that governs all de effects for all predictors and interactions. This division of generative structure reflects a prior assumption that the magnitude of the effect of one predictor might not be very informative regarding the magnitude of the effect of a different predictor. But, within a predictor, the magnitude of de effect produced by one level may inform the magnitude of de effect produced by other levels of that same predictor.¹

As was assumed in the case of oneway ANOVA, we will assume homogeneity of variance: The variability of the observed data is the same within each combination of predictors. This is indicated in Figure 19.2 by the single parameter σ_y that is used in likelihood function, regardless of the values of the predictors. As noted, there are two reasons for this assumption. First, the assumption is a natural simplification multiple regression on metric predictors, and ANOVA can be construed as a special case of multiple regression. Second, the assumption of equal variances is made in NHST ANOVA, and will also make it here in BANOVA to facilitate comparing across the techniques. There is no requirement in BANOVA to assume equal variances. If the situation suggests that different levels of the predictors produce radically different variances in the data, then the hierarchical prior can allow different variances.

19.1.3 An example in R and BUGS

Figure 19.3 shows the mean annual salaries of faculty in four departments at three levels of seniority. The four departments are business, finance, psychology, and educational psychology, chemistry, and theater. These departments are nominal levels of a predictor denoted dx_1 . The three levels of seniority are full professor, associate professor, and assistant professor. Assistant professors are usually within five years after completing their doctoral or post-doctoral studies. Associate professors are usually within about ten years of their doctoral or post-doctoral studies. Full professors are anywhere from 10 to 40 years post graduate school. Although seniority could, and perhaps should, be treated as an ordinal variable, we will treat it as a nominal predictor, denoted dx_2 . A glance at the means suggests that there are effects of department and of seniority. There appears also an interaction, meaning that the change in salary due to seniority depends on department. Our goal is to estimate the baseline salary, the main effect of department membership, the main effect of seniority, and the interaction of department and seniority.

The display of the means in Figure 19.3 obscures the fact that combinations of

¹By analogy to multiple regression, if there are many predictors included in a model, it is reasonable in principle to include a higher-level distribution across predictors such that the estimated variance of one predictor informs the estimated variance of another predictor. This would be especially useful if the application includes many nominal predictors, each with many levels. Such applications are rare.

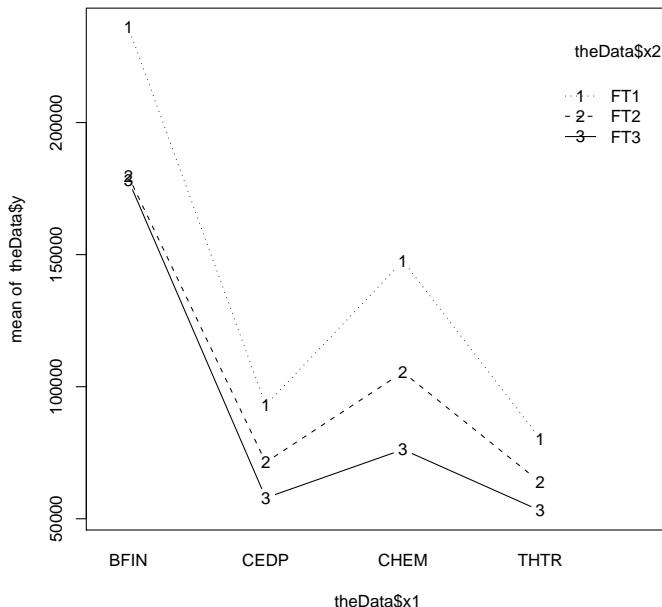


Figure 19.3: Mean annual salaries of faculty in four departments at three levels of seniority.

department and seniority had different numbers of data points. In other words, the number of associate professors in business finance was not necessarily equal to the number of full professors in theater. In traditional NHST ANOVA, this soft “unbalanced” design can cause serious computational difficulties (e.g., Maxwell & Delaney, 2004, pp. 320–343). But Bayesian ANOVA has no problem with unbalanced designs.

The model of Figure 19.2 was implemented in R and BRugs and is described in Section 19.3.1 ANOVAtwowayBRugs.R. Several tricks for running the model in BUGS are described in that section, before the program listing. The results, however, are much like the oneway ANOVA model of the previous chapter.

The results are shown in Figure 19.4. (The means and HDIs are displayed with only three significant digits, but more precise values can be obtained directly from the program.) The top left histogram shows that the baseline for these four departments is 111,381. Notice, however, that most of the data fall below this baseline because the overall data are skewed by the much higher salaries in one department. For salaries in the department of chemistry, the fourth histogram in the top row indicates that 2,164 should be subtracted from the baseline. For salaries of assistant professors, the first histogram in the bottom row indicates that 20,100 should be subtracted from the baseline. Thus, for assistant professors in the department of chemistry, the linearly-predicted salary is $111,381 - 2,164 - 20,100 = 89,117$. But there is a notable non-linear interaction component for that combination: The fourth histogram of the bottom row shows that 10,938 must be subtracted from the linear combination to get the mean for that combination, namely, 78,179.

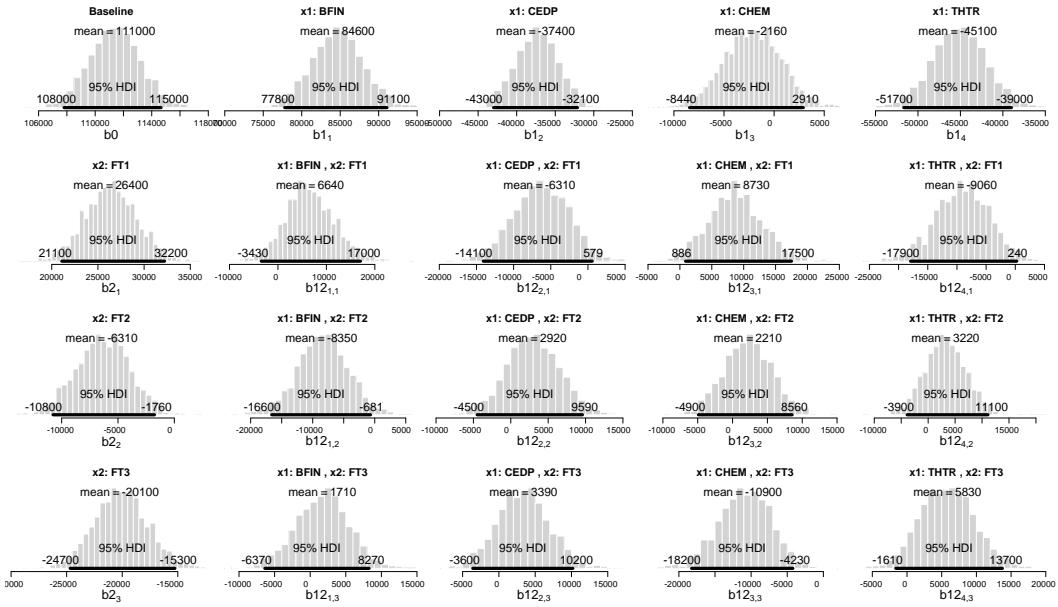


Figure 19.4: Posterior distribution for data in Figure 19.3. Baseline (x_0) is shown in upper left. Remainder of top row is main effect of x_1 (department). Remainder of left column is main effect of x_2 (seniority). Remaining cells show the interaction effects.

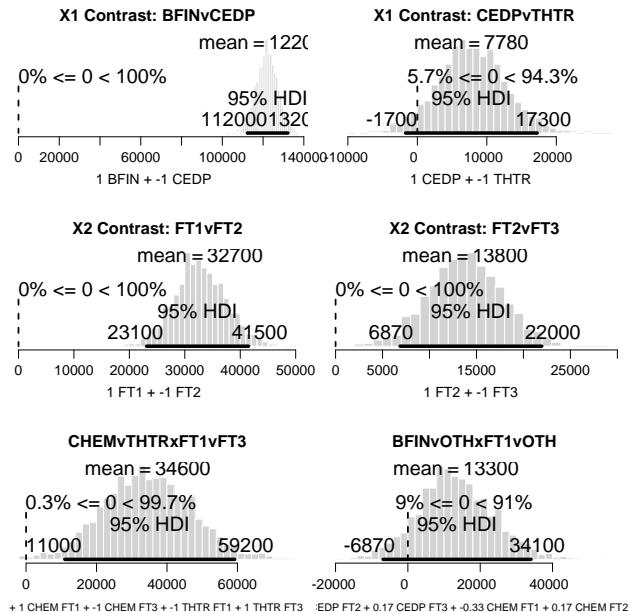


Figure 19.5: Selected contrasts for posterior in Figure 19.4.

19.1.4 Interpreting the posterior

In most applications, we are interested not only in estimated effects for each group, but we are also interested in deciding whether two groups are different. Just as we compared groups in oneway ANOVA in the previous chapter, we compare groups in multifactor ANOVA.

The top and middle rows of Figure 19.5 show selected contrast levels of the main effects. We may ask whether there is a credible difference in salaries, on average, between business finance (BFIN) and counseling and educational psychology (CEDP). The top left histogram indicates that the average difference is about \$122,000, and the 95% HDI falls far from zero. We may also ask whether there is a credible difference in salaries, on average, between CEDP and theater (THTR). The top right histogram indicates that the average difference is about \$7,780, but the 95% HDI spans zero, which indicates that the uncertainty in the estimated difference is fairly large relative to the estimated difference itself. The middle row of Figure 19.5 shows contrasts regarding levels of seniority: There is a credible difference between full professors (FT1) and associate professors (FT2), and between FT2 and assistant professors (FT3).

It is important to understand that the main effects of department and seniority are average effects, when the other factors are collapsed. For example, the contrast between FT2 and FT3 (middle row, right panel of Figure 19.5) is the average difference between FT2 and FT3, collapsed across all departments. But if you look at the data in Figure 19.3, you can see that the difference between FT2 and FT3 is not the same in every department. There is a fairly large difference in CHEM, but a very small difference in BFIN. The effect of changing from FT2 to FT3 depends on the department, which means that there is an interaction.

Main effects must be interpreted and described cautiously when there are interactions. It would be a mistake to say that “the” difference between FT2 and FT3 is 13,800. Instead, that is the average difference across departments. The actual difference within any particular department might be quite different. Similarly, it would be a mistake to say that “the” effect of FT3 is to subtract 20,100 from baseline, because there is an interaction of seniority with department.

19.1.4.1 Metric predictors and ANCOVA

Consider again the salary data in Figure 19.3. You can see that the mean salary for FT1's in Chemistry is much higher than in Theater. This difference might be attributable solely to being in one department or the other. But the difference might also be attributable to some other factor, such as years on the job. In other words, FT1's in Chemistry might happen to have been employed for decades, while the FT1's in Theater might happen to be relatively young. If we had the age of each employee, or, better yet, the number of years that the employee had been at the current level of seniority, we could include that information as an additional predictor of salary. We could then assess whether department membership contributed any predictiveness beyond number of years on the job.

When a nominal predictor, such as department membership, is combined with a metric predictor, such as years on the job, the model is sometimes referred to as analysis of covariance, or ANCOVA. The metric predictor is the “covariate”

Programming ANCOVA in BUGS is a trivial combination of the tools we've used for linear regression and ANOVA. Denote the nominal group membership for individual $asxNom[i]$, and denote the metric covariate value $asxCov[i]$. Then the core of the BUGS

model specification is

```
mu[i] <- a0 + a[ xNom[i] ] + bMet * xMet[i]
y[i]      dnorm( mu[i] , tau )
```

where a_0 is the deviation of each group from baseline, and a is the regression coefficient on the covariate.

19.1.4.2 Interaction contrasts

Just as we can ask whether differences among particular levels of predictors are credible, we can ask whether interactions among particular combinations of predictors are credible. Consider again the data in Figure 19.3. The difference between full professors (FT1) and assistant professors (FT3) appears to be large in the chemistry department (CHEM) but smaller in the theater department (THTR). Is the simple effect of seniority bigger in chemistry than it is in theater? In other words, is $(\text{CHEM} \cdot \text{FT1} - \text{CHEM} \cdot \text{FT3}) - (\text{THTR} \cdot \text{FT1} - \text{THTR} \cdot \text{FT3})$ credibly non-zero?

This sort of difference of differences is called an interaction contrast. In general, an interaction contrast is constructed by taking any set of coefficients on x_1 , and any set of contrast coefficients on x_2 , and computing their outer product. The outer product was described in Section 8.8.8 (TwoGrid.R), p. 144. Formally, the outer product of two vectors is denoted by the symbol “ \cdot ”. To provide an example of an interaction contrast as an outer product of main effects contrasts, we will re-cast the one we are presently considering, namely $(\text{CHEM} \cdot \text{FT1} - \text{CHEM} \cdot \text{FT3}) - (\text{THTR} \cdot \text{FT1} - \text{THTR} \cdot \text{FT3})$, in generic notation. Notice that CHEM is level 3 of predictor hence can be written as $x_{1;3}$. Writing the other components in the same fashion, the interaction contrast is $(x_{1;3} \cdot x_{2;1} - x_{1;3} \cdot x_{2;3}) - (x_{1;4} \cdot x_{2;1} - x_{1;4} \cdot x_{2;3})$. That can be algebraically re-arranged to highlight the coefficients on the particular combinations:

$$(+1)x_{1;3} \cdot x_{2;1} + (-1)x_{1;3} \cdot x_{2;3} + (-1)x_{1;4} \cdot x_{2;1} + (+1)x_{1;4} \cdot x_{2;3}$$

Those highlighted coefficients can be obtained as the outer product of main effects contrasts, namely the contrast $c_1 = h; 0; +1; -1$, which expresses CHEM minus THTR, and the contrast $c_2 = h+1; 0; -1$, which expresses FT1 minus FT3:

$$\begin{aligned} c_1 \cdot c_2 &= \begin{matrix} ! & x_{1;1} & ! & x_{1;2} & ! & x_{1;3} & ! & x_{1;4} \\ h & 0 & 0 & +1 & 1 & i & N & h \\ & & & & & & & +1 \\ & & & & & & & 0 \\ & & & & & & & 1 \\ & & & & & & & i \end{matrix} \\ &= \begin{matrix} ! & x_{1;1} & ! & x_{2;1} & ! & x_{2;2} & ! & x_{2;3} \\ ! & x_{1;2} & 2 & 0 & 0 & 0 & 3 \\ ! & x_{1;3} & +1 & 0 & 0 & 1 & 5 \\ ! & x_{1;4} & 1 & 0 & 0 & +1 & 6 \end{matrix} \end{aligned}$$

Notice that the coefficients in the matrix match the highlighted coefficients in the difference of differences that was expressed a few sentences previously. The prior of this interaction contrast is shown in the bottom left histogram of Fig 19.5. The mean of 34,600 indicates that the difference between FT1 and FT2 is about 34,600 greater for CHEM than for THTR. The 95% HDI clearly excludes zero, indicating that this interaction contrast is credibly non-zero.

Interaction contrasts can involve “complex” comparisons just as simply as pairwise comparisons. For example, suppose we are interested in comparing BFIN against the

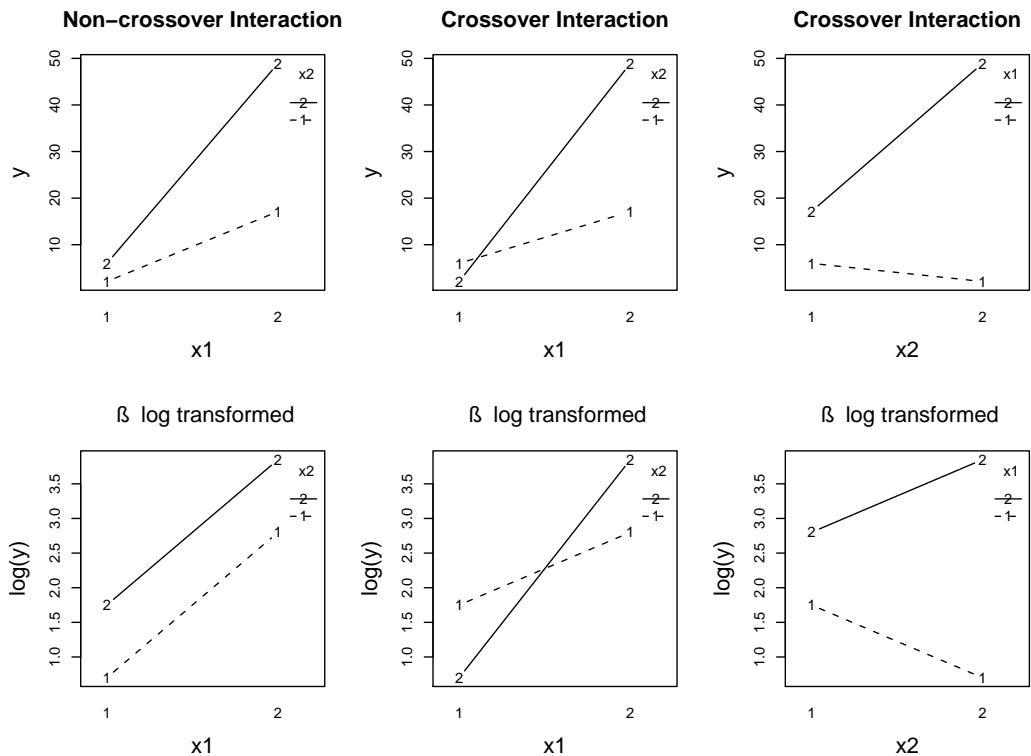


Figure 19.6: Top row shows means of original data; bottom shows means of logarithmically transformed data. Left column shows a non-crossover interaction; middle and right columns show the same crossover interaction plotted against x_1 or x_2 .

average of the other, non-business departments, specify for a contrast between FT1 and lesser ranks. This interaction contrast is expressed as $1=3; 1=3; 1=3i$, $1+1; 1+2; 1=2i$. The posterior of this contrast is shown in the bottom right histogram of Figure 19.5. (The label on the x-axis exceeds the margin of the figure because there are twelve combinations of levels involved in the contrast.) The result suggests that the larger difference, between FT1 and other ranks, in BFIN than in other departments has substantial uncertainty. Therefore we would not want to conclude that the interaction contrast is credibly non-zero. Exercise 19.2 gives you hands-on practice with specification of interaction contrasts.

19.1.5 Non-crossover interactions, rescaling, and homogeneous variances

When interpreting interactions, it can be important to consider the scale on which the data are measured. This is because an interaction means non-additivity of effects when measured in the current scale. If the data are nonlinearly transformed to a different scale, then the non-additivity can also change.

Consider an example, using utterly fictional numbers merely for illustration. Suppose the average salary of Democratic women is 10 monetary units, Democratic men it's 12 units, for Republican women it's 15 units, and for Republicans it's 18 units. These data indicate that there is a non-additive interaction of party and gender, because the

change in pay from women to men is 2 units for Democrats, but 3 for Republicans. Another way of describing the interaction is to notice that change in pay from Democrats to Republicans is 5 units for women but 6 units for men. A reader might be tempted to interpret the interaction as indicating some extra advantage retained by Republican men, or some special disadvantage faced by Democratic women. But such an interpretation may be inappropriate, because a mere rescaling of the data makes the interaction disappear, as will be described next.

Salary increases and comparisons are often measured by ratios, not by additive or subtractive differences. Consider the salary data in percentage terms. Among Democrats, men make 20% more than women. Among Republicans, men again make 20% more than the women. Among women, Republicans make 50% more than Democrats. Among men, Republicans again make 50% more than Democrats. In ratio terms, there is no interaction of gender and political party: Change from female to male predicts a 20% increase in salary regardless of party, and change from Democrat to Republican predicts a 50% increase in salary regardless of gender.

Equal ratios are transformed to equal distances by a logarithmic transformation. If we measure salary in terms of the logarithm of monetary units, the salary of Democratic women is $\log_0(10) = 1.000$, the salary of Democratic men is $\log(12) = 1.079$, the salary of Republican women is $\log(15) = 1.176$, and the salary of Republican men is $\log_{10}(18) = 1.255$. With this logarithmic scaling, the increase in salary of women to men is 0.079 for both parties, and the increase from Democrat to Republican is 0.176 for both genders. In other words, when salary is measured on a logarithmic scale, there is no interaction of gender and political party.

It may seem strange to measure salary on a logarithmic scale, there are many situations for which the scale is arbitrary. The pitch of a sound can be measured in terms of frequency (i.e., cycles per second), or in terms of *perceptible*, which is essentially the logarithm of the frequency. The magnitude of an earthquake can be measured by its energy, or by its value on the Richter scale, which is the logarithm of energy. The pace of a dragster on a race track can be measured by the average speed during the race, or by the duration from start to finish (which is the reciprocal of average speed). Thus, measurement scales are not unique, and are instead determined by convention.

The general issue is illustrated in Figure 19.6. Suppose predictor x_1 has two levels, as does predictor x_2 . Suppose we have three data points at each combination of levels, yielding twelve data points altogether. The means at each combination of levels are shown in the top-left graph of Figure 19.6. You can see that there is an interaction, with the effect of x_1 being bigger when $x_2 = 2$ than when $x_2 = 1$. But this interaction goes away when the data are transformed by taking the logarithm, as shown in the lower-left graph. Each individual data point was transformed, and then the mean for each cell were computed. Of course, the transformation can go the other way: Data with no interaction, as in the lower-left plot, can be made to have an interaction when they are scaled as in the upper-left plot, via an exponential transformation.

The transformability from interaction to non-interaction is only possible for non-crossover interactions. This terminology is merely a description of a graph: The lines do not cross over each other (and they have the same sign). In this situation, they-axis can have different portions stretched or shrunken so that the lines become parallel. If, however, the lines cross, as in the middle column of Figure 19.6, then there is no way to uncross the lines merely by stretching or shrinking them on the y-axis. The right column of Figure 19.6 shows the same data as the middle column, but with the roles of x_1

and x_2 exchanged. When plotted this way, the lines do not cross, ~~but they do have opposite-sign slopes (i.e., one slope is positive and the other slope is negative)~~. There is no way that stretching or shrinking the y -axis can change the signs of the slopes, hence the interaction cannot be removed merely by transforming the data. Because data have crossing lines when plotted as in the middle column, they are said to have a crossover interaction even when they are plotted as in the right column. (Test your understanding: Is the interaction in Figure 19.1 a crossover interaction?)

It is important to note that the transformation applies to individual raw data values, not to the means of the conditions. A consequence of transforming the data, therefore, is changes in the variances of the data within each condition. For example, suppose one condition has data values of 100, 110, and 120, while a second condition has data values of 1100, 1110, and 1120. For both conditions, the variance is ~~the same~~, there is homogeneity of variance. When the data are logarithmically transformed, the variance of the first group becomes 1.05×10^{-3} , but the variance of the second group becomes two orders of magnitude smaller, namely 1.02×10^{-5} ; i.e., there is not homogeneity of variance.

Therefore, when applying the hierarchical model of Figure 19.2, we must be aware that it assumes homogeneity of variance. If we transform the data, we are changing the variances within the levels of the predictors. The transformed variances might or might not be fairly homogeneous. If they are not, then either the data should be transformed in such a way as to respect homogeneity of variance, or the model should be changed to allow unequal variances.

The models we have been using also assume a normal likelihood function, which means that the data at any level of the predictors should be normally distributed. When the data are transformed to a different scale, the shape of their distribution also changes. If the distributions become radically non-normal, it may be misleading to use a model with a normal likelihood function. For a discussion of these issues, see Section 15.1.4, p. 326.

In summary, this section has made two main points. First, if you have a non-crossover interaction, be careful what you claim about it. A non-cross interaction merely means non-additivity in the scale you are using. If this scale is the only meaningful scale, or if this scale is the overwhelmingly dominant scale used in a field of research, then you can cautiously interpret the non-additive interaction with respect to that scale. But if transformed scales are reasonable, then keep in mind that they are scale-specific, and there might be no interaction in a different scale. With a crossover interaction, however, no rescaling can undo the interaction. Second, non-linear transformations change the within-cell variances and the shapes of the within-cell distributions. Be sure that the model you are using is appropriate to the homogeneity or non-homogeneity of variances in the data, and to the shapes of the distributions, on whatever scale you are using. Exercise 19.1 has you consider these issues “hands on”.

19.2 Repeated measures, a.k.a. within-subject designs

In many situations, a single “subject” contributes data to multiple levels of the predictors. For example, suppose we are interested in how quickly people press a button in response to a stimulus onset. The stimulus could appear in the visual modality as a light, or in the auditory modality as a tone. The subject could respond with their dominant hand, or with his/her non-dominant hand. Thus, there are two nominal predictors, namely modality and hand. The new aspect is that a single subject contributes to all combinations

of the predictors. On many successive trials, the subject either a tone or light, and is instructed to respond with either the dominant or non-dominant hand. Because every subject is measured many times, this situation is sometimes called a “repeated measures” design. Because the levels of the predictors change with subjects, this situation is also called a “within subject” design. I favor the latter terminology because it more explicitly connotes the essential aspect of the design, that the subject contributes data in more than one condition. Within-subject designs are contrasted with “between-subject” designs, in which different subjects contribute data to different levels of the predictors.

When every subject contributes many data points to every combination of predictors, then the model of the situation is a straight forward extension of the models we've already considered. We merely add “subject” as another predictor in the model, with each individual subject being a level of the predictor. If there is one predictor other than subject, the model becomes

$$y = \beta_0 + \beta_1 x_1 + \beta_s x_s + \beta_{1s} x_{1s}$$

This is exactly the two-predictor model we have already discussed, with the second predictor being subject. When there are two predictors other than subject, the model becomes

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_s x_s + \beta_{12} x_{12} + \beta_{1s} x_{1s} + \beta_{2s} x_{2s} + \beta_{12s} x_{12s}$$

This model includes all the two-way interactions of the first plus the three-way interaction. Again, subject merely plays the role of the third predictor.

The model above, that includes all the high-order interactions with subject, is fine in principle but may be overkill in practice. Unless you have specific theoretical motivations to seek out and interpret high-order interactions of subject with other predictors, there is little reason to model them, and difficulty making sense of them even if you did model them. Instead, if you have many data points from each subject in every cell, an alternative approach is to apply a Bayesian ANOVA model to each subject's data, and then put a higher-order prior across the subject parameter estimates, so that different subjects mutually inform each other's estimates and provide a stable group-level estimate. Thus, every subject has a baseline, β_0 , and there is a higher-order, group-level prior on the distribution of β_0 s across subjects. Each predictor also has subject-specific estimates, with the effect of the j^{th} level of predictor 1 denoted β_{1s_j} . Each of these effect parameters has a higher, group-level prior across subjects. (This was the modeling approach taken for repeated measures in simple linear regression in Section 16.3, p. 354.) Finally, the group-level effects have a hyperprior that provides shrinkage on the effects of a predictor. In other words, the shrinkage prior, on the effects of a predictor, is set at the group level, not at the subject level.

There are other situations, however, in which each subject contributes only one datum to a combination of the other predictors. For example, in the response-time study described above, perhaps we have only the median response of the subject in each combination of hand and modality. As another example, suppose the value to be predicted is IQ, as measured by a lengthy exam, with one predictor being noisy versus quiet exam environment, and the other predictor being paper versus computerized exam format. Although it is conceivable that subjects could be repeatedly tested in each condition, it would be challenging enough to get people to sit through all combinations even once. Thus, each subject would contribute one value to each condition.

In the situation when each subject contributes only one datum per condition, the models described above, with all the interaction terms, are not identifiable”, meaning that there are

more parameters than data points. The simplest case of this is trying to estimate the mean and variance of a normal distribution from a single point. A Bayesian analysis can still be conducted, but there will be high uncertainty in the parameter estimates, governed largely by the priors. Therefore, instead of attempting to estimate all the interactions of subjects with other predictors, we assume a simple model in which the only influence of subjects is on the baseline:

$$y = \beta_0 + \beta_S x_S + \beta_1 x_1 + \beta_2 x_2 + \beta_{12} x_{12}$$

In other words, we assume a main effect of subject, but no interaction of subject with other predictors. In this model, the subject effect (deception) is constant across treatments, and the treatment effects (deceptions) are constant across subjects. Notice that the model makes no requirement that every subject contributes a datum. This is a dummy condition. Indeed, the model allows zero or more than one datum per subject per treatment. As mentioned earlier, the computations in Bayesian ANOVA make no assumptions about requirements that the design is “balanced”. If you do have many observations per subject in every combination of predictors, then one of the previously described models may be considered.

19.2.1 Why use a within-subject design? And why not?

The primary reason to use a within-subject design is that you achieve much greater precision in the estimates of the effects than in a between-subject design. For example, suppose you are interested in measuring the mean response time of using the dominant versus non-dominant hand. Suppose there is a population of subjects from whom you could measure data. If we could measure every subject in every condition, we would know that for the first subject, his or her response times for dominant and non-dominant hands are 300 and 320 msec. For the second subject, the response times are 350 and 370. For the third subject, the response times are 400 and 420. For the fourth subject, the response times are 450 and 470. Thus, for every subject, the difference between dominant and non-dominant hands is 20 msec. Suppose we have the resources to measure only two data points in each condition. We measure response times from the dominant hands of two subjects. Should we measure response times from the non-dominant hands of the same two subjects, or the non-dominant hands of two other subjects? If we measure from the same two subjects, then the estimated effect for each subject is 20 msec, and we have high certainty in the magnitude of the effect. If we measure from two other subjects, then the estimated effect of dominant versus non-dominant hand is the average of the two subjects versus the average of the second two subjects, and the difference is badly affected by random sampling. The between-subject design yields low precision in the estimate of the effect. Exercise 19.3 has you examine, hands on, a case of improvement in precision.

Because of the gain in precision, it is desirable to use within-subject designs. But there are many dangers of within-subject designs that need to be considered before they are applied in any particular situation. The key problem is that in most situations, when you measure the subject you change the subject, and therefore subsequent measurements are not measuring the same subject. The simplest example is where mere fatigue or generic practice affects. In measures of response time, if you measure repeatedly the same subject, you will find improvement over the first several trials because of the subject gaining practice with the task, but after a while the subject tires, there will be a

decline in performance. The problem is that if you measured the dominant hand in the early trials, and the nondominant hand in the later trials, then the effect of practice or fatigue will contaminate the effect of handedness. The repeated measurement process and contaminates the measure that is supposed to be a significant predictor.

Practice and fatigue effects can be overcome by randomly distributing and repeating the conditions throughout the repeated measures. The practice and fatigue effects in influence all conditions equally. Thus, if practice improves the dominant and nondominant hand by 50 msec, then the difference between dominant and nondominant hands is unaffected by practice. But practice might affect the nondominant hand much more than the dominant hand. You can imagine that in complex designs with many predictors, each with many levels, it can become difficult to justify an assumption that repeated measures have comparable effects on all conditions.

Worse yet, in some situations there can be differential carryover effects from one condition to the next. For example, having just experienced a task in the visual modality with the nondominant hand might improve subsequent performance in the auditory modality with the nondominant hand, but might not improve subsequent performance in the visual modality with the dominant hand. Thus, the carryover effect is different for different subsequent conditions.

When you suspect strong differential carryover effects, you may be able to explicitly manipulate the ordering of the conditions and measure those effects, but this might be impossible mathematically and impractical, depending on the specifics of your situation. In this case, you must revert to a between subjects design and simply include many subjects to attenuate between-subject noise.

In general, all the models we have been using assume independence of observations. The probability of the collection of data is the product of probabilities of the individual data points. When we use repeated measures, this assumption is much less easy to justify. On the one hand, when we repeatedly flip a coin, we might be tempted to assume that its underlying bias does not change much from one flip to the next. But, on the other hand, when we repeatedly test the response time of a human subject, it is less easy to justify an assumption that the underlying response time remains constant by the previous trial. Researchers will often make the assumption of independence merely as an approximation of convenience, hoping that by arranging conditions randomly across many repeated measures, the differential carryover effects will be minimized.

19.3 R code

19.3.1 Bayesian two-factor ANOVA

Several implementation details of the program are the same as the oneway ANOVA program of the previous chapter:

Data are normalized so that prior constants can be more generic

Initialization of chains is based on the data. It is important to do this, otherwise burn-in can take forever.

Because there is nasty autocorrelation, we use a large number of chains and we also use multiple chains. For a reminder of the issues of burn-in and thinning, see Section 23.2, p. 510.

A new detail arises in how the uncentered parameter estimates are recentered to respect the sum-to-zero constraints. The uncentered estimates in BUGS are a_0 , $a_1[i]$, $a_2[j]$, and $a_{12}[i,j]$. By definition of the ANOVA model, the predicted mean of cell i,j is

$$m[i,j] = a_0 + a_1[i] + a_2[j] + a_{12}[i,j]$$

We use these predicted means to construct the zero-centered parameters. First b_0 is the mean across all the predicted means:

$$b_0 = \text{mean}(m[,])$$

Then the main effects are the marginal means minus the overall mean

$$\begin{aligned} b_1[i] &= \text{mean}(m[i,]) - b_0 \\ b_2[j] &= \text{mean}(m[,j]) - b_0 \end{aligned}$$

It is easy (honest!) to check that the those effects do sum to zero; i.e. $\sum(b_1) = 0$ and $\sum(b_2) = 0$. Finally, the interaction effects are the residuals after the additive effect of b_1 and b_2 is taken into account:

$$b_{12}[i,j] = m[i,j] - (b_0 + b_1[i] + b_2[j])$$

Again, it is easy to check that the rows and columns all sum to zero.

In the data section of the program, one option is to load data the article of Qian and Shen (2007). The program here uses a hierarchical structure similar to that used by Qian and Shen (2007), but their program did not re-center the parameters as is done here. It may be instructive to compare the results of the program with the results reported by Qian and Shen (2007).

BUGS for many factors. The program below applies only for cases of two nominal predictors. If you have many nominal predictors, along with two-way, three-way, and higher-order interactions, it becomes unwieldy to identify and separately name all the effect parameters. Instead, it can be more elegant to use a technique of “dummy coding”, whereby we essentially revert back to using vectors for coding the values of the predictors instead of integer indices. That is, level 2 is coded by the “dummy” vector $h[0; 1; 0; :; i]$ instead of by the integer index 2. Interactions are represented by matrices of dummy codes that have been attened into vectors. For an example of programming this technique in BUGS, see Ntzoufras (2009, Ch. 6). Unfortunately, those examples do not incorporate the higher-level prior structure emphasized in Figure 19.2.

(ANOVAtwowayBRugs.R)

```

1 graphics.off()
2 rm(list=ls(all=TRUE))
3 fnroot = "ANOVAtwowayBRugs"
4 library(BRugs)      # Kruschke, J. K. (2010). Doing Bayesian data analysis:
5                      # A Tutorial with R and BUGS. Academic Press / Elsevier.
6 #-----
7 # THE MODEL.
8
9 modelstring =
10 # BUGS model specification begins here...
11 model {
12   for ( i in 1:Ntotal ) {
13     y[i] ~ dnorm( mu[i] , tau )
14     mu[i] <- a0 + a1[x1[i]] + a2[x2[i]] + a1a2[x1[i],x2[i]]
}

```

```

15   }
16   #
17   tau <- pow( sigma , -2 )
18   sigma ~ dunif(0,10) # y values are assumed to be standardized
19   #
20   a0 ~ dnorm(0,0.001) # y values are assumed to be standardized
21   #
22   for ( j1 in 1:Nx1Lvl ) { a1[j1] ~ dnorm( 0.0 , a1tau ) }
23   a1tau <- 1 / pow( a1SD , 2 )
24   a1SD <- abs( a1SDunabs ) + .1
25   a1SDunabs ~ dt( 0 , 0.001 , 2 )
26   #
27   for ( j2 in 1:Nx2Lvl ) { a2[j2] ~ dnorm( 0.0 , a2tau ) }
28   a2tau <- 1 / pow( a2SD , 2 )
29   a2SD <- abs( a2SDunabs ) + .1
30   a2SDunabs ~ dt( 0 , 0.001 , 2 )
31   #
32   for ( j1 in 1:Nx1Lvl ) { for ( j2 in 1:Nx2Lvl ) {
33     a1a2[j1,j2] ~ dnorm( 0.0 , a1a2tau )
34   } }
35   a1a2tau <- 1 / pow( a1a2SD , 2 )
36   a1a2SD <- abs( a1a2SDunabs ) + .1
37   a1a2SDunabs ~ dt( 0 , 0.001 , 2 )
38 }
39 # ... end BUGS model specification
40 " # close quote for modelstring
41 # Write model to a file, and send to BUGS:
42 writeLines(modelstring,con="model.txt")
43 modelCheck( "model.txt" )

44 #
45 #-----#
46 # THE DATA.
47 # Specify data source:
48 dataSource = c( "QianS2007" , "Salary" , "Random" , "Ex19.3"      )[4]
49

50 # Load the data:
51 if ( dataSource == "QianS2007" ) {
52   fnroot = paste( fnroot , dataSource , sep="" )
53   datarecord = read.table( "QianS2007SeaweedData.txt" , he      ader=TRUE , sep="," )
54   # Logistic transform the COVER value:
55   # Used by Appendix 3 of QianS2007 to replicate Ramsey and Schafer (2002).
56   datarecord$COVER = -log( ( 100 / datarecord$COVER ) - 1 )
57   y = as.numeric(datarecord$COVER)
58   x1 = as.numeric(datarecord$TREAT)
59   x1names = levels(datarecord$TREAT)
60   x2 = as.numeric(datarecord$BLOCK)
61   x2names = levels(datarecord$BLOCK)
62   Ntotal = length(y)
63   Nx1Lvl = length(unique(x1))
64   Nx2Lvl = length(unique(x2))
65   x1contrastList = list( f_Effect=c( 1/2 , -1/2 , 0 , 1/2 , -1/2 ,          0 ) ,
66                         F_Effect=c( 0 , 1/2 , -1/2 , 0 , 1/2 , -1/2 ) ,
67                         L_Effect=c( 1/3 , 1/3 , 1/3 , -1/3 , -1/3 , -1/3 ) )
68   x2contrastList = NULL # list( vector(length=Nx2Lvl) )
69   x1x2contrastList = NULL # list( matrix( 1:(Nx1Lvl*Nx2Lvl) , nrow=Nx1Lvl ) )
70 }

71 if ( dataSource == "Salary" ) {
72   fnroot = paste( fnroot , dataSource , sep="" )

```

```

74 datarecord = read.table( "Salary.csv" , header=TRUE , sep=      "," )
75 y = as.numeric(datarecord$Salary)
76 if ( F ) { # take log10 of salary
77   y = log10( y )
78   fnroot = paste( fnroot , "Log10" , sep="" )
79 }
80 x1 = as.numeric(datarecord$Org)
81 x1names = levels(datarecord$Org)
82 x2 = as.numeric(datarecord$Post)
83 x2names = levels(datarecord$Post)
84 Ntotal = length(y)
85 Nx1Lvl = length(unique(x1))
86 Nx2Lvl = length(unique(x2))
87 x1contrastList = list( BFINvCEDP = c( 1 , -1 , 0 , 0 ) ,
88                         CEDPvTHTR = c( 0 , 1 , 0 , -1 ) )
89 x2contrastList = list( FT1vFT2 = c( 1 , -1 , 0 ) , FT2vFT3 = c(0,1      ,-1) )
90 x1x2contrastList = list(
91   CHEMvTHTRxFT1vFT3 = outer( c(0,0,+1,-1) , c(+1,0,-1) ) ,
92   BFINvOTHxFT1vOTH = outer( c(+1,-1/3,-1/3,-1/3) , c(+1,-1      /2,-1/2) ) )
93 }
94
95 if ( dataSource == "Random" ) {
96   fnroot = paste( fnroot , dataSource , sep="" )
97   set.seed(47405)
98   ysdtrue = 3.0
99   a0true = 100
100  a1true = c( 2 , 0 , -2 ) # sum to zero
101  a2true = c( 3 , 1 , -1 , -3 ) # sum to zero
102  a1a2true = matrix( c( 1,-1,0 , -1,1,0 , 0,0,0 , 0,0,0 ),# row an      d col sum to zero
103                      nrow=length(a1true) , ncol=length(a2true) , byrow=F )
104  npercell = 8
105  datarecord = matrix( 0, ncol=3 , nrow=length(a1true)*leng      th(a2true)*npercell )
106  colnames(datarecord) = c("y","x1","x2")
107  rowidx = 0
108  for ( x1idx in 1:length(a1true) ) {
109    for ( x2idx in 1:length(a2true) ) {
110      for ( subjidx in 1:npercell ) {
111        rowidx = rowidx + 1
112        datarecord[rowidx,"x1"] = x1idx
113        datarecord[rowidx,"x2"] = x2idx
114        datarecord[rowidx,"y"] = ( a0true + a1true[x1idx] + a2true      [x2idx]
115                               + a1a2true[x1idx,x2idx] + rnorm(1,0,ysdtrue) )
116      }
117    }
118  }
119  datarecord = data.frame( y=datarecord[, "y"] ,
120                         x1=as.factor(datarecord[, "x1"]) ,
121                         x2=as.factor(datarecord[, "x2"]) )
122 y = as.numeric(datarecord$y)
123 x1 = as.numeric(datarecord$x1)
124 x1names = levels(datarecord$x1)
125 x2 = as.numeric(datarecord$x2)
126 x2names = levels(datarecord$x2)
127 Ntotal = length(y)
128 Nx1Lvl = length(unique(x1))
129 Nx2Lvl = length(unique(x2))
130 x1contrastList = list( X1_1v3 = c( 1 , 0 , -1 ) ) #
131 x2contrastList = list( X2_12v34 = c( 1/2 , 1/2 , -1/2 , -1/2 ) ) #
132 x1x2contrastList = list(

```

```

133     IC_11v22 = outer( c(1,-1,0) , c(1,-1,0,0) ) ,
134     IC_23v34 = outer( c(0,1,-1) , c(0,0,1,-1) )
135   )
136 }
137
138 # Load the data:
139 if ( dataSource == "Ex19.3" ) {
140   fnroot = paste( fnroot , dataSource , sep="" )
141   y = c( 101,102,103,105,104, 104,105,107,106,108, 105,107 ,106,108,109, 109,108,110,111,112 )
142   x1 = c( 1,1,1,1,1, 1,1,1,1,1, 2,2,2,2,2, 2,2,2,2,2 )
143   x2 = c( 1,1,1,1,1, 2,2,2,2,2, 1,1,1,1,1, 2,2,2,2,2 )
144   # S = c( 1,2,3,4,5, 1,2,3,4,5, 1,2,3,4,5, 1,2,3,4,5 )
145   x1names = c("x1.1","x1.2")
146   x2names = c("x2.1","x2.2")
147   # Snames = c("S1","S2","S3","S4","S5")
148   Ntotal = length(y)
149   Nx1Lvl = length(unique(x1))
150   Nx2Lvl = length(unique(x2))
151   # NSLvl = length(unique(S))
152   x1contrastList = list( X1.2vX1.1 = c( -1 , 1 ) )
153   x2contrastList = list( X2.2vX2.1 = c( -1 , 1 ) )
154   x1x2contrastList = NULL # list( matrix( 1:(Nx1Lvl*Nx2Lvl) , nrow=Nx1Lvl ) )
155 }
156
157 # Specify the data in a form that is compatible with BRugs mode I, as a list:
158 ySDorig = sd(y)
159 yMorig = mean(y)
160 z = ( y - yMorig ) / ySDorig
161 dataList = list(
162   y = z ,
163   x1 = x1 ,
164   x2 = x2 ,
165   Ntotal = Ntotal ,
166   Nx1Lvl = Nx1Lvl ,
167   Nx2Lvl = Nx2Lvl
168 )
169 # Get the data into BRugs:
170 modelData( bugsData( dataList ) )
171
172 #-----
173 # INTIALIZE THE CHAINS.
174
175 # Autocorrelation within chains is large, so use several cha
176 # degree of thinning. But we still have to burn-in all the chai
177 # more time with more chains.
178 nchain = 10
179 modelCompile( numChains = nchain )
180
181 if ( F ) {
182   modelGenInits() # often won't work for diffuse prior
183 } else {
184   # initialization based on data
185   theData = data.frame( y=dataList$y , x1=factor(x1,labels      =x1names) ,
186                         x2=factor(x2,labels=x2names) )
187   a0 = mean( theData$y )
188   a1 = aggregate( theData$y , list( theData$x1 ) , mean )[,2] - a      0
189   a2 = aggregate( theData$y , list( theData$x2 ) , mean )[,2] - a      0
190   linpred = as.vector( outer( a1 , a2 , "+" ) + a0 )
191   a1a2 = aggregate( theData$y, list(theData$x1,theData$x2      ), mean)[,3] - linpred

```

```

192 genInitList <- function() {
193   return(
194     list(
195       a0 = a0 ,
196       a1 = a1 ,
197       a2 = a2 ,
198       a1a2 = matrix( a1a2 , nrow=Nx1Lvl , ncol=Nx2Lvl ) ,
199       sigma = sd(theData$y)/2 , # lazy
200       a1SDunabs = sd(a1) ,
201       a2SDunabs = sd(a2) ,
202       a1a2SDunabs = sd(a1a2)
203     )
204   )
205 }
206 for ( chainIdx in 1 : nchain ) {
207   modelInits( bugsInits( genInitList ) )
208 }
209 }
210
211 #-----
212 # RUN THE CHAINS
213
214 # burn in
215 BurnInSteps = 10000
216 modelUpdate( BurnInSteps )
217 # actual samples
218 samplesSet( c( "a0" , "a1" , "a2" , "a1a2" ,
219               "sigma" , "a1SD" , "a2SD" , "a1a2SD" ) )
220 stepsPerChain = ceiling(2000/nchain)
221 thinStep = 500 # 750
222 modelUpdate( stepsPerChain , thin=thinStep )
223
224 #-----
225 # EXAMINE THE RESULTS
226
227 source("plotChains.R")
228 source("plotPost.R")
229
230 checkConvergence = F
231 if ( checkConvergence ) {
232   sumInfo = plotChains( "a0" , saveplots=F , filenameroot=fn      root )
233   sumInfo = plotChains( "a1" , saveplots=F , filenameroot=fn      root )
234   sumInfo = plotChains( "a2" , saveplots=F , filenameroot=fn      root )
235   sumInfo = plotChains( "a1a2" , saveplots=F , filenameroot=fnroot )
236   readline("Press any key to clear graphics and continue")
237   graphics.off()
238   sumInfo = plotChains( "sigma" , saveplots=F , filenameroot=fnroot )
239   sumInfo = plotChains( "a1SD" , saveplots=F , filenameroot=fnroot )
240   sumInfo = plotChains( "a2SD" , saveplots=F , filenameroot=fnroot )
241   sumInfo = plotChains( "a1a2SD" , saveplots=F , filenameroot=t=fnroot )
242   readline("Press any key to clear graphics and continue")
243   graphics.off()
244 }
245
246 # Extract and plot the SDs:
247 sigmaSample = samplesSample("sigma")
248 a1SDSample = samplesSample("a1SD")
249 a2SDSample = samplesSample("a2SD")
250 a1a2SDSample = samplesSample("a1a2SD")

```

```

251 windows()
252 layout( matrix(1:4,nrow=2) )
253 par( mar=c(3,1,2,5,0) , mgp=c(2,0.7,0) )
254 plotPost( sigmaSample , xlab="sigma" , main="Cell SD" , bre     aks=30 )
255 plotPost( a1SDSample , xlab="a1SD" , main="a1 SD" , breaks= 30 )
256 plotPost( a2SDSample , xlab="a2SD" , main="a2 SD" , breaks= 30 )
257 plotPost( a1a2SDSample , xlab="a1a2SD" , main="Interacti    on SD" , breaks=30 )
258 dev.copy2eps(file=paste(fnroot,"SD.eps",sep=""))
259
260 # Extract a values:
261 a0Sample = samplesSample( "a0" )
262 chainLength = length(a0Sample)
263 a1Sample = array( 0 , dim=c( datalist$Nx1Lvl , chainLength )      )
264 for ( x1idx in 1:datalist$Nx1Lvl ) {
265   a1Sample[x1idx,] = samplesSample( paste("a1[",x1idx,"]",",sep="") )
266 }
267 a2Sample = array( 0 , dim=c( datalist$Nx2Lvl , chainLength )      )
268 for ( x2idx in 1:datalist$Nx2Lvl ) {
269   a2Sample[x2idx,] = samplesSample( paste("a2[",x2idx,"]",",sep="") )
270 }
271 a1a2Sample = array(0, dim=c( datalist$Nx1Lvl , datalist$N      x2Lvl , chainLength ) )
272 for ( x1idx in 1:datalist$Nx1Lvl ) {
273   for ( x2idx in 1:datalist$Nx2Lvl ) {
274     a1a2Sample[x1idx,x2idx,] = samplesSample( paste( "a1a2[      ",x1idx,",",x2idx,"]",",sep="" ) )
275   }
276 }
277 }
278
279 # Convert to zero-centered b values:
280 m12Sample = array( 0 , dim=c( datalist$Nx1Lvl , datalist$Nx      2Lvl , chainLength ) )
281 for ( stepIdx in 1:chainLength ) {
282   m12Sample[,stepIdx] = ( a0Sample[stepIdx]
283                         + outer( a1Sample[stepIdx] ,
284                                 a2Sample[stepIdx] , "+" )
285                         + a1a2Sample[,stepIdx] )
286 }
287 b0Sample = apply( m12Sample , 3 , mean )
288 b1Sample = ( apply( m12Sample , c(1,3) , mean )
289               - matrix(rep( b0Sample ,Nx1Lvl),nrow=Nx1Lvl,byrow=T) )
290 b2Sample = ( apply( m12Sample , c(2,3) , mean )
291               - matrix(rep( b0Sample ,Nx2Lvl),nrow=Nx2Lvl,byrow=T) )
292 linpredSample = array(0,dim=c(datalist$Nx1Lvl,datalis      t$Nx2Lvl,chainLength))
293 for ( stepIdx in 1:chainLength ) {
294   linpredSample[,stepIdx] = ( b0Sample[stepIdx]
295                             + outer( b1Sample[stepIdx] ,
296                                     b2Sample[stepIdx] , "+" ) )
297 }
298 b1b2Sample = m12Sample - linpredSample
299 # Convert from standardized b values to original scale b valu      es:
300 b0Sample = b0Sample * ySDorig + yMorig
301 b1Sample = b1Sample * ySDorig
302 b2Sample = b2Sample * ySDorig
303 b1b2Sample = b1b2Sample * ySDorig
304
305 # Plot b values:
306 windows((datalist$Nx1Lvl+1)*2.75,(datalist$Nx2Lvl+1)      )*2.0)
307 layoutMat = matrix( 0 , nrow=(datalist$Nx2Lvl+1) , ncol=(d      atalist$Nx1Lvl+1) )
308 layoutMat[1,1] = 1
309 layoutMat[1,2:(datalist$Nx1Lvl+1)] = 1:datalist$Nx1Lv      l + 1

```

```

310 layoutMat[2:(datalist$Nx2Lvl+1),1] = 1:datalist$Nx2Lv      | + (datalist$Nx1Lvl + 1)
311 layoutMat[2:(datalist$Nx2Lvl+1),2:(datalist$Nx1Lvl+      1)] = matrix(
312   1:(datalist$Nx1Lvl*datalist$Nx2Lvl) + (datalist$Nx2Lv      |+datalist$Nx1Lvl+1) ,
313   ncol=datalist$Nx1Lvl , byrow=T )
314 layout( layoutMat )
315 par( mar=c(4,0.5,2.5,0.5) , mgp=c(2,0.7,0) )
316 histinfo = plotPost( b0Sample , xlab=expression(beta * 0) ,           main="Baseline" ,
317                      breaks=30 )
318 for ( x1idx in 1:datalist$Nx1Lvl ) {
319   histinfo = plotPost( b1Sample[x1idx.] , xlab=bquote(beta      *1[(x1idx)]) ,
320                      main=paste("x1:",x1names[x1idx]) , breaks=30 )
321 }
322 for ( x2idx in 1:datalist$Nx2Lvl ) {
323   histinfo = plotPost( b2Sample[x2idx.] , xlab=bquote(beta      *2[(x2idx)]) ,
324                      main=paste("x2:",x2names[x2idx]) , breaks=30 )
325 }
326 for ( x2idx in 1:datalist$Nx2Lvl ) {
327   for ( x1idx in 1:datalist$Nx1Lvl ) {
328     histinfo = plotPost( b1b2Sample[x1idx,x2idx.] , breaks=3      0 ,
329                           xlab=bquote(beta*12[(x1idx)^*,(x2idx)]) ,
330                           main=paste("x1:",x1names[x1idx]," x2:",x2names[x2id      x]) )
331   }
332 }
333 dev.copy2eps(file=paste(fnroot,"b.eps",sep=""))
334
335 # Display contrast analyses
336 nContrasts = length( x1contrastList )
337 if ( nContrasts > 0 ) {
338   nPlotPerRow = 5
339   nPlotRow = ceiling(nContrasts/nPlotPerRow)
340   nPlotCol = ceiling(nContrasts/nPlotRow)
341   windows(3.75*nPlotCol,2.5*nPlotRow)
342   layout( matrix(1:(nPlotRow*nPlotCol),nrow=nPlotRow,n      col=nPlotCol,byrow=T) )
343   par( mar=c(4,0.5,2.5,0.5) , mgp=c(2,0.7,0) )
344   for ( cldx in 1:nContrasts ) {
345     contrast = matrix( x1contrastList[[cldx]],nrow=1 ) # make      it a row matrix
346     incldx = contrast!=0
347     histInfo = plotPost( contrast %*% b1Sample , compVal=0 , bre     aks=30 ,
348                           xlab=paste( round(contrast[incldx],2) , x1names[incldx      ] ,
349                           c(rep("+",sum(incldx)-1),"") , collapse=" " ) ,
350                           cex.lab = 1.0 ,
351                           main=paste( "X1 Contrast:", names(x1contrastList)[cldx      ] ) )
352   }
353   dev.copy2eps(file=paste(fnroot,"x1Contrasts.eps",se      p=""))
354 }
355 #
356 nContrasts = length( x2contrastList )
357 if ( nContrasts > 0 ) {
358   nPlotPerRow = 5
359   nPlotRow = ceiling(nContrasts/nPlotPerRow)
360   nPlotCol = ceiling(nContrasts/nPlotRow)
361   windows(3.75*nPlotCol,2.5*nPlotRow)
362   layout( matrix(1:(nPlotRow*nPlotCol),nrow=nPlotRow,n      col=nPlotCol,byrow=T) )
363   par( mar=c(4,0.5,2.5,0.5) , mgp=c(2,0.7,0) )
364   for ( cldx in 1:nContrasts ) {
365     contrast = matrix( x2contrastList[[cldx]],nrow=1 ) # make      it a row matrix
366     incldx = contrast!=0
367     histInfo = plotPost( contrast %*% b2Sample , compVal=0 , bre     aks=30 ,
368                           xlab=paste( round(contrast[incldx],2) , x2names[incldx      ] ) ,

```

```

369           c(rep("+",sum(inclidx)-1),"") , collapse=" " ) ,
370           cex.lab = 1.0 ,
371           main=paste( "X2 Contrast:", names(x2contrastList)[cldx] ] ) )
372     }
373   dev.copy2eps(file=paste(fnroot,"x2Contrasts.eps",sep=""))
374 }
375 #
376 nContrasts = length( x1x2contrastList )
377 if ( nContrasts > 0 ) {
378   nPlotPerRow = 5
379   nPlotRow = ceiling(nContrasts/nPlotPerRow)
380   nPlotCol = ceiling(nContrasts/nPlotRow)
381   windows(3.75*nPlotCol,2.5*nPlotRow)
382   layout( matrix(1:(nPlotRow*nPlotCol),nrow=nPlotRow,n      col=nPlotCol,byrow=T) )
383   par( mar=c(4,0.5,2.5,0.5) , mgp=c(2,0.7,0) )
384   for ( cldx in 1:nContrasts ) {
385     contrast = x1x2contrastList[[cldx]]
386     contrastArr = array( rep(contrast,chainLength) ,
387                           dim=c(NROW(contrast),NCOL(contrast),chainLength) )
388     contrastLab = ""
389     for ( x1idx in 1:Nx1Lvl ) {
390       for ( x2idx in 1:Nx2Lvl ) {
391         if ( contrast[x1idx,x2idx] != 0 ) {
392           contrastLab = paste( contrastLab , "+" ,
393                               signif(contrast[x1idx,x2idx],2) ,
394                               x1names[x1idx] , x2names[x2idx] )
395         }
396       }
397     }
398     histInfo = plotPost( apply( contrastArr * b1b2Sample , 3 , su      m ) ,
399                          compVal=0 , breaks=30 , xlab=contrastLab , cex.lab = 0.75 ,
400                          main=paste( names(x1x2contrastList)[cldx] ) )
401   }
402   dev.copy2eps(file=paste(fnroot,"x1x2Contrasts.eps",sep=""))
403 }
404 #
405 #=====
406 # Do NHST ANOVA:
407
408 theData = data.frame( y=y , x1=factor(x1,labels=x1names) ,
409                       x2=factor(x2,labels=x2names) )
410 windows()
411 interaction.plot( theData$x1 , theData$x2 , theData$y , ty      pe="b" )
412 dev.copy2eps(file=paste(fnroot,"DataPlot.eps",sep=""))
413 aovresult = aov( y ~ x1 * x2 , data = theData )
414 cat("\n-----\n")
415 print( summary( aovresult ) )
416 cat("\n-----\n")
417 print( model.tables( aovresult , type = "effects" , se = TRUE ) , digits=3 )
418 cat("\n-----\n")
419 #
420 #=====
```

19.4 Exercises

Exercise 19.1. [Purpose: Inspecting an interaction for transformed data] Consider the data plotted in Figure 19.3, p. 426.

(A) Is the interaction a crossover interaction or not? Briefly explain your answer.

(B) Suppose we are interested in salaries thought of in terms of percentage (i.e., ratio) differences rather than additive differences. Therefore we take the logarithm, base 10, of the individual salaries (the R code has this option built into the data section, where the salary data are loaded). Run the analysis on the transformed data producing the results and contrasts analogous to those in Figures 19.4 and 19.5. Do the conclusions change?

(C) Examine the within-cell variances in the original and in the transformed data. (Hint: Try using the aggregate function on the data. As a guide, see how the function is used to initialize a1a2. Instead of applying the mean to the aggregated data, apply the standard deviation. The result is the within-cell standard deviation. Are they all roughly the same?) Do the original or the transformed data better respect the assumption of homogeneous variances?

Exercise 19.2. [Purpose: Investigating a case of two factor Bayesian ANOVA] In the data specification of the program in Section 19.3 ANOVAatowayBRugs, you can load data from Qian and Shen (2007), regarding how quickly seaweed regenerates in the presence of different types of grazers. Data were collected from eight different tidal areas of the Oregon coast; this predictor (x_2) is referred to as the “Block”. In each of the eight Blocks, there were six different combinations of seaweed grazers established by researchers. This predictor (x_1) had six levels: Control, with no grazers allowed; only small sh allowed; small and large sh allowed; only limpets allowed; limpets and small sh allowed; and, all three grazers allowed. The predicted variable was the percentage of the plot covered by regenerated seaweed, logarithmically transformed.

(A) Load the data and run the program. You will find that there are many levels of the two predictors to fit all the posterior histograms into a single multi-panel display. Therefore, modify the plotPost.R program so that it produces the mean and HDI limits, marked by a horizontal bar with a circle at the mean (without histogram), and perhaps without a main title. Name your program something other than plotPost.R, use it in the plotting section at the end of the program instead of plotPost.R, show your results. (A secondary goal of this part of the exercise is to give you practice modifying graphics in R to suit your own purposes.) Hints: There are many ways to do this, but here are some options. To suppress plotting of the histogram, just type argument in the hist function: plot=F . To suppress a title on a plot, just use the argument main="" . To adjust the font size, specify the “character expansion” for text, cex.lab for axis labels, and so forth. To reduce the margins around a plot, so there is more room for the plot itself, try variations of these margin specifications: par(mar=c(2,0.5,1,0.5), mgp=c(0.5,0,0)) . The par command needs to be called before the plots are made.

(B) The program already includes contrasts that consider whether there is an effect of small sh, and effect of large sh, and an effect of limpets. What conclusions do you reach from the posteriors of these contrasts?

(C) Construct a contrast of the average of Blocks 3 and 4 versus the average of Blocks 1 and 2. Show your specification, the graph of the posterior distribution, and state your conclusion.

(D) Is the effect of limpets different in Block 6 than in Block 7? To answer this question,

construct an interaction contrast using an outer product: (refer to the already-coded L_{effect} for the contrast that specifies the effect of limpets). Is the effect of small shell diameter in Blocks 1 and 7 than in Blocks 3 and 5? For both questions, show the contrasts vectors that you constructed and show the posterior of the contrast; state your conclusion.

Exercise 19.3 [Purpose: Notice that within-subject designs can be more sensitive (hence more powerful) than between-subject designs.] Consider these data:

	x ₁	x ₂	y	S
1	1	101	1	
1	1	102	2	
1	1	103	3	
1	1	105	4	
1	1	104	5	
1	2	104	1	
1	2	105	2	
1	2	107	3	
1	2	106	4	
1	2	108	5	
2	1	105	1	
2	1	107	2	
2	1	106	3	
2	1	108	4	
2	1	109	5	
2	2	109	1	
2	2	108	2	
2	2	110	3	
2	2	111	4	
2	2	112	5	

(A) Ignoring the last column, which indicates the subject who generated the data, conduct a Bayesian ANOVA using x_1 and x_2 as predictors of y . Show the code you used to load the data, and show the resulting posterior histograms, $\pi_{0;j}$, $\pi_{2;k}$, and $\pi_{1;2;jk}$. Also show the posterior of the contrast $\pi_{2-1;1}$ (i.e., the marginal difference between levels 1 and 2 of factor 1, also called the main effect of factor 1), and the posterior of the contrast $\pi_{2;2-2;1}$ (i.e., the marginal difference between levels 1 and 2 of factor 2, also called the main effect of factor 2).

(B) Now include the subject as a predictor, by expanding the model to include a deviation from baseline due to subject. (Do not include any subject interaction terms.) Again show the posteriors of the's requested in the previous part. Are the certainties on the estimates and contrasts different than in the previous part? In what way, and why?

Hint regarding the answer: Figure 19.7 shows posterior histograms for the main effect of factor 2, when the data are considered to be between-subjects or within-subject. Notice that the means are (essentially) the same in both histograms, but the uncertainties are very different!

Programming hints: The model specification without a subject factor is

```
mu[i] <- a0 + a1[x1[i]] + a2[x2[i]] + a1a2[x1[i],x2[i]]
```

but with a subject factor becomes

```
mu[i] <- a0 + a1[x1[i]] + a2[x2[i]] + a1a2[x1[i],x2[i]] + aS[S[i]]
```

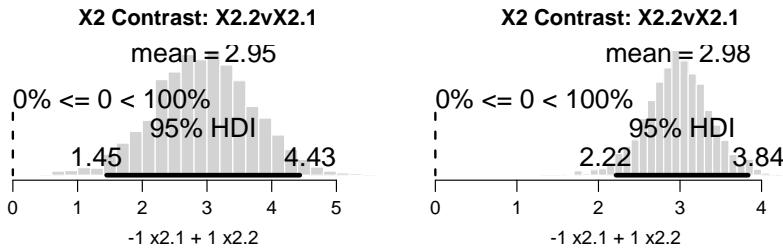


Figure 19.7: For Exercise 19.3. Left: Posterior for difference between levels of factor 2 when data are considered to be between-subject. Right: Posterior for difference between levels of factor 2 when data are considered within-subject. Notice that the means are (essentially) the same in both plots, but the uncertainties are very different!

where $s[i]$ is the subject number for the i^{th} datum, and $s[]$ are the deflections from baseline for each subject. You must, of course, specify a prior $a[1]$ analogous to the prior on $a1[]$.

The conversion of the $s[]$ values to zero-centered $b[]$ values proceeds analogously to what was explained at the beginning of Section 19.3. The code merely needs to be expanded to include the additional subject. Here is a guide: (ANOVAtwowayBRugsWithinSubj.R)

```

209 # Convert the a values to zero-centered b values.
210 # m12Sample is predicted cell means at every step in MCMC chain:
211 m12Sample = array( 0, dim=c( datalist$Nx1Lvl , datalist$Nx2Lvl ,
212                           datalist$NSLvl , chainLength ) )
213 for ( stepIdx in 1:chainLength ) {
214   for ( a1idx in 1:Nx1Lvl ) {
215     for ( a2idx in 1:Nx2Lvl ) {
216       for ( aSidx in 1:NSLvl ) {
217         m12Sample[ a1idx , a2idx , aSidx , stepIdx ] = (
218           a0Sample[stepIdx]
219           + a1Sample[a1idx,stepIdx]
220           + a2Sample[a2idx,stepIdx]
221           + a1a2Sample[a1idx,a2idx,stepIdx]
222           + aSSample[aSidx,stepIdx] )
223       }
224     }
225   }
226 }
227
228 # b0Sample is mean of the cell means at every step in chain:
229 b0Sample = apply( m12Sample , 4 , mean )
230 # b1Sample is deflections of factor 1 marginal means from b0Sample:
231 b1Sample = ( apply( m12Sample , c(1,4) , mean )
232             - matrix(rep( b0Sample ,Nx1Lvl),nrow=Nx1Lvl,byrow=T) )
233 # b2Sample is deflections of factor 2 marginal means from b0Sample:
234 b2Sample = ( apply( m12Sample , c(2,4) , mean )
235             - matrix(rep( b0Sample ,Nx2Lvl),nrow=Nx2Lvl,byrow=T) )
236 # bSSample is deflections of factor S marginal means from b0Sample:
237 bSSample = ( apply( m12Sample , c(3,4) , mean )
238             - matrix(rep( b0Sample ,NSLvl),nrow=NSLvl,byrow=T) )
239 # linpredSample is linear combination of the marginal effects:
240 linpredSample = 0*m12Sample

```

```

241 for ( stepIdx in 1:chainLength ) {
242   for ( a1idx in 1:Nx1Lvl ) {
243     for ( a2idx in 1:Nx2Lvl ) {
244       for ( aSidx in 1:NSLvl ) {
245         linpredSample[a1idx,a2idx,aSidx,stepIdx] = (
246           b0Sample[stepIdx]
247           + b1Sample[a1idx,stepIdx]
248           + b2Sample[a2idx,stepIdx]
249           + bSSample[aSidx,stepIdx] )
250       }
251     }
252   }
253 }
254 # b1b2Sample is the interaction deflection, i.e., the difference
255 # between the cell means and the linear combination:
256 b1b2Sample = apply( m12Sample - linpredSample , c(1,2,4) , mean )
257
258 # Convert from standardized b values to original scale b values:
259 b0Sample = b0Sample * ySDorig + yMorig
260 b1Sample = b1Sample * ySDorig
261 b2Sample = b2Sample * ySDorig
262 bSSample = bSSample * ySDorig
263 b1b2Sample = b1b2Sample * ySDorig

```

Exercise 19.4. [Purpose: Power analysis for Bayesian ANOVA, for within-subject versus between-subject designs.] Conduct power analyses for the between-subject and ~~within~~ subject versions of the previous exercise. Specifically, suppose that ~~it~~ for the 95% HDI of the contrast on factor 2 to have a width of 2.0 or less. Conduct ~~respective~~ power analysis for this goal, for the within-subject version and the ~~between~~ subject version. Caution: This exercise demands a lot of programming and could be time ~~consuming~~, but the results drive home the point that within-subject designs can be more ~~powerful~~ than between-subject designs.

Chapter 20

Dichotomous Predicted Variable

Contents

20.1	Logistic regression50
20.1.1	The model	451
20.1.2	Doing it in R and BUGS	451
20.1.3	Interpreting the posterior	452
20.1.4	Perils of correlated predictors	454
20.1.5	When there are few 1's in the data	445
20.1.6	Hyperprior across regression coe nts nts	454
20.2	Interaction of predictors in logistic regression	455
20.3	Logistic ANOVA	456
20.3.1	Within-subject designs	845
20.4	Summary	458
20.5	R code	459
20.5.1	Logistic regression code	945
20.5.2	Logistic ANOVA code	463
20.6	Exercises	468

Fortune and Favor make ~~ckle~~ decrees, it's
Heads or it's tails with no middle degrees.
Flippant commandments decreed by law gods, have
Reasons so rare they have minus log odds.

There are many situations in which the value to be predicted is dichotomous (instead of metric). For example, we might want to predict whether ~~is~~ a patient is cured or not (a dichotomous outcome), on the basis of the dosage of drug ~~is~~ administered and age of the patient (two metric predictors). What the model does, ~~is~~ ~~to~~ ~~use~~ ~~the~~ ~~function~~ ~~to~~ ~~generate~~ ~~a~~ ~~prediction~~ ~~of~~ ~~the~~ ~~probability~~ ~~that~~ ~~a~~ ~~patient~~ ~~is~~ ~~cured~~, given the specified dosage and ~~other~~ ~~situations~~, ~~the~~ ~~predictors~~ ~~might~~ ~~be~~ ~~nominal~~. For example, we could ~~predict~~ ~~versus~~ ~~not-cured~~ on the basis of type of drug and gender of patient. The model ~~specifies~~ ~~the~~ probability that a patient is cured, given the drug and the gender. We will ~~introduce~~ ~~such~~ situations in this chapter.

The formal framework for this situation was presented in the second row of Table 14.1 on p. 312. As you may recall, the link function, which maps ~~any~~ ~~a~~ ~~linear~~ combination of predictors to an outcome tendency, is a logistic function. ~~Therefore~~ ~~the~~ models we use in this chapter are referred to as cases of logistic regression ~~or~~ ANOVA.

Kruschke, J. K. (2010) Doing Bayesian Data Analysis: A Tutorial with R and BUGS. Academic Press Elsevier. Copyright 2010 by John K. Kruschke. Draft of May 11, 2010. Please do not distribute this preliminary draft. If you report Bayesian analyses based on this book, please do cite it!_ "

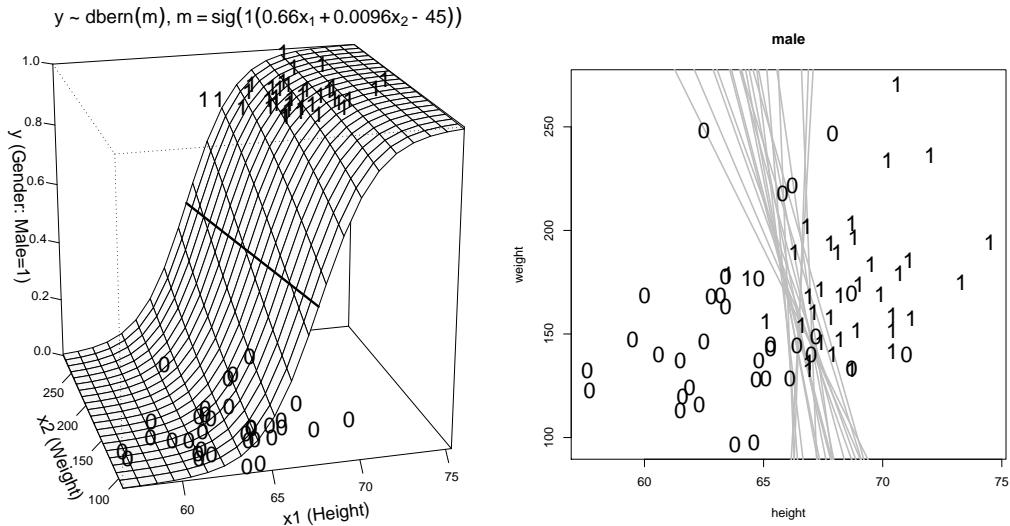


Figure 20.1: Gender (`male`, `female`=0), height, and weight. Left panel shows data in perspective, with the MLE logistic surface, and $p(y=1) = .5$ level contour marked as a dark line. The data points are all at $y=0$, not on the logistic surface. Right panel shows the data “from above” with several grey lines that indicate a smattering of credibility $p(y=1) = .5$ level contours derived from Bayesian logistic regression.

20.1 Logistic regression

Suppose we have metric predictors and a dichotomous ~~predictable~~ variable. As a concrete example, suppose we want to predict the gender of a person (`male`=1 and `female`=0) on the basis of the person's height and weight. Figure 20.1 shows realistic data from 70 people. Males are plotted as numeral 1's, and females as numeral 0's. You can see that there is notable variation of heights and weights ~~in~~ people, and, importantly, there is considerable overlap in the distributions of ~~males~~ females. Thus, at any specific combination of height and weight, the best we can do is predict the probability that a person with that height and weight is male.

The model we will use to formulate the probability that $y=1$, as a function of the two predictors, is the logistic transform of a linear combination of the predictors. Formal details will be reviewed below, but a graphical preview is provided in the left panel of Figure 20.1. The smoothly ascending surface plots the probability that $y=1$ as a function of x_1 and x_2 . The parameters of the model control the “shape”, specifically its orientation, location, and steepness.

The maximum likelihood estimate (MLE) of the model parameters provides a single logistic surface that summarizes the distribution of the data. The MLE surface is shown in the left panel of Figure 20.1. Also plotted on the surface, a dark line, is the level contour at which $p(y=1) = 0.5$. This level contour marks the points that are halfway up the logistic curve. The level contour is a convenient summary of the orientation and location of the logistic surface, but the level contour loses information about the steepness of the curve.

The right panel of Figure 20.1 shows the same data viewed from above. Also plotted are several credible $p(y=1) = 0.5$ level contours, from a Bayesian analysis. Unlike the single MLE contour in the left panel, the distribution of credible contours reveals our uncertainty

in the parameter estimates of the logistic-regression model. The mechanics of the Bayesian analysis are described in the next sections.

20.1.1 The model

Recall that for multiple linear regression, the central tendency of the predicted value is a linear combination of the predictors, as in Equation 14.207), and as explained at length in Chapter 17. In multiple logistic regression, the linear combination is transformed by a logistic squashing function, and the resulting value, between zero and one, is used as the probability that the predicted value is one. Formally, at the second row of Table 14.1 on p. 312, logistic regression can be written as

$$\begin{aligned} &= \text{sig } _0 + _1 x_1 + _2 x_2 + \dots \\ y &\sim \text{dbern}() \end{aligned}$$

where sig refers to the sigmoid function, which is merely the name for the logistic function:

$$\text{sig}(x) = \frac{1}{1 + \exp(-x)} \quad (20.1)$$

Graphs of the logistic (a.k.a. sigmoid) function were shown in Figure 14.6 (p. 305) for a single predictor, Figure 14.7 (p. 306) for two predictors, and Figure 14.10 (p. 310) for two predictors with dichotomous predicted values superimposed. Please do review those figures now!

The logit function is the inverse of the logistic function. or finally, for $0 < p < 1$, $\text{logit}(p) = \ln(p/(1-p))$. In applications to logistic regression, the "is the probability that $y = 1$, and therefore we can write $\text{logit}(y=1) = \ln(p(y=1)/p(y=0))$, where the logarithm is the natural logarithm, i.e., the inverse of the exponential function. The ratio, $p(y=1)/p(y=0)$, is called the odds of outcome 1 to outcome 0. The logistic regression model can be written in terms of the logit function, also called the log-odds form, like this:

$$\begin{aligned} \text{logit}() &= _0 + _1 x_1 + _2 x_2 + \dots \\ y &\sim \text{dbern}() \end{aligned} \quad (20.2)$$

I prefer the expression in terms of the logistic (i.e., sigmoid) function in Equation 20.1 because it provides an explicit specification of that is natural to depict in hierarchical diagrams. But the logit form in Equation 20.2 is useful for interpreting the regression coefficients, as will be described in Section 20.1.3.

The hierarchical diagram for the model is shown in Figure 220. It is just like the one for multiple linear regression in Figure 17.4, p. 375, except it includes a sigmoid and a change in the likelihood function at the bottom of the diagram from a normal distribution to a Bernoulli distribution.

To reiterate, what the model does is take the individual's predictor values x_{ji} , and generate the probability $y = 1$ for that individual. The Bayesian estimation yields values of the parameters θ_j and β_j , that respect the data and the prior.

20.1.2 Doing it in R and BUGS

The implementation of logistic regression in R, BRugs, and JAGS is straight forward, with only a few modifications of the linear regression program presented earlier in the book. The full logistic regression program is listed in Section 20.5.1

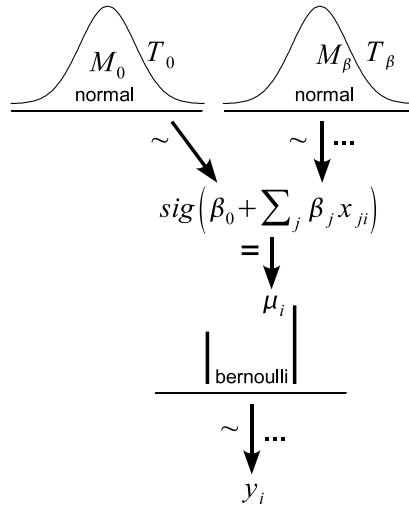


Figure 20.2: Hierarchical model for logistic regression analogous to the hierarchical model for linear regression in Figure 17.4, p. 375. A prior across regression coefficients can be added, as in Figure 17.6, p. 379.

(`MultipleLogisticRegressionBruks.R`). The only notable changes, other than the model, are (i) the use of R's `glm` function to initialize the chains at the MLE, and (ii) a `skip` formula for conversion of standardized data to original-scale. A more thorough explanation is provided before the listing in Section 20.5 (`MultipleLogisticRegressionBruks.R`).

20.1.3 Interpreting the posterior

Consider again the height, weight, and gender data in Figure 20.1. Aspects of the posterior distribution are shown in Figure 20.3. Often we are primarily interested in knowing how big an influence a predictor has on the predicted value, and, specifically, whether that influence is credibly non-zero.

The middle panel of Figure 20.3 indicates that height provides credibly non-zero predictiveness for gender. The mean value of the effect on height is 0.721, which means that the log odds, $\log(p(\text{male})/p(\text{female}))$, increases by 0.721 when height increases by one inch.

This notion of constant increase in log odds is not instantly intuitive. Recall from Equation 20.2 that the regression equation can be expressed in odds form as

$$\text{logit}(p) = \log \frac{p(y=1)}{p(y=0)} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots$$

Therefore, when x_1 increases one unit, the log-odds increase β_1 units. Some numerical examples might help. Consider a hypothetical person who weighs 160 pounds. We will compute what happens to the log-odds when the person's height increases by 1 inch. First, consider an increase from 63 inches tall to 64 inches tall. According to the logistic function with parameters set at the mean of the posterior, if the 160-pound person were 63 inches tall, then the probability of being male is 0.0728, and if the person were 64 inches tall, the probability of being male is 0.1390. That increase of 6.62 percentage points corresponds to a change in log odds of $\log(0.1390/(1 - 0.1390)) - \log(0.0728/(1 - 0.0728))$, which equals 0.721. Next, consider an increase from 66 inches tall to 67 inches tall. If the person

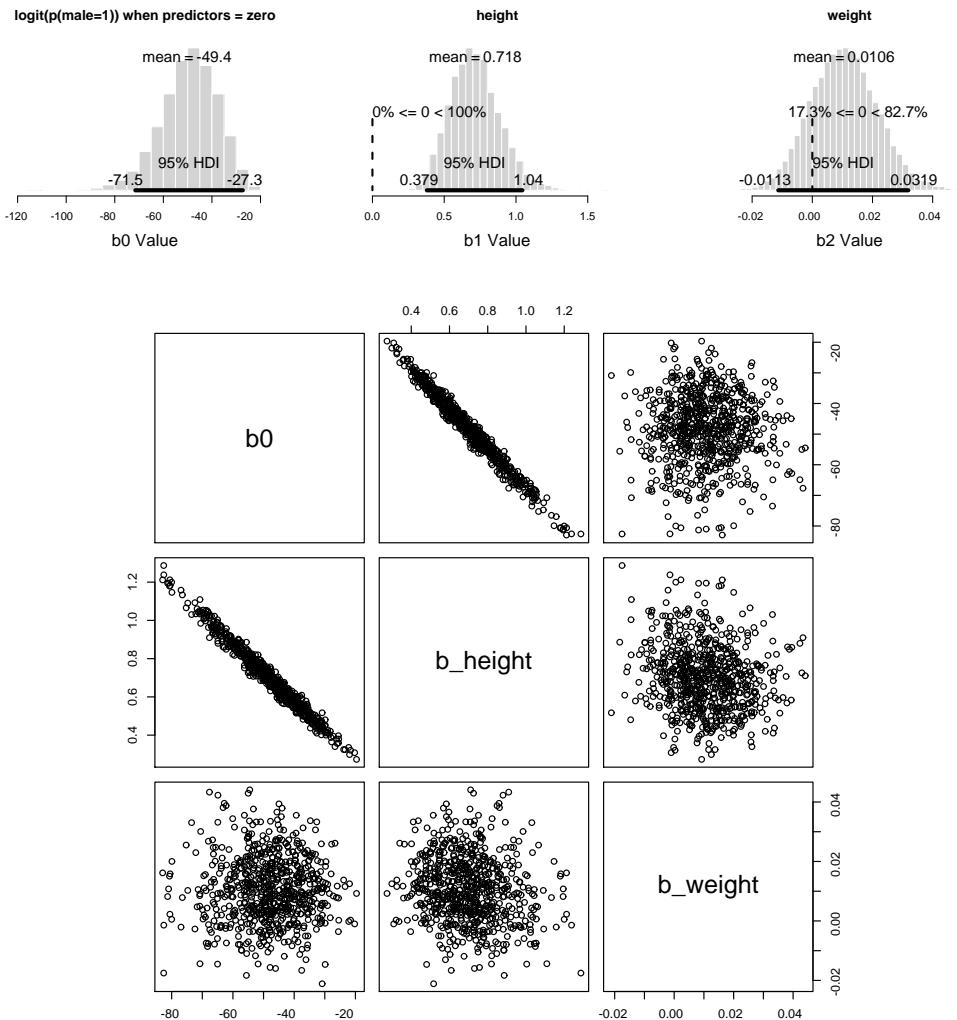


Figure 20.3: Posterior distribution of Bayesian logistic regression of gender on height and weight.

were 66 inches tall, the probability of being male is 0.4056, if the person were one inch taller, at 67 inches, then the probability of being male is 0.5839. The increase of 17.83 percentage points corresponds to a change in log odds $\log(0.5839/(1 - 0.5839)) - \log(0.4056/(1 - 0.4056))$, which again equals 0.721. Thus, for any increase of one inch in height, the increase in predicted log-odds of being male is the same, namely 0.721.

The right panel of Figure 20.3 indicates that weight does not provide a credibly non-zero predictiveness for gender, because zero is well among the values in the posterior. The mean value of the coefficient on weight is 0.0106, which means that the log odds, $\log(p(\text{male})/p(\text{female}))$, increases by only 0.0106 when weight increases by one pound.

The lower part of Figure 20.3 shows pairwise scatterplots of the 3-dimensional posterior. Notice that the intercept and the coefficient on height are strongly negatively correlated. This correlation results from two properties of the dataset: height is predictive of gender, and therefore the intercept is constrained by height. Since these original-scale data are not standardized at heights of zero, so if the coefficient on height is increased, the intercept

must be decreased to keep the logistic at the appropriate position.

20.1.4 Perils of correlated predictors

When predictors are correlated, care must be taken in fitting their regression coefficients. This issue was emphasized in Section 17.1.1 in the context of multiple linear regression, but analogous perils arise in logistic regression.

Consider a situation with two predictors. When the predictors are uncorrelated, the posterior will have relatively high certainty regarding the regression coefficients. But when the predictors are strongly correlated, the posterior will have relatively low certainty about the regression coefficients, because many different logistic surfaces can fit the data fairly well. Figure 20.7 in Exercise 20.1 provides an example.

Another peril arising from correlated predictors arises when one of the correlated predictors is not included in the model. Suppose that one predictor is included in the model, but another predictor is not included, perhaps because the excluded predictor were lost or unobtainable or not even considered in the first place. Suppose further that the true regression coefficient on the included predictor is zero but the true regression coefficient on the excluded predictor is large positive. The posterior estimate of the included regression coefficient will be non-zero if the included predictor is correlated with the excluded predictor. The reason is intuitively straight forward: As the included predictor increases, the excluded predictor changes (because it is correlated with the included predictor), and therefore the outcome changes, even though the included predictor has no direct predictive value for the outcome. Exercise 20.1 provides an example of this situation.

20.1.5 When there are few 1's in the data

In linear regression, the predicted variable is on a metric scale. Usually the data are distributed over a range of values, without being severely clustered over a single value. In logistic regression, however, the predicted variable is dichotomous, and the data can sometimes consist of mostly 1's or mostly 0's. For example, if the predicted variable is the occurrence of a rare disease, then, by definition, there are few 1's. As another example, the predicted variable might be the occurrence of a defect on an assembly line, which is expected to be rare.

The problem with data that have only a few 1's or a few 0's is that the estimate of the regression coefficients is usually relatively uncertain. This makes inference less precise, because it is only the transition from 0's to 1's that constrains the regression coefficients. Exercise 20.2 provides an example.

20.1.6 Hyperprior across regression coefficients

When there are many candidate predictors, we can use prior knowledge to mutually constrain the estimates of the coefficients. Because all the predictors come from a pool of factors that might have some remotely plausible prediction with the dichotomous value, we could reasonably put a hyperprior on the regression coefficients that expresses the assumption that most coefficients are near zero, but some may be notably farther from zero. A distribution that captures this prior knowledge is a distribution, with a precision that is estimated from the data. This scheme was described in the context of multiple linear regression in Section 17.2, p. 378. A diagram of the hierarchical prior was displayed

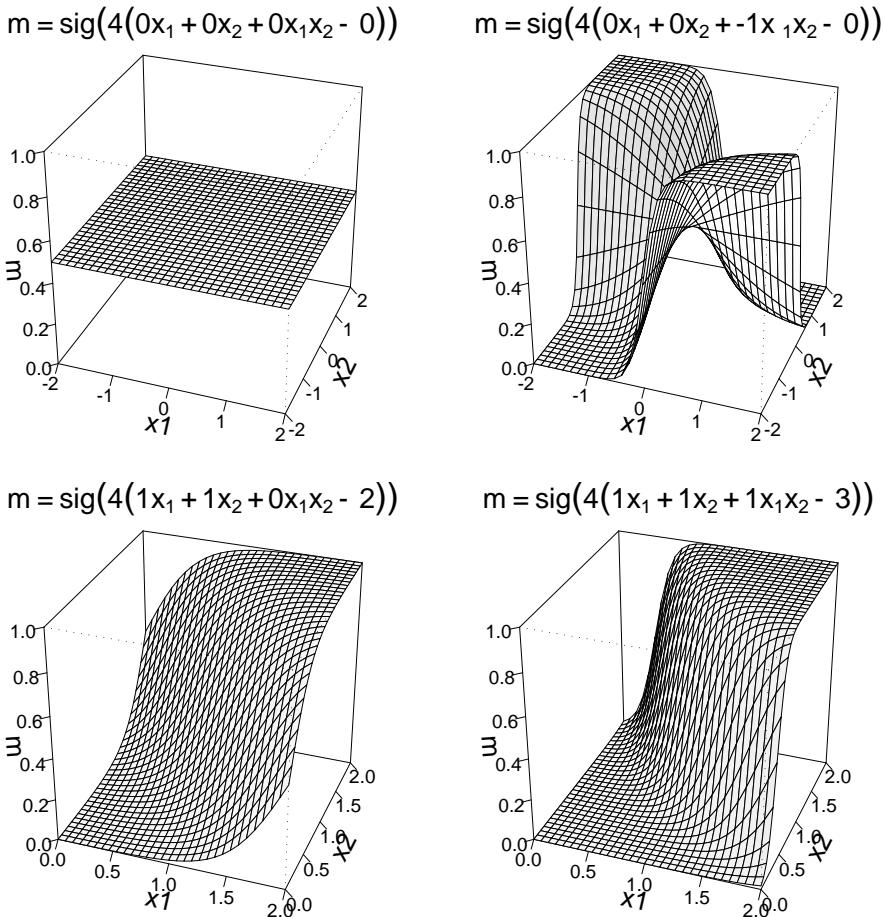


Figure 20.4: Right panels: Logistic regression with interaction. Left panels show corresponding function without any interaction component. Notice that the upper panels include negative values of the predictors, whereas lower panels do not. In the lower-right graph, when both x_1 and x_2 are sufficiently large negative values (not shown), then the surface rises up toward 1.

in Figure 17.6, p. 379. The same scheme can be applied to logistic regression, merely by substituting the Bernoulli likelihood function, as in Figure 20.2.

The knowledge expressed in the hyperprior provides shrinkage on the coefficients in logistic regression. When several predictors have estimate coefficients near zero, they drive plausible values of the precision of the overestimating toward higher values, thereby shrinking estimates of other predictor's coefficients. The shrinkage attenuates false alarms, as discussed at length in Section 17.2.

20.2 Interaction of predictors in logistic regression

Logistic regression can include a multiplicative interaction just as when the predicted variable is metric (recall Figures 14.3 and 17.8). The interpretation of such an interaction requires care.

Suppose that the predicted variable is a two-valued spirit of happiness: The respon-

dent says he is happy or not happy. The two predictors are annual income and healthiness assayed on a metric scale. Intuitively, we can imagine that happiness depends on the conjunction of good health and high income; either one alone may be insufficient for happiness. To put it another way, a person can be happy either because of poor health or because of poverty. Thus, the probability of being happy in a conjunctive interaction of health and wealth.

As another example, consider predicting the success of a high-air vehicle. One predictor is the intensity of heat under the balloon. A second predictor is the amount of hydrogen in the balloon. Intuitively, the balloon will rise if there is sufficient heat or sufficient hydrogen but not both. This is called an exclusive-OR interaction between the predictors.

Figure 20.4 shows graphs of multiplicative interactions of predictors in a logistic function. The top-right panel shows an exclusive-OR: The predicted probability, p , is high only if x_1 is large or x_2 is large but not both. The top-left panel shows absence of interaction as a comparison and reminder that the baseline for a logistic is 0.5.

The lower-right panel of Figure 20.4 shows a conjunctive interaction: The predicted probability is high only if x_1 and x_2 are large. The lower-left panel shows a comparison that has zero interaction. The two lower panels are subtler: Level contours (not shown) on the zero-interaction curve are straight lines, but level contours (not shown) on the interactive curve are curved, resembling hyperbolas.

Multiplicative interactions can be incorporated into logistic regression models in the same way they are incorporated into metric-predictor models. As was emphasized in the chapter on multiple linear regression, the coefficients on the predictors must be interpreted with care, especially when interactions are involved.

20.3 Logistic ANOVA

There are many situations in which the variable to be predicted is dichotomous, and the predictors are nominal. As an example, recall the Italian condensation experiment introduced in Section 9.3.1, p. 178, and used many times thereafter. The predicted variable was accuracy on each trial, i.e., correct or incorrect. The predictor was a nominal variable that indicated which of four category structures the learner experienced. As another example, recall the relevance-shift experiment of Exercise 13.4, p. 282. Again the predicted variable was accuracy on each trial, and the predictor was a nominal variable indicating which of four relevance shifts the learner experienced. If the predictor were metric instead of dichotomous, then these cases could be analyzed with a traditional ANOVA as in Chapter 18.

The logistic function can be used to adapt ANOVA models to dichotomous outcomes. We start with the ANOVA model to generate a predicted value for each condition, then pass that value through a sigmoid so it becomes the predicted probability that $y = 1$ for that condition. We also assume that subjects can vary around a condition's predicted probability, just as in the ANOVA model. Figure 20.5 shows the logistic ANOVA model for a single factor, i.e., "oneway" logistic ANOVA. The basic ANOVA model comes directly from the top part of Figure 18.1, p. 403. The beta distribution for subjects within conditions comes directly from the left side of Figure 9.17, p. 183. Note that this model assumes equal variances for all conditions, like traditional ANOVA, as far as the same value is used for all groups simultaneously.

Figure 20.6 shows the results of running logistic oneway ANOVA on the Itra-

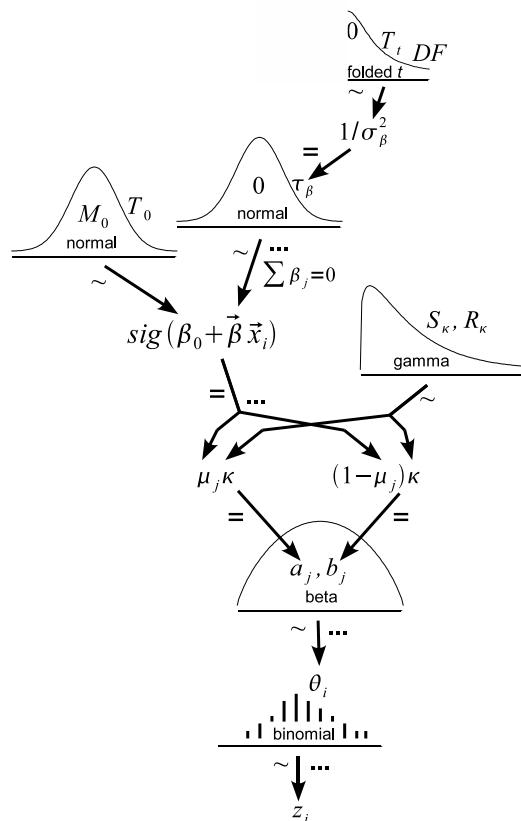


Figure 20.5: Bayesian logistic ANOVA. The ANOVA model, in the upper parts, generates the μ_j for the groups through the sigmoidal transform. The beta distribution, in the lower part of the diagram, describes the within-group variability across subjects. The same is assumed for all groups, which is analogous to homogeneity of variance in metric ANOVA.

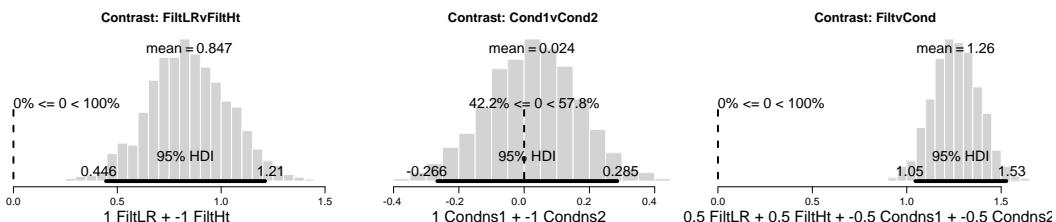


Figure 20.6: Results of logistic ANOVA applied to Itration/condensation experiment. Compare with Figure 9.16, p. 181, which did not use ANOVA model for these data. The scale on these histograms is the β values, not difference of β values as in Figure 9.16.

tion/condensation data. The comparisons are differences of the parameters in the ANOVA model. The parameters are not probabilities, but the parameters can be mapped to probabilities via the logistic function. The comparisons of the parameters from different conditions yield the same conclusions as the analysis from ~~earlier~~ chapters, as in Figure 9.16, p. 181, which did not use an ANOVA model for these data.

If the ANOVA model leads to the same conclusions as the no ~~ANOVA~~ approach, what are the advantages or disadvantages of the two models? A disadvantage of the logistic ANOVA model is that it tends to have larger autocorrelation ~~its~~ parameters, and therefore takes longer to run because more thinning is needed. Another advantage of the logistic ANOVA model is that the parameters are not instantly interpretable as probabilities, but must be converted through the logistic. But the ANOVA approach has great advantages when the situation to be modeled becomes more complex. In particular, when there are two or more predictors, the ANOVA model provides a natural representation for the main effects and interactions of the predictors. This generalizability of the logistic ANOVA model is its trump card.

Some analysts of dichotomous data will treat N_i as if it were a metric value, and apply standard ANOVA. Because standard ANOVA assumes that the data come from normal distributions, the analyst will transform the values to $\frac{z_i}{N_i}$ (e.g., Winer, Brown, & Michels, 1991, p. 356), to approximate normality. Unfortunately the approximation is not very good for small N_i or for values of $z_i = N_i$ that are very close to 0 or 1. Most importantly, the ratio z_i/N_i loses information about the magnitude of N_i . If all we know is that $z_i/N_i = 0.75$, we don't know whether $z_i = 3$ and $N_i = 4$, or $z_i = 3000$ and $N_i = 4000$. This matters because a proportion based on a large sample should influence our beliefs more than a proportion based on a small sample. Therefore, ~~approximating~~ proportions fails in unbalanced designs, i.e., when there are different sample sizes for different subjects. In the Bayesian approach used here, every individual observation is an equal influence, and the Bayesian analysis automatically handles unbalanced designs appropriately.

20.3.1 Within-subject designs

Fortunately, there is no news when it comes to repeated measures within-subject designs. All the concepts of Section 19.2, p. 432, apply the same way to logistic ANOVA. There were two main messages of that section. First, in terms of model structure, if every subject contributes many observations to every combinations of factors, then it makes sense to apply a logistic ANOVA model to every individual subject, with hyperdistributions on the ANOVA coefficients to allow mutually informed estimation across subjects. If every subject contributes only one observation to any combination of items, then it makes sense to include a subject factor in the model, with no interaction of that subject factor with the other factors. Second, in terms of whether to use a repeated measures design at all, there must be a good rationale for assuming that repeated measures of the same subject are reasonably independent of each other, so that any of the designs can be applied.

20.4 Summary

The extension of linear regression and ANOVA to logistic regression and logistic ANOVA is straight forward. The hierarchical models of Figures 22.1 and 20.5 illustrated how the core linear model of regression or ANOVA is passed through a logistic function, and the

resulting value is used as the mean for a Bernoulli (or binomial) likelihood function. All the basic structural concepts of linear regression and ANOVA applied to logistic regression and ANOVA.

Only a few new technical concepts introduced in this chapter were the notion of level contour on the logistic surface. The contour at level $\pi = .5$ is useful for marking the boundary between x values for which $p(y=1) < .5$ and x values for which $p(y=1) > .5$.

Another new technical concept was interpreting a coefficient in logistic regression in terms of the increase in log-odds of the outcome when the predictor increases by one unit. Although it was not previously mentioned, the same idea applies to logistic ANOVA: The j value indicates the increase in log-odds of the outcome relative to the baseline, 0 . This interpretation must be qualified, however, when the model includes interaction terms, whether it is regression or ANOVA.

The remaining technical issues arise only in the mechanics of computing the analyses. The next section, which includes the R code, explains these. There are two main issues, one regarding standardizing the predictors (for regression only) and transforming the parameters back to the original scales, and the other dealing with how to initialize the chains by starting with the maximum likelihood estimate (MLE).

20.5 R code

20.5.1 Logistic regression code

The program for logistic regression is much like the analogous program for linear regression, with only a few changes. One obvious change is that the likelihood function is Bernoulli, not normal. Beyond the model specification, it is important to understand the following details.

The predictor values are standardized, to reduce autocorrelation in MCMC sampling. Notice, for example, the highly correlated original-scale parameter values in Figure 20.3, which would be difficult to navigate by most MCMC sampling algorithms. Unlike linear regression, however, the predicted values are not standardized, because the predicted values must be zero or one. Because the predicted values are not standardised, the conversion of the standardized parameter values to original-scale parameter values is less complicated:

$$\begin{aligned} \text{logit}(z_{ji}) &= \beta_0 + \sum_j \beta_j z_{ji} \\ &= \beta_0 + \sum_j \frac{x_{ji} - E[x]_j}{\sqrt{\text{Var}[x]_j}} \\ &= \beta_0 + \sum_j \frac{\beta_j E[x]_j}{\sqrt{\text{Var}[x]_j}} + \sum_j \frac{\beta_j}{\sqrt{\text{Var}[x]_j}} x_{ji} \\ &\quad | \quad \{z_{ji}\} \quad | \quad \{\beta_j\} \end{aligned}$$

This conversion is executed at every step in the MCMC chain, and is performed at the end of the program.

Initialization of the chains is based on parameter values that are plausible given the data. Specifically, the values of the MLE are used, because R has a built-in function for determining the MLE of logistic regression, namely the `glm` function with the `family` argument specified as `binomial(logit)`.

The final lines of the program perform an MLE of the logistic regression, just for comparison with the Bayesian results. The MLE is typically ~~argue~~ most credible values in the Bayesian posterior. But the posterior provides more information about correlations of parameters and the range of uncertainty in each parameter.

(MultipleLogisticRegressionBruks.R)

```

1 graphics.off()
2 rm(list=ls(all=TRUE))
3 fname = "MultipleLogisticRegressionBruks"
4 library(BRugs)          # Kruschke, J. K. (2010). Doing Bayesian data analysis:
5                                # A Tutorial with R and BUGS. Academic Press / Elsevier.
6 #-----
7 # THE MODEL.
8 modelstring = "
9 # BUGS model specification begins here...
10 model {
11   for( i in 1 : nData ) {
12     y[i] ~ dbern( mu[i] )
13     mu[i] <- 1/(1+exp(- b0 + inprod( b[], x[i,] ))))
14   }
15   b0 ~ dnorm( 0 , 1.0E-12 )
16   for ( j in 1 : nPredictors ) {
17     b[j] ~ dnorm( 0 , 1.0E-12 )
18   }
19 }
20 # ... end BUGS model specification
21 " # close quote for modelstring
22 # Write model to a file:
23 writeLines(modelstring,con="model.txt")
24 modelCheck( "model.txt" )
25
26 #-----
27 # THE DATA.
28
29 dataSource = c( "HtWt" , "Cars" , "HeartAttack" )[3]
30
31 if ( dataSource == "HtWt" ) {
32   fname = paste( fname , dataSource , sep="" )
33   # Generate random but realistic data:
34   source( "HtWtDataGenerator.R" )
35   dataMat = HtWtDataGenerator( nSubj = 70 , rndsd=474 )
36   predictedName = "male"
37   predictorNames = c( "height" , "weight" )
38   nData = NROW( dataMat )
39   y = as.matrix( dataMat[,predictedName] )
40   x = as.matrix( dataMat[,predictorNames] )
41   nPredictors = NCOL( x )
42 }
43
44 if ( dataSource == "Cars" ) {
45   fname = paste( fname , dataSource , sep="" )
46   dataMat = read.table(file="Lock1993data.txt",header=T      ,sep=" ")
47   predictedName = "AirBag"
48   predictorNames = c( "MidPrice" , "RPM" , "Uturn" )
49   nData = NROW( dataMat )
50   y = as.matrix( as.numeric( dataMat[,predictedName] > 0 ) ) # 0,1,2 to 0,1
51   x = as.matrix( dataMat[,predictorNames] )
52   nPredictors = NCOL( x )

```

```

53 }
54
55 if ( dataSource == "HeartAttack" ) {
56   fname = paste( fname , dataSource , sep="" )
57   dataMat = read.table(file="BloodDataGeneratorOutput.t      xt",header=T,sep=" ")
58   predictedName = "HeartAttack"
59   predictorNames = c( "Systolic", "Diastolic", "Weight", "C      holesterol",
60                      "Height", "Age" )
61 #  predictorNames = c( "Systolic", "Diastolic" )
62 nData = NROW( dataMat )
63 y = as.matrix( dataMat[,predictedName] )
64 x = as.matrix( dataMat[,predictorNames] )
65 nPredictors = NCOL( x )
66 }
67
68 # Prepare data for BUGS:
69 # Re-center data at mean, to reduce autocorrelation in MCMC s   ampling.
70 # Standardize (divide by SD) to make initialization easier.
71 standardizeCols = function( dataMat ) {
72   zDataMat = dataMat
73   for ( colIdx in 1:NCOL( dataMat ) ) {
74     mCol = mean( dataMat[,colIdx] )
75     sdCol = sd( dataMat[,colIdx] )
76     zDataMat[,colIdx] = ( dataMat[,colIdx] - mCol ) / sdCol
77   }
78   return( zDataMat )
79 }
80 zx = standardizeCols( x )
81 zy = y # y is not standardized; must be 0,1
82
83 # Get the data into BUGS:
84 dataList = list(
85   x = zx ,
86   y = as.vector( zy ) , # BUGS does not treat 1-column mat as vecto   r
87   nPredictors = nPredictors ,
88   nData = nData
89 )
90 modelData( bugsData( dataList ) )
91
92
93 #-----#
94 # INTIALIZE THE CHAINS.
95
96 nchain = 3
97 modelCompile( numChains = nchain )
98
99 genInitList <- function() {
100   glmInfo = glm( dataList$y ~ dataList$x , family=binomial(l      ogit) ) # R func.
101   show( glmInfo ) ; flush.console() # display in case glm() has      troubles
102   b0Init = glmInfo$coef[1]
103   b1Init = glmInfo$coef[-1]
104   return( list(
105     b0 = b0Init ,
106     b = b1Init
107   ) )
108 }
109 for ( chainIdx in 1 : nchain ) {
110   modellnits( bugsInits( genInitList ) )
111 }
```

```

112
113 #-----
114 # RUN THE CHAINS
115
116 # burn in
117 BurnInSteps = 1000
118 modelUpdate( BurnInSteps )
119 # actual samples
120 samplesSet( c( "b0" , "b" ) )
121 stepsPerChain = ceiling(5000/nchain)
122 thinStep = 50 # check autocorrelation and increase as needed
123 modelUpdate( stepsPerChain , thin=thinStep )
124
125 #-----
126 # EXAMINE THE RESULTS
127
128 source("plotChains.R")
129 source("plotPost.R")
130
131 # Check chains for mixing
132 checkConvergence = T
133 if ( checkConvergence ) {
134   b0Sum = plotChains( "b0" , saveplots=F , filenameroot=fnam e )
135   bSum = plotChains( "b" , saveplots=F , filenameroot=fname )
136 }
137
138 # Extract chain values:
139 zb0Sample = matrix( samplesSample( "b0" ) )
140 chainLength = length(zb0Sample)
141 zbSample = NULL
142 for ( j in 1:nPredictors ) {
143   zbSample = cbind( zbSample , samplesSample( paste("b[,j,    "],sep="" ) ) )
144 }
145
146 # Convert to original scale:
147 x = dataMat[,predictorNames]
148 y = dataMat[,predictedName]
149 My = mean(y)
150 SDy = sd(y)
151 Mx = apply(x,2,mean)
152 SDx = apply(x,2,sd)
153 b0Sample = 0 * zb0Sample
154 bSample = 0 * zbSample
155 for ( stepIdx in 1:chainLength ) {
156   b0Sample[stepIdx] = ( zb0Sample[stepIdx]
157                         - sum( Mx / SDx * zbSample[stepIdx,] ) )
158   for ( j in 1:nPredictors ) {
159     bSample[stepIdx,j] = zbSample[stepIdx,j] / SDx[j]
160   }
161 }
162
163 # Examine sampled values, z scale:
164 windows()
165 thinIdx = ceiling(seq(1,chainLength,length=700))
166 pairs( cbind( zb0Sample[thinIdx] , zbSample[thinIdx,] ) ,
167         labels=c( "zb0" , paste("zb",predictorNames,sep="" ) ) )
168 # Examine sampled values, original scale:
169 windows()
170 pairs( cbind( b0Sample[thinIdx] , bSample[thinIdx,] ) ,

```

```

171     labels=c( "b0", paste("b_",predictorNames,sep="") ) )
172 dev.copy2eps(file=paste(fname,"PostPairs.eps",sep=""))
173
174 # Display the posterior :
175 windows(3.5*(1+nPredictors),2.75)
176 layout( matrix(1:(1+nPredictors),nrow=1) )
177 histInfo = plotPost( b0Sample , xlab="b0 Value" , compVal=N      ULL , breaks=30 ,
178                      main=paste( "logit(p(", predictedName ,
179                               "=1)) when predictors = zero" , sep="" ) )
180 for ( blidx in 1:nPredictors ) {
181   histInfo = plotPost( bSample[blidx] , xlab=paste("b",blidx      x," Value",sep="" ) ,
182                      compVal=0.0 , breaks=30 ,
183                      main=paste(predictorNames[blidx]) )
184 }
185 dev.copy2eps(file=paste(fname,"PostHist.eps",sep=""))
186
187 # Plot data with .5 level contours of believable logistic sur      faces.
188 # The contour lines are best interpreted when there are only t      wo predictors.
189 for ( p1idx in 1:(nPredictors-1) ) {
190   for ( p2idx in (p1idx+1):nPredictors ) {
191     windows()
192     xRange = range(x[,p1idx])
193     yRange = range(x[,p2idx])
194     # make empty plot
195     plot( NULL , NULL , main=predictedName , xlim=xRange , ylim= yRange ,
196           xlab=predictorNames[p1idx] , ylab=predictorNames[p2id      x] )
197     # Some of the 50% level contours from the posterior sample.
198     for ( chainIdx in ceiling(seq( 1 , chainLength , length=20 )) ) {
199       abline( -( b0Sample[chainIdx]
200               + if (nPredictors>2) {
201                 bSample[chainIdx,c(-p1idx,-p2idx)]*Mx[c(-p1idx,-p2i      dx)]
202               } else { 0 } )
203               / bSample[chainIdx,p2idx] ,
204               -bSample[chainIdx,p1idx]/bSample[chainIdx,p2idx] ,
205               col="grey" , lwd = 2 )
206     }
207     # The data points:
208     for ( yVal in 0:1 ) {
209       rowIdx = ( y == yVal )
210       points( x[rowIdx,p1idx] , x[rowIdx,p2idx] , pch=as.cha      cter(yVal) ,
211               cex=1.75 )
212     }
213     dev.copy2eps(file=paste(fname,"PostContours",p1idx,      p2idx,".eps",sep=""))
214   }
215 }
216
217 #-----#
218
219 # MLE logistic regression:
220 glmRes = glm( datalist$y ~ as.matrix(x) , family=binomial(      logit) )
221 show( glmRes )

```

20.5.2 Logistic ANOVA code

This program implements the model shown in Figure 20.5.

(LogisticOnewayAnovaBrugs.R)

¹ graphics.off()


```

61   x = CondOfSubj
62   n = nTrlOfSubj
63   z = nCorrOfSubj
64   Ntotal = length(z)
65   xnames = c("Rev","Rel","Irr","Cmp")
66   NxLvl = length(unique(x))
67   contrastList = list( REVvREL = c(1,-1,0,0) , RELvIRR = c(0,1      ,-1,0) ,
68                      IRRvCMP = c(0,0,1,-1) , CMPvOneRel = c(0,-1/2,-1/2,1) ,
69                      FourExvEightEx = c(-1,1/3,1/3,1/3) ,
70                      OneRelvTwoRel = c(-1/2,1/2,1/2,-1/2) )
71 }
72
73 if ( dataSource == "Random" ) {
74   fnroot = paste( fnroot , dataSource , sep="" )
75   #set.seed(47405)
76   a0true = -0.5
77   atrue = c( 0.8 , -0.3 , -0.5 ) # sum to zero
78   ktrue = 100
79   subjPerCell = 50
80   nPerSubj = 100
81   datarecord = matrix( 0, ncol=3 , nrow=length(atrue)*subjP      erCell )
82   colnames(datarecord) = c("x","z","n")
83   rowidx = 0
84   for ( xidx in 1:length(atrue) ) {
85     for ( subjidx in 1:subjPerCell ) {
86       rowidx = rowidx + 1
87       datarecord[rowidx,"x"] = xidx
88       mu = sigmoid(a0true+atru

```

```

120     NxLvl = NxLvl
121   )
122 # Get the data into BRugs:
123 modelData( bugsData( dataList ) )
124
125 #-----
126 # INTIALIZE THE CHAINS.
127
128 # Autocorrelation within chains is large, so use several cha      ins to reduce
129 # degree of thinning. But we still have to burn-in all the chai      ns, which takes
130 # more time with more chains.
131 nchain = 10
132 modelCompile( numChains = nchain )
133
134 if ( F ) {
135   modelGenInits() # often won't work for diffuse prior
136 } else {
137   # initialization based on data
138   theData = data.frame( pr=.01+.98*datalist$z/datalist$n      ,
139                         x=factor(x,labels=xnames) )
140   a0 = mean( logit(theData$pr) )
141   a = aggregate( logit(theData$pr) , list( theData$x ) , mean )      [,2] - a0
142   mGrp = aggregate( theData$pr , list( theData$x ) , mean )[2]
143   sdGrp = aggregate( theData$pr , list( theData$x ) , sd )[2]
144   kGrp = mGrp*(1-mGrp)/sdGrp^2 - 1
145   k = mean(kGrp)
146   genInitList <- function() {
147     return(
148       list(
149         a0 = a0 ,
150         a = a ,
151         aSDunabs = sd(a) ,
152         theta = theData$pr ,
153         k = k
154       )
155     )
156   }
157   for ( chainIdx in 1 : nchain ) {
158     modelInits( bugsInits( genInitList ) )
159   }
160 }
161
162 #-----
163 # RUN THE CHAINS
164
165 # burn in
166 BurnInSteps = 10000
167 modelUpdate( BurnInSteps )
168 # actual samples
169 samplesSet( c( "a0" , "a" , "aSD" , "k" ) )
170 stepsPerChain = ceiling(2000/nchain)
171 thinStep = 750
172 modelUpdate( stepsPerChain , thin=thinStep )
173
174 #-----
175 # EXAMINE THE RESULTS
176
177 source("plotChains.R")
178 source("plotPost.R")

```

```

179
180 checkConvergence = T
181 if ( checkConvergence ) {
182   sumInfo = plotChains( "a0" , saveplots=F , filenameroot=fn      root )
183   sumInfo = plotChains( "a" , saveplots=F , filenameroot=fnr     oot )
184   sumInfo = plotChains( "aSD" , saveplots=F , filenameroot=f    nroot )
185   sumInfo = plotChains( "k" , saveplots=F , filenameroot=fnr     oot )
186 }
187
188 # Extract and plot the SDs:
189 aSDSample = samplesSample("aSD")
190 windows()
191 par( mar=c(3,1,2.5,0) , mgp=c(2,0.7,0) )
192 histInfo = plotPost( aSDSample , xlab="a SD" , main="a SD" , b     reaks=30 )
193 dev.copy2eps(file=paste(fnroot,"SD.eps",sep=""))
194
195 # Extract a values:
196 a0Sample = samplesSample( "a0" )
197 chainLength = length(a0Sample)
198 aSample = array( 0 , dim=c( datalist$NxLvl , chainLength ) )
199 for ( xidx in 1:datalist$NxLvl ) {
200   aSample[xidx,] = samplesSample( paste("a[",xidx,"]",se      p="") )
201 }
202
203 # Convert to zero-centered b values:
204 mSample = array( 0 , dim=c( datalist$NxLvl , chainLength ) )
205 for ( stepIdx in 1:chainLength ) {
206   mSample[stepIdx ] = ( a0Sample[stepIdx] + aSample[,stepI     dx] )
207 }
208 b0Sample = apply( mSample , 2 , mean )
209 bSample = mSample - matrix(rep( b0Sample ,NxLvl),nrow=NxL vl,byrow=T)
210
211 # Plot b values:
212 windows(datalist$NxLvl*2.75,2.5)
213 layout( matrix( 1:datalist$NxLvl , nrow=1 ) )
214 par( mar=c(3,1,2.5,0) , mgp=c(2,0.7,0) )
215 for ( xidx in 1:datalist$NxLvl ) {
216   plotPost( bSample[xidx,] , breaks=30 ,
217             xlab=bquote(beta[.(xidx)]) ,
218             main=paste(xnames[xidx]) )
219 }
220 dev.copy2eps(file=paste(fnroot,"b.eps",sep=""))
221
222 # Consider parameter correlations:
223 kSample = samplesSample("k")
224 windows()
225 pairs( cbind( b0Sample , t(bSample) , kSample ) , labels=c("      b0",xnames,"k") )
226
227 # Display contrast analyses
228 nContrasts = length( contrastList )
229 if ( nContrasts > 0 ) {
230   nPlotPerRow = 5
231   nPlotRow = ceiling(nContrasts/nPlotPerRow)
232   nPlotCol = ceiling(nContrasts/nPlotRow)
233   windows(3.75*nPlotCol,2.5*nPlotRow)
234   layout( matrix(1:(nPlotRow*nPlotCol),nrow=nPlotRow,n      col=nPlotCol,byrow=T) )
235   par( mar=c(4,0.5,2.5,0.5) , mgp=c(2,0.7,0) )
236   for ( clIdx in 1:nContrasts ) {
237     contrast = matrix( contrastList[[clIdx]],nrow=1) # make it      a row matrix

```

```

238     inclIdx = contrast!=0
239     histInfo = plotPost( contrast %*% bSample , compVal=0 , brea    ks=30 ,
240                           xlab=paste( round(contrast[inclIdx],2) , xnames[inclIdx] ,
241                                     c(rep("+",sum(inclIdx)-1),"") , collapse=" " ) ,
242                           cex.lab = 1.5 ,
243                           main=paste( "Contrast:", names(contrastList)[cldx] ) )
244   }
245   dev.copy2eps(file=paste(fnroot,"xContrasts.eps",sep      =""))
246 }
247
248 #=====
249 # Do NHST ANOVA:
250
251 theData = data.frame( y=z/n , x=factor(x,labels=xnames) )
252 aovresult = aov( y ~ x , data = theData )
253 cat("\n-----\n")
254 print( summary( aovresult ) )
255 cat("\n-----\n")
256 print( model.tables( aovresult , "means" ) , digits=4 )
257 windows()
258 boxplot( y ~ x , data = theData )
259 cat("\n-----\n")
260 print( TukeyHSD( aovresult , "x" , ordered = FALSE ) )
261 windows()
262 plot( TukeyHSD( aovresult , "x" ) )
263 if ( F ) {
264   for ( xIdx1 in 1:(NxLvl-1) ) {
265     for ( xIdx2 in (xIdx1+1):NxLvl ) {
266       cat("\n-----\n")
267       cat( "xIdx1 = " , xIdx1 , ", xIdx2 = " , xIdx2 ,
268            " , M2-M1 = " , mean(score[x==xIdx2])-mean(score[x==xIdx1] ) ,
269            "\n" )
270       print( t.test( score[ x == xIdx2 ] , score[ x == xIdx1 ] ) )
271     }
272   }
273 }
274 cat("\n-----\n")
275
276 #=====

```

20.6 Exercises

Exercise 20.1. [Purpose: Correlated predictors in logistic regression.] For this exercise, we'll use some fictitious data regarding whether or not a patient had a heart attack within a year after a check up that included measurements of systolic and diastolic blood pressures, cholesterol, weight, height, and age. The data are completely fabricated and fictional, for illustrative purposes only. The data are generated by the `BloodDataGenerator.R`, which is available at the author's web site.

(A) For this part, we will generate two data sets. In one set, the predictors have correlations of zero with each other. In the second set, the two predictors have a very strong positive correlation, but still zero correlation with all the other predictors. For both sets, the true regression coefficient on the first predictor is zero, but the true regression coefficient on the second predictor is positive. Use the program `BloodDataGenerator.R` to generate these two data sets, by selecting the relevant `correlationMatrix` at the top of the pro-

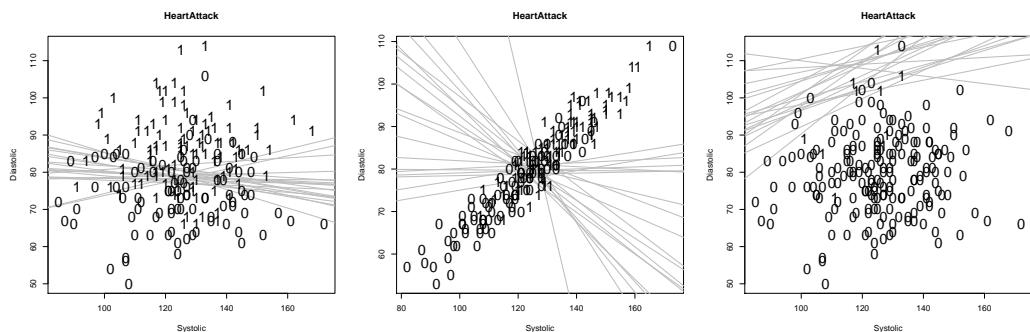


Figure 20.7: In all panels, the true regression coefficient on Systolic is zero, and the true regression coefficient on Diastolic is positive (non-zero). In the left panel the predictors are uncorrelated, and the credible $p = .05$ level contours are fairly certain. In the middle panel, the predictors are highly correlated, and the credible $p = .05$ level contours are very uncertain. In the right panel, only a small percentage of the data has value 1, and again the level contours are very uncertain. See Exercises 20.1 and 20.2.

gram. For each data set, conduct a Bayesian analysis using `gridogram` in Section 20.5.1 (`MultipleLogisticRegressionBrugs.R`). Be sure to specify that the only predictors are the first two, namely systolic and diastolic pressure. Include posterior contour plots, which should look similar to those in Figure 20.7. Be sure also to include the histograms of the marginal posterior distributions. In which case are the HDI's narrower?

(B) For this part, we generate a different fictitious data set, in which the first two predictors (systolic and diastolic) have zero correlation with each other, but they are correlated with another predictor (weight). The first predictor (systolic) has a true regression coefficient of zero, but the second (diastolic) and third (weight) predictors have positive regression coefficients. Run the Bayesian logistic regression analysis only the first two predictors. What is the estimate of the regression coefficient on the first predictor? Include the histograms of the marginals of the posterior. Now re-run Bayesian logistic regression analysis including all six predictors. What is the estimate of the regression coefficient on the first predictor? Include the histograms of the marginals of the posterior (for which you might need to manually stretch the graph window or change character sizes in the plot). Why are the estimates of the regression coefficient for the first predictor different in the two analyses?

Exercise 20.2. [Purpose: Small proportion of 1's in logistic regression.] For this exercise, we again use fictitious data generated by the program `gridDataGenerator.R`. Set all the inter-predictor correlations to zero, by selecting the appropriate correlation matrix at the top of the program. Change the `proportionOnes` constant to 0.02. This causes the proportion of 1's in the data to be very small. In other words, the threshold for sigmoidal function is set very high, so that the weighted sum of the predictors must be large before $y=1$ gets very large. This sort of situation, in which there are few, is not unusual in real data. In heart attack data, for example, there will be few people who have a heart attack in the year following a routine check up. Run a Bayesian logistic regression on the data, using only the first two predictors. What are the estimates of the regression coefficients? Include the histograms of the marginals of the posterior, and include a posterior contour plot, which should look similar to the right panel of Figure 20.7. Compare the widths of the HDI's to

the results of the first part of the previous exercise.

Exercise 20.3. [Purpose: Logistic ANOVA without assuming homogeneity of variance.] Extend the program of Section 20.5.2(LogisticOneWayAnovaBrugs.R) so that it allows different certainty () parameters for each condition, with a hyperdistributionas conditions, just as in the model on the right side of Figure 9.17, p. 183. Apply the model to the Iteration-condensation data, and compare the results to the results of Exercise 9.2, p. 191. If you already did Exercise 18.3, you can simply modify your program for that Exercise. Programming hint: Here's what the model section might look like: (LogisticOneWayAnovaHeteroVarBrugs.R)

```

11 model {
12   for ( i in 1:Ntotal ) {
13     z[i] ~ dbin( theta[i] , n[i] )
14     theta[i] ~ dbeta( aBeta[x[i]] , bBeta[x[i]] )I(0.001,0.99      9)
15   }
16   for ( j in 1:NxLvl ) {
17     aBeta[j] <- mu[j] * k[j]
18     bBeta[j] <- (1-mu[j]) * k[j]
19     mu[j] <- 1 / ( 1 + exp( -( a0 + a[j] ) ) )
20     a[j] ~ dnorm( 0.0 , atau )
21     k[j] ~ dgamma( skappa , rkappa )
22   }
23   a0 ~ dnorm( 0 , 0.001 )
24   atau <- 1 / pow( aSD , 2 )
25   aSD <- abs( aSDunabs ) + .1
26   aSDunabs ~ dt( 0 , 0.001 , 2 )
27   skappa <- pow(mg,2)/pow(sg,2)
28   rkappa <- mg/pow(sg,2)
29   mg ~ dunif(0,50)
30   sg ~ dunif(0,30)
31 }
```

Chapter 21

Ordinal Predicted Variable

Contents

21.1	Ordinal probit regression	472
21.1.1	What the data look like	472
21.1.2	The mapping from metric to ordinal	472
21.1.3	The parameters and their priors	74 4
21.1.4	Standardizing for MCMC efficiency	475
21.1.5	Posterior prediction	547
21.2	Some examples	476
21.2.1	Why are some thresholds outside the data?	478
21.3	Interaction	480
21.4	Relation to linear and logistic regression	481
21.5	R code	481
21.6	Exercises	486

The winner is first, and that's all that he knows, whether
Won by a mile or won by a nose. But
Second recalls every inch of that distance, in
Vivid detail and with haunting persistence.

Very often the predicted variable is ordinal, such as a rating on a scale from 1 to 5: Rate how much you agree with this statement: “Bayesian methods are most informative, useful, and rational way to analyze data”, where 5 = strongly agree, 4 = moderately agree, 3 = neither agree nor disagree, 2 = moderately disagree, 1 = strongly disagree. This sort of response scale is often called Likert (pronounced LICK-ert) scale (Likert, 1932), and it typically has 5, 7 or more levels (Dawes, 2008). Vast amounts of survey data use Likert rating scales.

We often want to predict the ordinal measure on the basis of predictors. For example, predictors might include age of respondent, number of years of education, CPU speed of the computer used by respondent, and so forth. To predict values from metric predictors, we need a way of mapping a weighted linear combination of predictors onto ordinal scale values. The linking function we will use is `alprobit` function with thresholds. This type of analysis is therefore called ordinal probit regression. It appeared in the bottom row of Table 14.1, p. 312.

Kruschke, J. K. (2010) Doing Bayesian Data Analysis: A Tutorial with R and BUGS. Academic Press Elsevier. Copyright 2010 by John K. Kruschke. Draft of May 11, 2010. Please do not distribute this preliminary draft. If you report Bayesian analyses based on this book, please do cite it!_

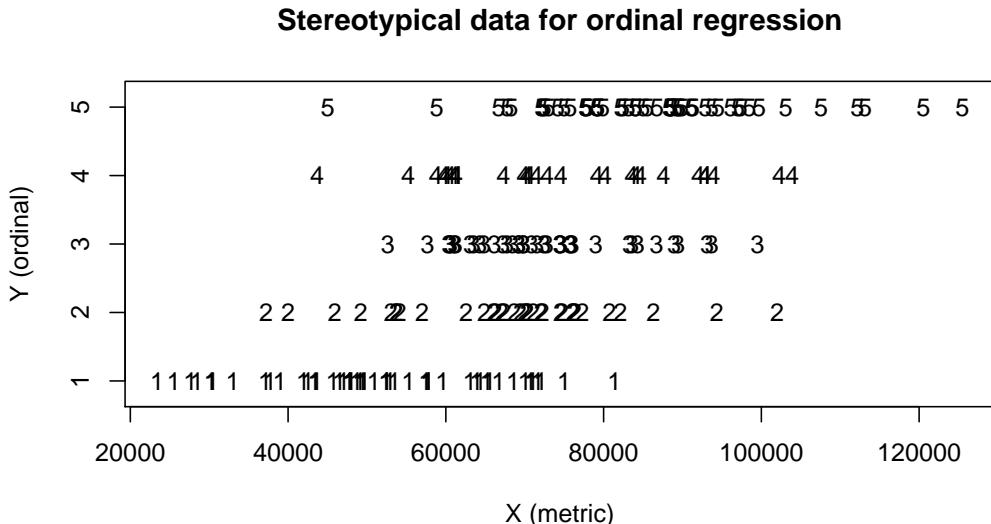


Figure 21.1: Some stereotypical (ctitious) data of the type modeled by ordinal regression. Points are labeled with their value instead of plotted merely as dots. Although their values are plotted at equal intervals on the x-axis, we really only know their order, not their separation. The conceptual emphasis is on vertical slices through the scatterplot: At any given value of x , what is the probability distribution across the discrete values?

21.1 Ordinal probit regression

21.1.1 What the data look like

Some stereotypical (ctitious) data that could be modeled by ordinal probit regression are shown in Figure 21.1. The predicted variable, y , is from a 5-point ordinal scale. The predictor variable, x , is metric. If you like concrete examples, imagine that y is a rating of overall happiness, and x is annual income in dollars. The key aspect of y that makes it ordinal, not metric, is that we don't know whether the distance from 1 to 2 is the same as the distance from 2 to 3, or from 3 to 4, etc. For example, is an increase in happiness from 1 to 2 the same amount of potential happiness as an increase in happiness from 2 to 3? It may be, but we don't know. All we know is the order of the ratings, not the distances between them.

Regarding Figure 21.1, notice that when x is relatively small, around 40,000, then the y values tend to be 1's or 2's. When x is relatively large, around 100,000, then the y values tend to include more 4's and 5's. But notice also that there is a lot of random variability, in the sense that there is a range of values at any given x value.

21.1.2 The mapping from metric x to ordinal y

Our goal is to construct a model for this sort of data. To accomplish the goal, we start with linear regression on x and then link the continuous prediction to ordinal values via a thresholded cumulative-normal function. Figure 21.2 shows the underlying mappings from x to y . The right side of the diagram shows simple linear regression on the metric

Observed
ordinal
value

$$\text{"5"} \quad p('5'|\mu) = 1 - \Phi((\theta_4 - \mu)/\sigma)$$

$$\text{"4"} \quad p('4'|\mu) = \Phi((\theta_4 - \mu)/\sigma) - \Phi((\theta_3 - \mu)/\sigma)$$

$$\text{"3"} \quad p('3'|\mu) = \Phi((\theta_3 - \mu)/\sigma) - \Phi((\theta_2 - \mu)/\sigma)$$

$$\text{"2"} \quad p('2'|\mu) = \Phi((\theta_2 - \mu)/\sigma) - \Phi((\theta_1 - \mu)/\sigma)$$

$$\text{"1"} \quad p('1'|\mu) = \Phi((\theta_1 - \mu)/\sigma) - 0$$

Underlying metric value

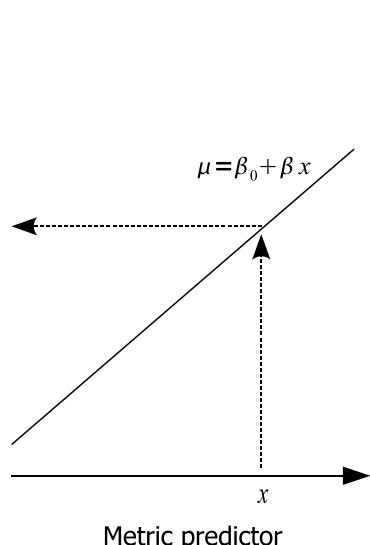


Figure 21.2: The underlying mappings in ordinal regression. The lower right shows the predictor value, which gets mapped to an underlying metric value via a linear relation. There is normally distributed noise in the underlying metric value, which has mean and standard deviation. The underlying metric value is mapped to an ordinal value depending on which threshold falls between. The thresholds are not evenly spaced. The probability falling between thresholds is given by the cumulative normal function,

predictor value x , gets mapped to an underlying metric value via the simple linear relationship, $\mu = \beta_0 + \beta_1 x$. The underlying metric value is noisy, distributed around a normal distribution with standard deviation, which is drawn sideways in the diagram. So far, this is merely simple linear regression.

The new part, on the left of Figure 21.2, gets us from the underlying metric value to an observed ordinal value. To achieve this mapping, the underlying metric scale is carved into intervals by thresholds, located at scale values $\theta_1, \theta_2, \dots$, and so on, as shown. The threshold values are estimated from the data. When the underlying metric value falls between θ_n and θ_{n+1} , an ordinal value n is generated. Because the underlying metric value is normally distributed, the probability that it falls in any interval is given by the cumulative normal distribution.

As you may recall from Section 14.1.7.2 and Figure 14.8, the cumulative normal function is denoted $\Phi(z)$. The value of $\Phi(z)$ indicates how much of the normal distribution falls to the left of z , where z is a standardized value. For example, $\Phi(0) = .5$, because half of the standardized normal distribution falls to the left of 0. The cumulative normal function is $\text{pnorm}(z)$. In BUGS, the cumulative normal function is $\Phi(z)$.

To determine how much of the normal distribution falls between two values z_n and z_{n+1} , we simply compute the difference, $\Phi(z_{n+1}) - \Phi(z_n)$. We apply this result to the thresholds in Figure 21.2, but first standardize the thresholds relative to the mean and standard deviation of the noise distribution. Thus, the mass of the normal distribution between z_n and z_{n+1} is $\Phi(z_{n+1}) - \Phi(z_n)$. This mass is the probability of ordinal value n .

summary, and as written in Figure 21.2 for each interval,

$$p("n"j) = \max(0; ((n_0) =) - ((n_{-1}) =)) \quad (21.1)$$

where ‘‘h’’, in quotes, means that the values emitted by the subject. A more traditional, but awkward, way of writing $p("n"j) = \dots$ is $p(y=n_j) = \dots$. For the lowest and highest ordinal values, the outside thresholds are at negative infinity in nity. The area under the normal distribution to the left of positive infinity is 1 hence Equation 21.1 becomes 1 $((n_{-1}) =)$ at the highest ordinal value.

Because the mapping from underlying metric value to ordinal value uses a cumulative normal function, this type of regression might be called “cumulative normal regression” or “phi regression”, by analogy to “logistic regression”. But not. Instead, it’s called “probit regression” to indicate that a cumulative normal is used as the link function, and it’s called “ordinal probit regression” to indicate that the predicted value is ordinal.

21.1.3 The parameters and their priors

Inspection of Figure 21.2 reveals the parameters that are used in ordinal probit regression. There are the usual parameters for linear regression including α , β_j for each predictor, and σ . Also, there are the threshold parameters γ_L through γ_1 , where L is the number of levels in the ordinal scale. It turns out that the intercept trades off with the thresholds, because we can add any constant to the intercept and subtract that constant from the thresholds and have an equivalent model. Therefore we can set the intercept to an arbitrary convenient constant. Also, the noise standard deviation trades off with the slopes and thresholds, because if we multiple the noise by a constant we can divide all the slopes and thresholds by that constant and end up with an equivalent model (also translating the thresholds if the intercept is not zero). In summary, we can set the slopes on the predictors, β_j , and the thresholds, γ_n , but we can set the intercept α and the noise σ at convenient constants. (Which constants are convenient? This is answered in the description of the R code in Section 21.5(`ordinalProbitRegressionBrugs.R`).)

To estimate these parameters, we must define our prior ~~uniform~~. The prior on the regression coefficients is the same as for linear or logistic regression. For example, each regression coefficient could have a normal prior. If there are several predictors, then a hyper-prior could be posited to express the mutual informativeness of the regression coefficients, as in Figure 17.6.

The prior on the thresholds can be (mildly) informed by the scale of the predicted value and by the fact that the thresholds need to be in order, ~~they~~ it’s best to be that $\gamma_n > \gamma_{n-1}$. Because the y values are 1 through L , and we will set up the underlying metric value to match that scale, it is reasonable to believe that the thresholds ~~should~~ be approximately 1.5, 2.5, ..., $L-0.5$. Our uncertainty in those values is moderate, not extreme, therefore we may opt to put a normal prior on each threshold with a standard deviation about 1 unit. This sort of prior, in which each threshold has a normal distribution that is independent of the others, makes it quite possible to have thresholds that are out of order. There are various ways to express prior knowledge that the thresholds must be ordered. One simple way is to put the knowledge into the likelihood function itself, by changing Equation 21.1 to this:

$$p("n"j) = \max(0; ((n_0) =) - ((n_{-1}) =))$$

Notice that if $n < n_1$, then the probability of " n " is zero. That would be fine if the value n never occurs in the data, but if it does, then threshold ~~sizes~~ will never be found to be credible.

The prior on the regression coefficients is explained in Exercise 21.2. In that exercise you will also explore the robustness of the posterior to ~~change~~ in the prior (and you'll find that the posterior is very robust).

21.1.4 Standardizing for MCMC efficiency

To reduce correlation among regression coefficients, and thereby make the MCMC sampling more efficient, the predictor values are standardized before the parameters are estimated. In other words, the x_j values are transformed according to Equation 16.1: $x_j^{(z)} = (x_j - \bar{M}_j)/SD_j$. For the formulas in this section, the standardized parameters for the standardized predictors all have a superscript, like this $x_j^{(z)}$, $\beta_j^{(z)}$, and so on. After estimation, the parameters must be transformed back to the original scales of the predictors. To do this back-transformation, consider the following system of equalities:

$$p("n^l j x) = \left(\begin{smallmatrix} (z) & (z) \\ n & n \end{smallmatrix} \right) = \left(\begin{smallmatrix} (z) \\ n-1 \end{smallmatrix} \right) = \left(\begin{smallmatrix} (z) \\ (z) \end{smallmatrix} \right)$$

that's Equation 21.1 with superscripts

$$p("n!jx) = \binom{z}{n} = \frac{z!}{n!(z-n)!}$$

notice “n” on the left-hand side!

where all the parameters without superscripts are on the original scale. Why is θ_0 0? Because it's arbitrary: we could set it to any value, adding to the sum on the far right and subtracting it out of the thresholds. Similarly, the θ_{ref} is arbitrary; we could set it to any value if we multiplied all the thresholds and regression coefficients by that value, and then divided by before passing into the function.

21.1.5 Posterior prediction

The slope and threshold values from Equation 21.2 are arbitrary determined only relative to 0 and 1. Their usefulness is being able to easily compute probabilities

from the original scale values, as follows:

$$p(n^j | x) = \frac{1}{\pi} \left(\frac{\exp(-x_j)}{1 + \exp(-x_j)} \right)^{n^j} \left(\frac{\exp(x_j)}{1 + \exp(x_j)} \right)^{m - n^j} \quad (21.3)$$

This expression (in Equation 21.3) is merely Equation 21.2 written with the original-scale parameters set where the underbraces indicate they be. Exercise 21.3 has you use this equation to make predictions.

21.2 Some examples

Consider the data shown in the top panel of Figure 21.3. There are two metric predictors, labeled X_1 and X_2 in the graph, and the data to be predicted, ranging from 1 to 7. Each data point is labeled with its value. These data were randomly generated but with very little noise in the value, so that the transitions from one level to the next would be very clear. Notice that as X_1 increases, the value tends to decrease. Hence the regression coefficient on X_1 should be negative. Notice also that as X_2 increases, the values tends to increase. Hence the regression coefficient on X_2 should be positive.

The program in Section 21.5 (`ordinalProbitRegressionBrugs.R`), which does Bayesian ordinal probit regression, was run with these data. Aspects of the resulting posterior are shown in the middle and lower portions of Figure 21.3. The first histogram indicates credible values of α_1 , which are all negative, as inspection of the trends in the data suggested. The second histogram indicates credible values of α_2 , which are all positive, again as inspection of the trends in the data suggested.

The next six histograms show credible values of the thresholds. One thing to notice about the distributions of the thresholds is that they appear to overlap tremendously: The histogram for α_1 covers much of the same range as the histogram for α_2 , and so on. This overlap would seem to suggest a violation of the necessary ordering of the thresholds, because it seems α_1 could easily be larger than α_2 . What's wrong? Nothing, because credible values of the thresholds are very strongly correlated. The last part of Figure 21.3 shows pairwise scatterplots of the parameter values, and reveals extreme correlation. In fact, if you plot histograms of the posterior $\alpha_{1,2}$, you will find that the differences do not overlap zero at all; i.e., there is no violation of ordering the thresholds. Exercise 21.1 provides more details regarding how to do this.

The plot of the data, in the top panel of Figure 21.3, also includes examples of level contours from the posterior. A set of six (i.e., 1) level contours comes from a single step in the MCMC chain. The six level contours correspond to six thresholds, for the α_1 and α_2 values at that step in the chain. The level contours within each have the same slope because they all correspond to the same α_2 values. Different sets of contours have different slopes because they have different α_1 and α_2 values. Because these data were generated with very little noise in the values, there is little uncertainty in the slopes or thresholds. One indicator of the lack of noise is that fact that between level contours in the plot there is mostly just one value of

The level contours in the top panel of Figure 21.3 were defined as follows. Consider the line that marks the transition from the "1" region to the "0" region. The transition occurs when $p(n^j | x)$ is 0.5, which means that $\text{prob}(n^j | x) = 0.5$. From Equation 21.3

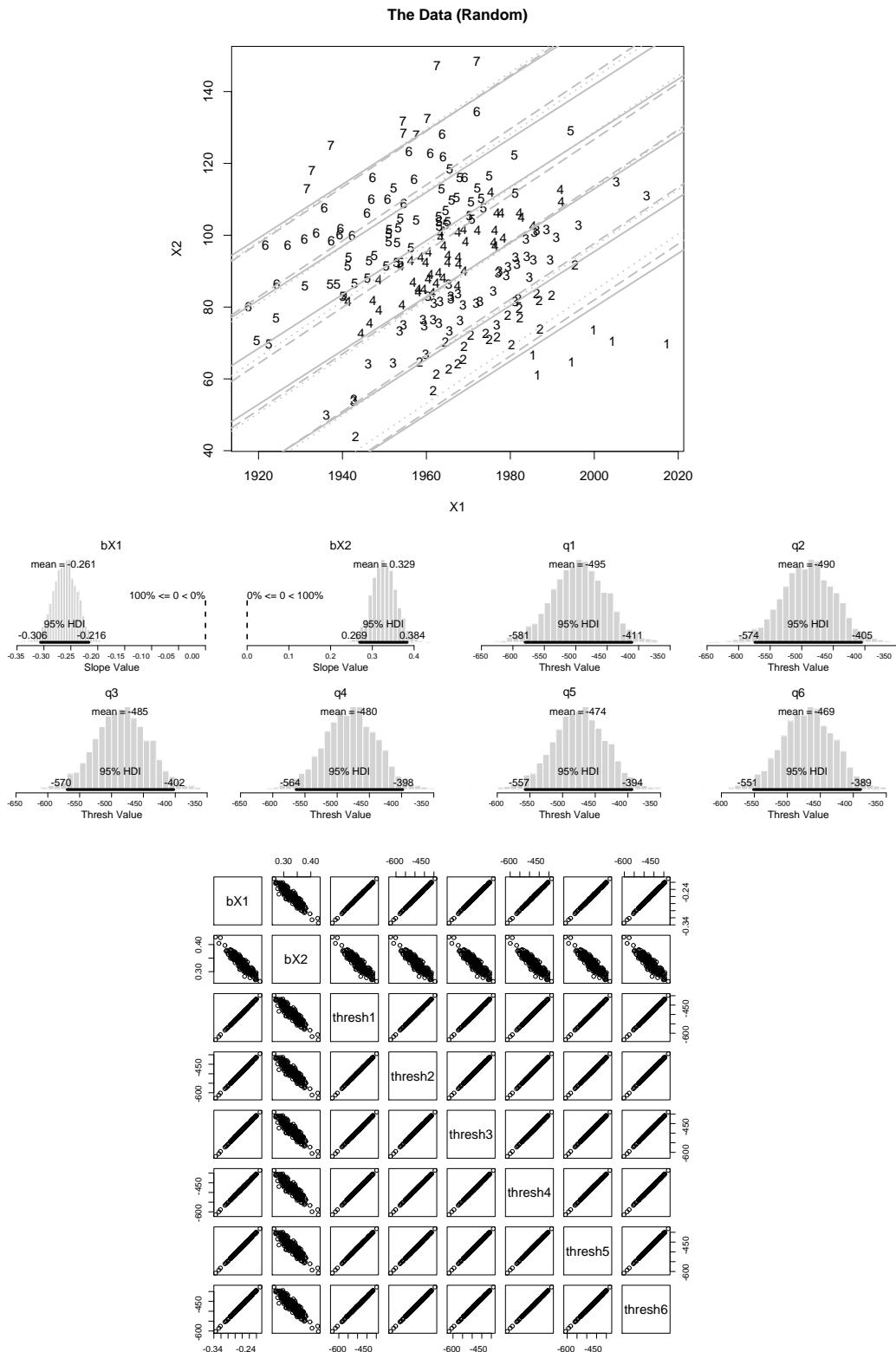


Figure 21.3: Upper panel: Random ordinal data plotted against predictors, superimposed with a smattering of threshold lines from the prior. Histograms and lower scatterplots show posterior on the regression coefficients and thresholds.

(andor Eqn. 21.2), we know

$$\text{probit } p(\text{``} n^j x \text{''}) = \sum_j x_j$$

It follows that the transition line occurs for values that satisfy

$$0 = \sum_j x_j$$

$0 = n_1 x_1 + n_2 x_2$ when there are two predictors

$$x_2 = \frac{n_1}{2} + \frac{1}{2} x_1$$

The level contours in the top panel of Figure 21.3 used that; they are applied in the last few lines of the code in Section 21.5 (`OrdinalProbitRegressionBrugs.R`).

The previous example used fictitious data so that the concept could be clearly illustrated. Realistic data are typically much messier. Consider, for example, ratings of movies (Moore, 2006), as shown in the top panel of Figure 21.4. The predictor on the abscissa is the year the movie was made, and the predictor on the ordinate is length, i.e., duration, of the movie in minutes. Notice that ratings are rather noisy: when any pair of contour lines there occur many different rating values.

Despite the noise in the data, the regression coefficients on the two predictors are very credibly non-zero, as can be seen in the first two histograms in Figure 21.4. Ratings clearly tend to decrease as year of production increases, and ratings clearly tend to increase as length increases. Nevertheless, the degree of uncertainty in the regression coefficients is much larger than the previous example. This uncertainty is apparent in the smattering of posterior level contours superimposed on the data: These show the level contours vary greatly from one set to another.

The posterior distribution of thresholds again shows strong correlations, in the lower scatterplots of Figure 21.4. The thresholds never violate proper order, although exactly where they should be is rather uncertain.

21.2.1 Why are some thresholds outside the data?

In Figure 21.4, the upper panel shows a scatterplot of data with a few of the posterior contour lines superimposed. Notice that some of the extre threshold contours fall outside all the data points. For example, the lower right of the plot shows that the first and second threshold contours have almost no data below them. How can thresholds be among the most credible?

The answer lies in gaining a better intuition for the predictions of the ordinal probit model. Figure 21.5 shows the predictions of the model for six variable and thresholds set at 1 through 5. The upper panel shows the predictions; the noise is very small compared to the separation between thresholds, $\gamma_{\text{max}} 0:1$. In this case, the predictions match intuition: When x is between thresholds i and $i+1$, the most probable response is $i+1$. It's only when x is very close to a threshold that an adjacent response gains probability.

Moving down Figure 21.5, the successive panels show greater amounts of noise, i.e., larger values of γ . Consider what happens to the predicted response for $x=2:2$. When γ is small, the highest probability response is 3. As the noise increases, the dominance of this

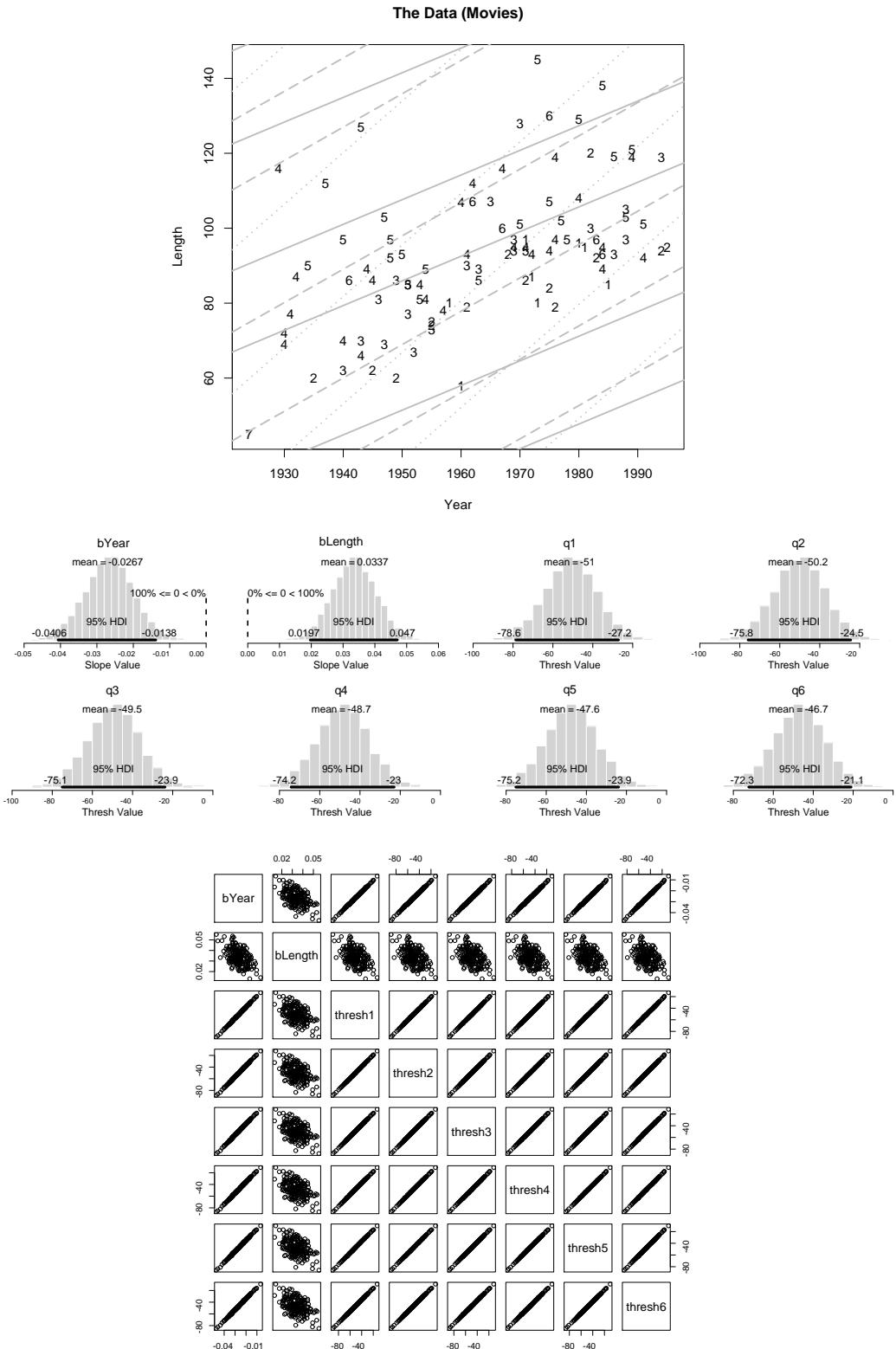


Figure 21.4: Upper panel: Movie ratings from Moore (2006) plotted against Length and Year, superimposed with a smattering of thresholds from the posterior. Histograms and lower scatterplots show posterior distributions of the regression coefficients and thresholds.

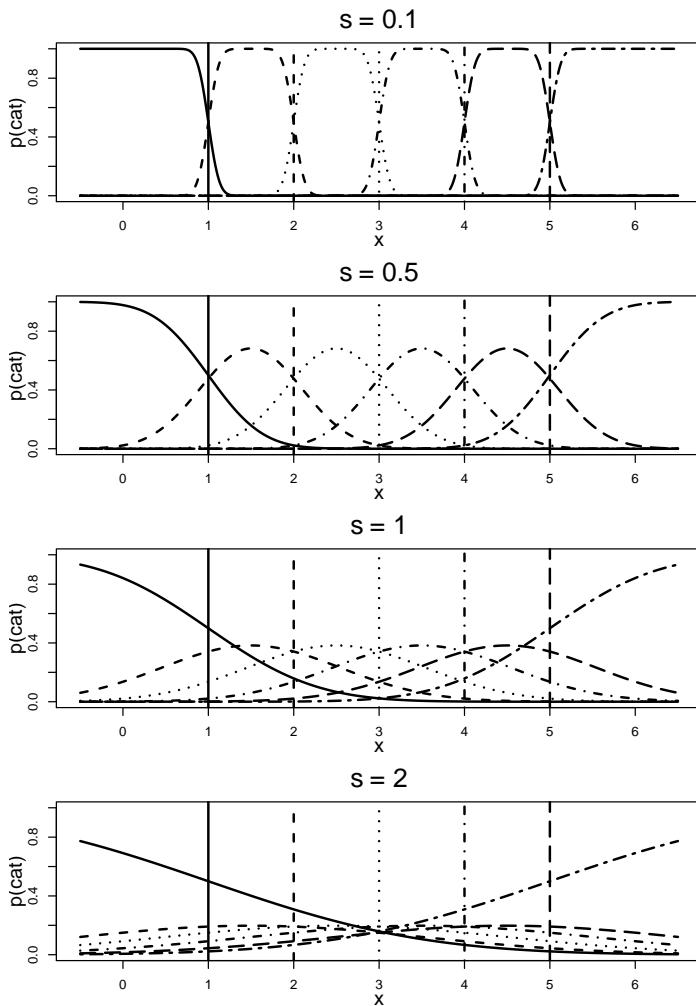


Figure 21.5: Each panel shows the probability of an ordinal response as a function of the stimulus value x . The upper panel is for small, and the lower panel is for large s .

response becomes weaker. Finally, in the bottom panel, with noise, 3 is no longer the dominant response, and instead 1 is! Thus, even though 2.2, which falls between the 2nd and 3rd thresholds, the highest probability response is

Suppose that the data consist of values scattered between $x = 2$ and $x = 4$. Suppose also that responses are very noisy, with $s = 2$, as in the lowest panel of Figure 21.5. For a point near $x = 2$, the most probable response is 1, even though the point is above the threshold for 1. For a point near $x = 4$, the most probable response is 6, even though the point is far below the threshold for 6. It is this sort of pattern that we see in the movie data of Figure 21.4, in which the responses are very noisy.

21.3 Interaction

Just as in multiple linear regression and multiple logistic regression, there can be multiplicative interaction of the predictors in ordinal probit regression. Much as the logistic function

squashes the predictions, the cumulative normal function ~~pushes~~ pushes the predictions. Hence the prediction surface, for ordinal probit regression ~~without~~, is much like the prediction surface for logistic regression, as was shown in Fig 20.4, p. 455. Ordinal probit regression is analogous to logistic regression when there ~~are~~ only two levels of the ordinal variable ~~but~~, because logistic regression has a predicted variable ~~with~~ two levels. In fact, in ordinal regression, the levels can be groups into " $< n$ " and " $> n$ ", and therefore every transition between adjacent levels can be thought of as analogous to logistic regression on two levels. Thus, when there is multiplicative interaction ~~of~~ predictors, the interaction can express ideas such as conjunction: The value is " $> n$ " when x_1 is large and x_2 is large, but not otherwise. A multiplicative interaction can also express an exclusive-OR: The value is " $> n$ " when x_1 is large or x_2 is large, but not both.

21.4 Relation to linear and logistic regression

This chapter took the core model of multiple linear regression and extended it to ordinal predicted values by mapping the linear prediction through the thresholded cumulative normal function. The only real news regarded the mathematical details dealing with the cumulative normal function, and implementing the model in R and JAGS. Otherwise, the basic concepts of multiple regression apply to ordinal probit regression.

When dealing with ordinal predicted values, some analysts ~~will~~ assume, for convenience, that the ordinal scale values are equally spaced ~~on~~ interval scale. With this assumption, the problem reverts to ordinary linear regression. Often the results of such an analysis will be very similar to ordinal regression, in terms of the ~~probabilities~~ of each ordinal level, the ordinal regression gives a more direct answer (via Eqn 21.3).

When the predicted variable has only two levels, it is tantamount to an ordinal scale with two levels (because order can be assigned arbitrarily, " $W > 0$ " or " $0 > 1$ ", without changing the results). In this case, ordinal probit regression will yield results essentially equivalent to logistic regression. The only difference will be subtle because of the slight difference in shape between the logistic function and the cumulative normal function. In fact, when there are only two levels of the predicted variable, analysts prefer probit (cumulative normal) regression to logistic regression. Because results tend to be so similar, the choice between link function usually is governed by meaningfulness for the domain of application. If the link function is to be thought of as ~~describing~~ normally distributed noise on an underlying linear relation, then the probit is more meaningful. If the analyst finds it meaningful to discuss regression coefficients in terms of log-odds (see Section 20.1.3, p. 452), then the logistic link function is more useful.

21.5 R code

Initialization of the chains uses the MLE from linear regression as a heuristic approximation to start the slope estimates for ordinal regression. The function `Im()` is used for this purpose. The MLE also returns an estimate of the intercept β_0 which should be close to the mean of the values, because the values are standardized), which is used as the arbitrary fixed value for β_0 in the model. The MLE has a residual error, which is used to obtain a corresponding standard deviation that is used as the standard deviation for β_0 .

(OrdinalProbitRegressionBrugs.R)

```

1 graphics.off()
2 rm(list=ls(all=TRUE))
3 fname = "OrdinalProbitRegressionBrugs"
4 library(BRugs)      # Kruschke, J. K. (2010). Doing Bayesian dat      a analysis:
5                      # A Tutorial with R and BUGS. Academic Press / Elsevier.
6 #-----
7 # THE MODEL.
8
9 modelstring = "
10 # BUGS model specification begins here...
11 model {
12     for( i in 1 : nData ) {
13         y[i] ~ dcat( pr[i,1:nYlevels] )
14         pr[i,1] <- phi( ( thresh[1] - mu[i] ) / sigma )
15         for ( k in 2:(nYlevels-1) ) {
16             pr[i,k] <- max( 0 , phi( ( thresh[ k ] - mu[i] ) / sigma )
17                           - phi( ( thresh[k-1] - mu[i] ) / sigma ) )
18         }
19         pr[i,nYlevels] <- 1 - phi( ( thresh[nYlevels-1] - mu[i] ) / si      gma )
20         mu[i] <- b0 + inprod( b[1:nPredictors] , x[i,1:nPredictors] )
21     }
22     bPrec <- pow( nYlevels/4 , -2 ) # max plausible slope is 1SD
23     for ( j in 1:nPredictors ) {
24         b[j] ~ dnorm(0,bPrec) # modest precision because of normali      zed x,y values
25     }
26     threshPriorPrec <- 1
27     for ( k in 1:(nYlevels-1) ) {
28         threshPriorMean[k] <- k+0.5
29         thresh[k] ~ dnorm( threshPriorMean[k] , threshPriorPrec )
30     }
31 }
32 # ... end BUGS model specification
33 " # close quote for modelstring
34 writeLines(modelstring,con="model.txt")
35 modelCheck( "model.txt" )

36 #
37 #-----
38 # THE DATA.
39
40 dataSource = c("Random","Movies")[2]
41
42 # The loading of data must produce a matrix called dataMat tha      t has
43 # one row per datum, where the first column is the ordinal pred      icted value
44 # and the 2nd - last columns are the predictor values. The colu      mns should
45 # be named.
46
47 if ( dataSource=="Random" ) {
48     fname = paste( fname , dataSource , sep="" )
49     # Generate some random toy data.
50     source( "OrdinalProbitDataGenerator.R" )
51     nYlevels = 7
52     dataMat = OrdinalProbitDataGenerator( nData = 200 ,
53                                         normPrec=200 , slope=c(-1.0,1.26) , # c(-1.0,1.26) matches Movies
54                                         thresh=c(-Inf,seq(-1.2,1.2,length=nYlevels-1),Inf) ,
55                                         nYlevels=nYlevels , makePlots=F , rndSeed=47405 )
56     # Change x values to arbitrary non-standardized scales:
57     dataMat[,2] = 1963.64 + 18.13 * dataMat[,2]
58     dataMat[,3] = 92.87 + 18.26 * dataMat[,3]

```

```

59 }
60
61 if ( dataSource=="Movies" ) {
62   fname = paste( fname , dataSource , sep="" )
63   dataFram = read.table( "Moore2006data.txt" , header=T )
64   rateVals = sort( unique( dataFram[,"Rating"] ) )
65   rankVals = match( dataFram[,"Rating"] , rateVals ) # conver      t to ranks
66   dataMat = cbind( rankVals , dataFram[,"Year"] , dataFram[,
67   colnames(dataMat) = c("Rating","Year","Length") )
68 }
69
70 # Rename for use by generic processing later:
71 nData = NROW(dataMat)
72 x = dataMat[, -1]
73 predictorNames = colnames(dataMat)[-1]
74 nPredictors = NCOL(x)
75 y = as.matrix(dataMat[, 1])
76 predictedName = colnames(dataMat)[1]
77 nYlevels = max(y)
78
79 # Re-center x values at mean, to reduce autocorrelation in MC      MC sampling.
80 # Standardize (divide by SD) to make prior-setting easier.
81 standardizeCols = function( dataMat ) {
82   zDataMat = dataMat
83   for ( colIdx in 1:NCOL( dataMat ) ) {
84     mCol = mean( dataMat[, colIdx] )
85     sdCol = sd( dataMat[, colIdx] )
86     zDataMat[, colIdx] = ( dataMat[, colIdx] - mCol ) / sdCol
87   }
88   return( zDataMat )
89 }
90 zx = standardizeCols( x )
91 # Don't standarize y because they must be integers, 1 to nYlev      els
92
93 lmlInfo = lm( y ~ zx ) # R function returns MLE
94 b0Init = lmlInfo$coeff[1]
95 b1Init = lmlInfo$coeff[-1]
96 sigmaInit = sqrt(sum(lmlInfo$res^2)/nData)
97
98 # Get the data into BUGS:
99 dataList = list(
100   x = zx ,
101   y = as.vector( y ) , # BUGS does not treat 1-column mat as vector
102   nPredictors = nPredictors ,
103   nData = nData ,
104   nYlevels = nYlevels ,
105   sigma = sigmaInit , # fixed, not estimated
106   b0 = b0Init # fixed, not estimated
107 )
108 modelData( bugsData( dataList ) )
109
110 #-----
111 # INTIALIZE THE CHAINS.
112
113 nChain = 3
114 modelCompile( numChains = nChain )
115
116 genInitList <- function() {
117   list(

```

```

118      b = bInit , # from lm(y~zx), above
119      thresh = 1:(nYlevels-1)+.5
120    )
121  }
122  for ( chainIdx in 1 : nChain ) {
123    modellnits( bugsInits( genInitList ) )
124  }
125
126 #-----
127 # RUN THE CHAINS
128
129 # burn in
130 BurnInSteps = 2000
131 modelUpdate( BurnInSteps )
132 # actual samples
133 samplesSet( c( "b" , "thresh" ) )
134 stepsPerChain = ceiling(5000/nChain)
135 thinStep = 20
136 modelUpdate( stepsPerChain , thin=thinStep )
137
138 #-----
139 # EXAMINE THE RESULTS
140
141 source("plotChains.R")
142 source("plotPost.R")
143
144 checkConvergence = T
145 if ( checkConvergence ) {
146   bSum      = plotChains( "b"      , saveplots=F , filenameroot=fname )
147   threshSum = plotChains( "thresh" , saveplots=F , filenameroot= fname )
148 }
149
150 # Extract chain values:
151 zbSamp = NULL
152 for ( j in 1:nPredictors ) {
153   zbSamp = cbind( zbSamp , samplesSample( paste("b[",j,"]", sep="") ) )
154 }
155 chainLength = NROW(zbSamp)
156 zthreshSamp = NULL
157 for ( j in 1:(nYlevels-1) ) {
158   zthreshSamp = cbind( zthreshSamp ,
159                         samplesSample(paste("thresh[",j,"]",sep="")) )
160 }
161
162
163 # Convert to original scale:
164 bSamp = zbSamp * matrix( 1/(sigmalnits*apply(x,2,sd)) , byrow=TRUE ,
165                           ncol=nPredictors , nrow=chainLength )
166 threshSamp = (1/sigmalnits) * ( zthreshSamp - b0Init +
167                                 rowSums( zbSamp * matrix( apply(x,2,mean)/apply(x,2,sd) ,
168                                               byrow=TRUE , ncol=nPredictors ,
169                                               nrow=chainLength ) ) )
170 b0 = 0
171 sigma = 1
172
173 # Scatter plots of parameter values, pairwise:
174 if ( (nPredictors+nYlevels) <= 10 ) { # don't display if too many
175   windows()
176   thinIdx = ceiling(seq(1,chainLength,length=200))

```

```

177 pairs( cbind( zbSamp[thinIdx,] , zthreshSamp[thinIdx,] ) ,
178   labels=c( paste("zb",predictorNames,sep="") ,
179             paste("zthresh",1:nYlevels,sep="") ) )
180 windows()
181 pairs( cbind( bSamp[thinIdx,] , threshSamp[thinIdx,] ) ,
182   labels=c( paste("b",predictorNames,sep="") ,
183             paste("thresh",1:nYlevels,sep="") ) )
184 dev.copy2eps(file=paste(fname,"PostPairs.eps",sep=""))
185 }
186
187 # Display the posterior:
188 nPlotPerRow = 5
189 nPlotRow = ceiling((nPredictors+nYlevels-1)/nPlotPerR      ow)
190 nPlotCol = ceiling((nPredictors+nYlevels-1)/nPlotRow)
191 windows(3.5*nPlotCol,2.25*nPlotRow)
192 layout( matrix(1:(nPlotRow*nPlotCol),nrow=nPlotRow,n      col=nPlotCol,byrow=T) )
193 par( mar=c(4,3,2.5,0) , mgp=c(2,0.7,0) )
194 for ( sIdx in 1:nPredictors ) {
195 histInfo = plotPost( bSamp[,sIdx] , xlab="Slope Value" , co      mpVal=0.0 ,
196   breaks=30 ,
197   main=bquote( b *.(predictorNames[sIdx]) ) ,
198   cex.main=1.67 , cex.lab=1.33 )
199 }
200 for ( sIdx in 1:(nYlevels-1) ) {
201 histInfo = plotPost( threshSamp[,sIdx] , xlab="Thresh Val      ue" , compVal=NULL ,
202   breaks=30 ,
203   main=bquote( theta * .(sIdx) ) ,
204   cex.main=1.67 , cex.lab=1.33 )
205 }
206 dev.copy2eps(file=paste(fname,"PostHist.eps",sep=""))
207
208 # Plot the data
209 if ( nPredictors == 2 ) {
210 windows()
211 plot( x[,1] , x[,2] , xlab=colnames(x)[1] , ylab=colnames(      x)[2] ,
212   main=paste( "The Data (" , dataSource , ")" , sep="" ) ,
213   pch=as.character(y) )
214 for ( chainIdx in round(seq(1,chainLength,len=3)) ) {
215   for ( threshIdx in 1:(nYlevels-1) ) {
216     abline( threshSamp[chainIdx,threshIdx]/bSamp[chainId      x,2] ,
217       -bSamp[chainIdx,1]/bSamp[chainIdx,2] ,
218       lwd = 2 , lty=chainIdx , col="grey" )
219   }
220 }
221 dev.copy2eps(file=paste(fname,"Data.eps",sep=""))
222
223 } # end if nPredictors == 2
224
225 # Posterior prediction.
226 xProbe = c( 1991 , 94 ) # Note order of values: x1 is year and x2 is duration.
227 # Set up a matrix for storing the values of p(y|xProbe) at each step in chain.
228 py = matrix( 0 , nrow=chainLength , ncol=nYlevels )
229 # Step through chain and compute p(y|xProbe) at each step:
230 for ( chainIdx in 1:chainLength ) {
231   yValue = 1
232   py[chainIdx,yValue] = (
233     pnorm( threshSamp[chainIdx,yValue]
234           - sum( bSamp[chainIdx,] * xProbe ) ) )
235   for ( yValue in 2:(nYlevels-1) ) {

```

```

236     py[chainIdx,yValue] = (
237         pnorm( threshSamp[chainIdx,yValue]
238             - sum( bSamp[chainIdx,] * xProbe ) )
239         - pnorm( threshSamp[chainIdx,yValue-1]
240             - sum( bSamp[chainIdx,] * xProbe ) ) )
241     }
242     yValue = nYlevels
243     py[chainIdx,yValue] = ( 1 -
244         pnorm( threshSamp[chainIdx,yValue-1]
245             - sum( bSamp[chainIdx,] * xProbe ) ) )
246   }
247 # Now average across the chain:
248 pyAve = colMeans( py )

```

21.6 Exercises

Exercise 21.1. [Purpose: Investigate posterior differences of thresholds.] Run the program in Section 21.5 (`OrdinalProbitRegressionBrugs.R`) on the Movies data (Moore, 2006). At the end of the program, notice that the posterior thresholds are stored in a matrix named `threshSamp`, one column per threshold.

(A) Make histograms of the differences between every pair of adjacent thresholds, i.e., `threshSamp[,n]-threshSamp[,n-1]` , for n in `2:nYlevels` . (Hint: Use `plotPost` with `compVal=0`) Are any differences between adjacent thresholds close to zero?

(B) Make a scatterplot of $_4$ against $_3$. (Hint: `plot(threshSamp[,4],threshSamp[,3])`) Superimpose a line that marks $_3 = _4$. (Hint: `abline(0,1)`) What is the relation of this scatterplot to the histogram you made in the previous part? In particular, what is the relation of the line in the scatterplot to the histogram?

Exercise 21.2. [Purpose: Examine robustness when prior is changed.] In the program in Section 21.5 (`OrdinalProbitRegressionBrugs.R`), the priors on the regression coefficients and thresholds are mildly informed. Consider the precision of the normal prior for the regression coefficients. Because the predictors are standardized, the precision is approximately 2 to +2. If we have a prior assumption that the steepest possible relationship between x and y maps the lowest value of x to the lowest value of y and the highest value of x to the highest value of y , then the steepest regression line would rise $nYlevels/4$, and the program makes that value the standard deviation of the prior on each regression coefficient. Consider now the prior on the thresholds. Because the thresholds are separated in the `prob` unit, and the thresholds should not violate consecutive ordering, it is reasonable to set their standard deviations no larger than 1.

(A) Set the prior on the regression coefficients so that the precision is very small, say `1.0E-6`. Run the program on the Movies data. Is the posterior much different?

(B) Set the prior on the thresholds to that the precision is very small, say `1.0E-6`. (Set the prior on the regression coefficients back to the original setting.) Run the program on the Movies data. Is the posterior much different?

Exercise 21.3. [Purpose: Predicted values for novel predictor values.] Run the program in Section 21.5 (`OrdinalProbitRegressionBrugs.R`) on the movies data (Moore, 2006). Consider a movie that was not included in the rated movies. Suppose it was made in 1991, and had a length of 94 minutes. What is the probability of each rating the movie? In other words,

what is the probability that it would be rated 7, what is the probability that it would be rated 6, and so on. Answer the same question for a movie that length of 94 minutes but was made in 1931. Hint: Use Equation 21.3, averaging across all the steps in the MCMC chain. Extra special bonus hint: Code for doing this is already included at the end of the program; your job is to understand what it's doing.

Chapter 22

Contingency Table Analysis

Contents

22.1 Poisson exponential ANOVA	49
22.1.1 What the data look like	490
22.1.2 The exponential link function	904
22.1.3 The Poisson likelihood	492
22.1.4 The parameters and the hierarchical prior	494
22.2 Examples	494
22.2.1 Credible intervals on cell probabilities	495
22.3 Log linear models for contingency tables	496
22.4 R code for Poisson exponential model	497
22.5 Exercises	504

Count me the hours that we've been together, I'll
Count you the hours I'm light as a feather, but
'Cause every hour you're all that I see, there's
No telling if there's a contingency.

Consider a situation in which we observe two nominal values for every item measured. For example, suppose there is an election, and we randomly selected people regarding their political party affiliation (which is a nominal variable) and their candidate preference (which is also a nominal variable). Presumably, the two variables are not independent, which is to say that the proportion of people who prefer candidate A should differ from one political party to another. This is the type of situation addressed in this chapter.

Traditionally, this situation comes under the rubric of chi-square tests of independence. But we will, of course, take a Bayesian approach, and consequently have no need to compute chi-square values (except for comparison with the Bayesian conclusions). There are many advantages of a Bayesian approach. As usual, a significant advantage is never having to compute a p value. Better yet is that the Bayesian analysis provides credible intervals on the conjoint probabilities and on any desired comparisons.

I will call our modeling framework Poisson exponential ANOVA because it uses a Poisson likelihood distribution with an exponential link function from an underlying ANOVA model. This terminology is not conventional, and it might be misleading if readers mistakenly infer from the term "ANOVA" that there is a metapredicted variable involved. Nevertheless, the terminology is highly descriptive of the structural elements of the model.

Kruschke, J. K. (2010) Doing Bayesian Data Analysis: A Tutorial with R and BUGS. Academic Press Elsevier. Copyright 2010 by John K. Kruschke. Draft of May 11, 2010. Please do not distribute this preliminary draft. If you report Bayesian analyses based on this book, please do cite it! "

Table 22.1: Frequencies of different combinations of hair color and eye color. Data from Snee (1974).

Hair Color	Eye Color			
	Blue	Brown	Green	Hazel
Black	20	68	5	15
Blond	94	7	16	10
Brunette	84	119	29	54
Red	17	26	14	14

22.1 Poisson exponential ANOVA

22.1.1 What the data look like

To motivate the model, we need first to understand the structure of the data. An example of the sort of data we'll be dealing with is shown in Table 22.1. The data come from a classroom poll of students at the University of Delaware (Snee 1974). Respondents reported their hair color and eye color, with each variable split into four nominal levels as indicated in Table 22.1. The cells of the table indicate the frequency with which each combination occurred in the sample. Each respondent falls into one cell of the table.

The data to be predicted are the cell frequencies. The variables are the nominal variables. This situation is analogous to two-way ANOVA, which had two nominal predictors, but had several metric values in each cell instead of a single frequency.

For data like these, we can ask a number of questions. We consider about one variable at a time, and ask questions such as, "Are there more brown-eyed people than non-brown-eyed people?", or, "Are there more blonde-haired people than non-brunette haired people?" Those questions are analogous to main effects in ANOVA. But usually we display the conjoint frequencies specifically because we're interested in the relationship between the variables. We would like to know if the distribution of frequencies across one variable depends on, or is contingent upon, the level of the other variable. For example, does the distribution of hair colors depend on eye color, and, specifically, is the proportion of blonde-haired people the same for brown-eyed and blue-eyed? These questions are analogous to interaction contrasts in ANOVA.

22.1.2 The exponential link function

The model can be motivated two different ways. One way is simply to start with the two-way ANOVA model and find a way to map the predicted value to frequency data. The predicted value from the ANOVA model could be any value from negative infinity, but frequencies are non-negative. Therefore we must transform the ANOVA predictions to non-negative values, while preserving order. The natural way to do this, mathematically, is via the exponential transformation. But this transformation only gets us to an underlying continuous predicted value, not to the probability of an integer frequency. A natural candidate for the needed likelihood distribution is the Poisson (described below), which takes a non-negative value and gives a probability for each integer from zero to infinity. But this motivation may seem a bit arbitrary, even if there's nothing wrong with it in principle.

A different motivation was introduced in Section 14.2.1, p. 310. Start by treating the cell frequencies as representative of underlying cell probabilities, and then asking whether

Table 22.2: Example of exponentiated linear model with ~~interaction~~. Margins show the values of the's, and cells show $r_c = \exp(_0 + _r + _c)$. Notice that every row has the same relative probabilities, namely, 10, 100, 1. In other words, the row and column attributes are independent.

$0 = 4:60517$	$c_1 = 0$	$c_2 = 2:30259$	$c_3 = 2:30259$
$r_1 = 0$	100	1000	10
$r_2 = 2:30259$	1000	10000	100
$r_3 = 2:30259$	10	100	1

the two nominal variables are independent of each other.~~(The definition of independence way back in Section 3.4.3, p. 39.)~~ For example, in Table 22.2, there's a particular marginal probability that hair color is black, and a particular marginal probability that eye color is brown. If hair color and eye color are independent, then the conjoint probability of black hair and brown eyes is the product of the marginal probabilities. The attributes of hair color and eye color are independent if that relationship holds for every cell in the table. Independence of hair color and eye color means that the proportion of black hair among brown-eyed people is the same as the proportion of black hair among blue-eyed people, and so on for all hair and eye colors.

To check for independence of attributes, we need to estimate the marginal probabilities of the attributes. Denote the marginal (i.e., total) frequency of the r^{th} row as $f(r)$, and denote the marginal frequency of the c^{th} column as $f(c)$. Then the estimated marginal probabilities are $\hat{f}(r) = N$ and $\hat{f}(c) = N$, where N is the total of the entire table. If the attributes are independent, then the predicted conjoint probability, $\hat{p}(r; c)$, should equal the product of the marginal probabilities, which means $\hat{p}(r; c) = p(r) \cdot p(c)$, hence $\hat{f}(r; c) = N = f(r) \cdot f(c) = N$. Because the models we've been dealing with involve additive combinations, not multiplicative combinations, we convert the multiplicative expression of independence into an additive expression by using the fact that $\log(a) = \log(a) + \log(b)$, as follows:

$$\begin{aligned} \hat{f}(r; c) &= f(r) \cdot f(c) = N \\ \hat{f}(r; c) &= f(r) \cdot f(c) - 1 = N \\ \hat{f}(r; c) &= \exp \left[\log(f(r)) + \log(f(c)) + \log(1/N) \right] \end{aligned} \quad (22.1)$$

If we abstract the form of Equation 22.1 away from the specific frequencies, we get the equation $r_c = \exp(r + c + _0)$. The idea is that whatever are the values of the resulting 's will obey multiplicative independence. An example is shown in Table 22.2. The choice of 's is shown in the margins of the table, and the resulting are shown in the cells of the table. Notice that every row has the same relative probabilities, namely, 10, 100, and 1. In other words, the row and column attributes are independent.

We have dealt before with additive combinations of row and column influences, in the context of ANOVA. In ANOVA, $_0$ is a baseline representing the overall central tendency, and $_r$ is a deviation away from baseline due to being in the row, and $_c$ is a deviation away from baseline due to being in the column. The deviations are constrained to sum to zero, and the example in Table 22.2 respects this constraint.

In ANOVA, when the cell data are not captured by an additive combination of row and column effects, we include an interaction term, denoted here as r_c . The interaction

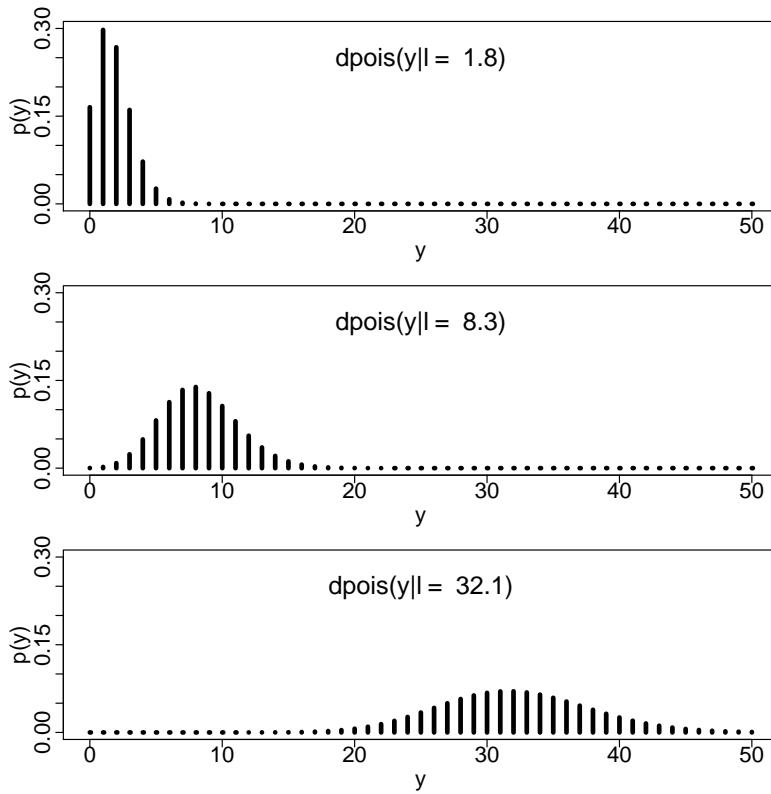


Figure 22.1: Examples of the Poisson distribution includes integers from zero to positive infinity.

terms are constrained so that every row and every column ~~sums to zero~~: $\sum_{r=0}^R \sum_{c=0}^C r_{rc} = 0.8c$ and $\sum_{r=0}^R \sum_{c=0}^C c_{rc} = 0.8r$. The key idea is that the interaction term in the model, ~~which indicates violations of additivity in standard ANOVA, indicates violations of multiplicative independence in exponentiated ANOVA~~, To estimate the magnitude of the ~~interaction terms~~ in exponentiation ANOVA. To summarize, the model of the cell tendencies is

$$X_{rc} = \exp(X_0 + X_r + X_c + X_{rc}) \quad \text{with the constraints} \\ r = 0 \text{ and } c = 0 \text{ and } r_{rc} = 0.8c \text{ and } c_{rc} = 0.8r \quad (22.2)$$

If the researcher is interested in violations of independence, then the interest is on the magnitudes of the r_{rc} interaction terms. The model is especially convenient for this purpose, because arbitrary interaction contrasts can be investigated to determine in more detail where the non-independence is arising.

22.1.3 The Poisson likelihood

The value of r_{rc} in Equation 22.2 is a cell tendency, not a predicted frequency. In particular, the value of r_{rc} can be any non-negative real value, but frequencies can only be integers. What we need to complete the model is a likelihood function that maps the parameter value r_{rc} to a probability of possible frequencies. The Poisson distribution is a

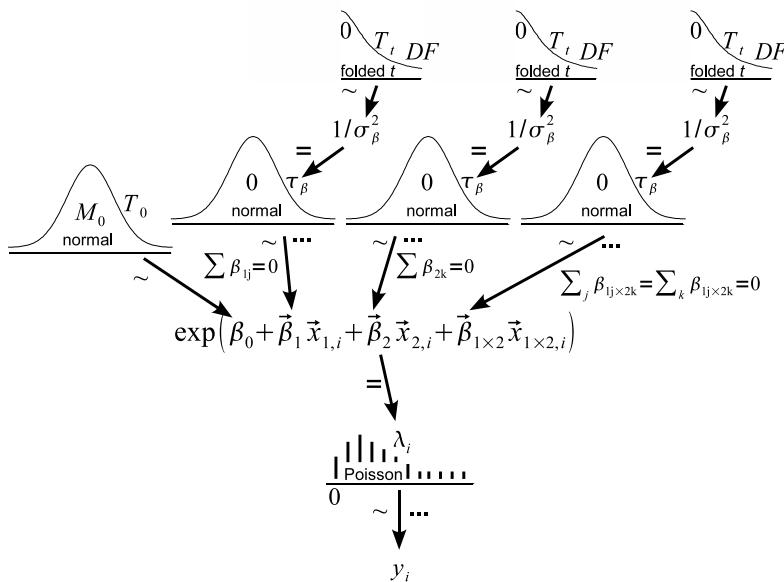


Figure 22.2: Hierarchical dependencies for Poisson-exponential model of a two-way frequency table. Compare with Figure 19.2, p. 424, for two-way ANOVA.

natural choice. The Poisson distribution is named after French mathematician Simon-Denis Poisson (1781–1840), and is defined as

$$p(y) = \frac{y}{y!} e^{-\lambda} \lambda^y \quad (22.3)$$

where y is a non-negative integer and λ is a non-negative real number. The mean of the Poisson distribution is λ . Importantly, the variance of the Poisson distribution is λ . Thus, in the Poisson distribution, the variance is λ . Examples of the Poisson distribution are shown in Figure 22.1. Notice that the distribution is discrete, having masses only at non-negative integer values. Notice that the visual central tendency of the distribution does indeed correspond with λ . And notice that as λ increases, the width of the distribution also increases. The examples in Figure 22.1 use non-integer values of λ to emphasize that it is not necessarily an integer, even though it is an integer.

The Poisson distribution is often used to model discrete occurrences in time (or across space) when the probability of occurrence is the same at any point in time (e.g., Sadiku & Toghi, 1999). For example, suppose that customers arrive at a retail store at an average rate of 35 people per hour. Then the Poisson distribution with $\lambda = 35$, is a model of the probability that any particular number of people will arrive in an hour. As another example, suppose that 11.2% of the students at the University of Delaware in the early 1970's had black hair and brown eyes, and suppose that an average of 600 students per term will fill out a survey. That means an average of 67.2 (11.2% of 600) students per term will indicate they have black hair and brown eyes. The Poisson distribution ($\lambda = 67.2$), gives the probability that any particular number of people will give that response.

We will use the Poisson distribution as the likelihood function for modeling the probability of the observed frequency $f(r; c)$, given the mean λ_c , from Equation 22.2. The idea is that each particular c combination has an underlying average rate of occurrence, λ_c . We collect data for a certain period of time, during which we happen to observe particular frequencies $f(r; c)$, of each combination. From the observed frequencies, we can estimate the parameters of the Poisson distribution.

Table 22.3: Some toy data for illustrating that interaction contrasts can be more sensitive to interaction than individual cells. (Posterior appears in Figure 22.4.)

Row Attribute	Column Attribute			
	C1	C2	C3	C4
R1	20	20	10	10
R2	20	20	10	10
R3	10	10	20	20
R4	10	10	20	20

underlying average rates.

22.1.4 The parameters and the hierarchical prior

The parameters of the Poisson exponential ANOVA model recall the parameters from standard ANOVA. The prior on those parameters is the same as standard ANOVA. Figure 22.2 shows the Poisson-exponential model, with its hierarchical prior. Above the linear core, the hierarchy is identical to the two-way ANOVA model in Figure 19.2, p. 424. The diagram retains the original notation from standard ANOVA (e.g., showing $\mu_{r;c}$ instead of $\mu_{r;c}$), but hopefully the correspondence of notation is easy to see.

Below the linear core in Figure 22.2, the diagram simply shows the exponentiation of the ANOVA summation to specify, and the use of the Poisson distribution to specify $p(y)$. This portion of the diagram replaces the normal likelihood of standard ANOVA by the prior on the normal distribution's precision. The Poisson, of course, has no separate precision parameter.

The model is easily implemented in R, BRugs, and BUGS, as in Section 22.4 (PoissonExponentialBRugs.R). The code is a straightforward modification of the program for two-factor ANOVA (that was listed in Section 19.3 ANOVAtwoWayBRugs.R). The modifications are described immediately before the listing.

22.2 Examples

Consider again the hair color and eye color data in Table 20.1. When the program in Section 22.4 (PoissonExponentialBRugs.R) is run with these data, the posterior has marginal histograms as shown in Figure 22.3. Inspection of the histograms of the individual cell interactions parameters reveals robustly non-zero interactions. For example, blue eyes and black hair has an interaction parameter that is credibly below zero, which means simply that the combination of blue eyes with black hair happens less frequently than would be expected if eye color and hair color were independent. As another example, brown eyes and black hair has an interaction parameter that is credibly above zero, which means simply that this combination happens more frequently than would be expected if eye color and hair color were independent. Exercise 22.2 has you consider other data and how to interpret main effects and interaction contrasts.

Another way to investigate interaction is by examining specific interaction contrasts. Interaction contrasts can be more sensitive and revealing than single-cell interaction terms. To illustrate this point, consider the toy data in Table 22.3. The upper-left and lower-right cells have higher frequencies than the upper-right and lower-left cells. When a Bayesian

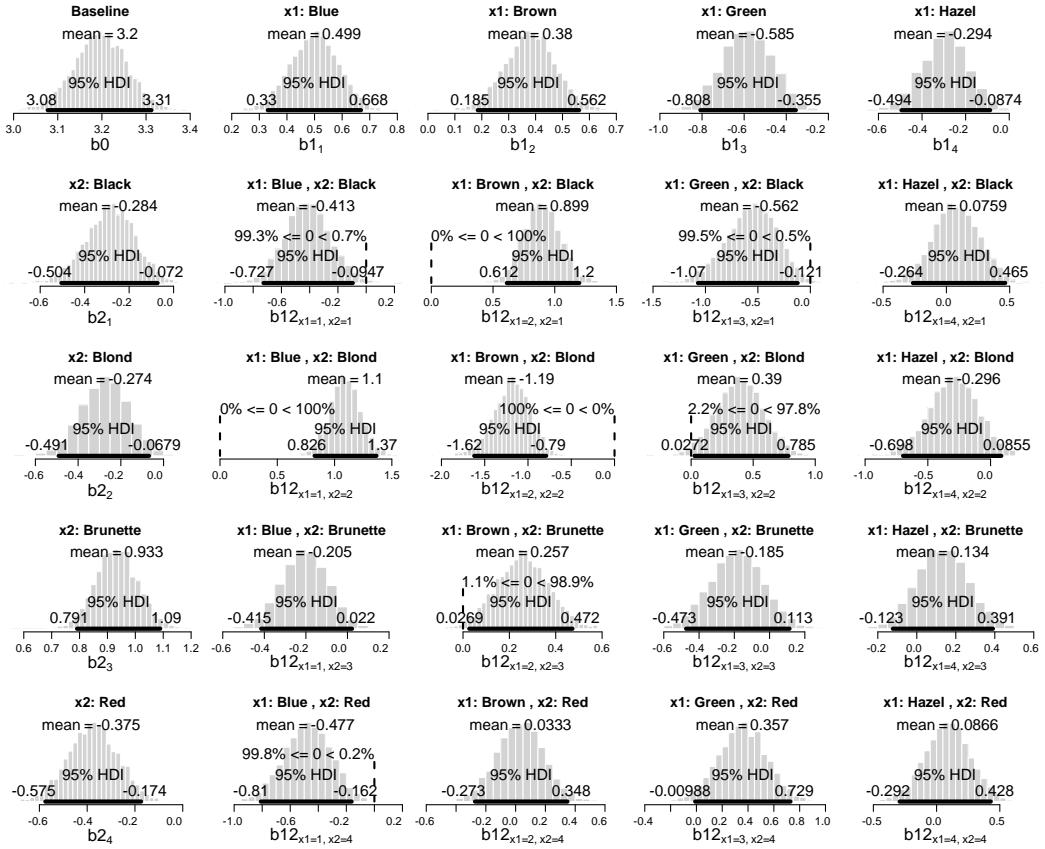


Figure 22.3: Posterior of Poisson exponential model applied to the data of Table 22.1.

analysis is conducted on these data, the resulting posteriors like Figure 22.4. Notice that not one of the single-cell interaction coefficients excludes zero from its HDI. The bottom histogram in Figure 22.4 shows the posterior distribution of an interaction contrast, namely $\text{ht}(1; 1; 1; 1) - \text{ht}(1; 1; 2; 2)$. This contrast takes the average of the interaction coefficients in the top-left and bottom-right cells, and subtracts the average of the interaction coefficients in the top-right and bottom-left cells. The histogram of the contrast magnitude clearly excludes zero from the credible values. Therefore, we would conclude that the attributes in Table 22.3 are not independent. Exercise 22.3(b) investigate the pattern of this toy example with different total sample sizes, and compare the Bayesian results with classical chi-square tests.

22.2.1 Credible intervals on cell probabilities

Although the posterior distributions of the beta coefficients are useful for making inferences about differences between cells, the magnitudes of the beta coefficients do not explicitly tell us about the credible cell proportions. To get those proportions, we need to use Equation 22.2 to compute the predicted cell mean frequency and then divide by their total to get the predicted cell proportions. Thus, at every step in the MCMC chain, we compute $P_{ij} = \exp(\alpha_0 + \beta_i + \gamma_j + \delta_{ij}) / \sum_{k=1}^K \sum_{l=1}^L \exp(\alpha_0 + \beta_k + \gamma_l + \delta_{kl})$, and then we divide, again at each step in the chain, by $\sum_{k=1}^K \sum_{l=1}^L \exp(\alpha_0 + \beta_k + \gamma_l + \delta_{kl})$. The resulting posterior distribution of normalized proportions shows the credible proportions of each combination of attributes. Figure 22.5 shows the resulting estimated

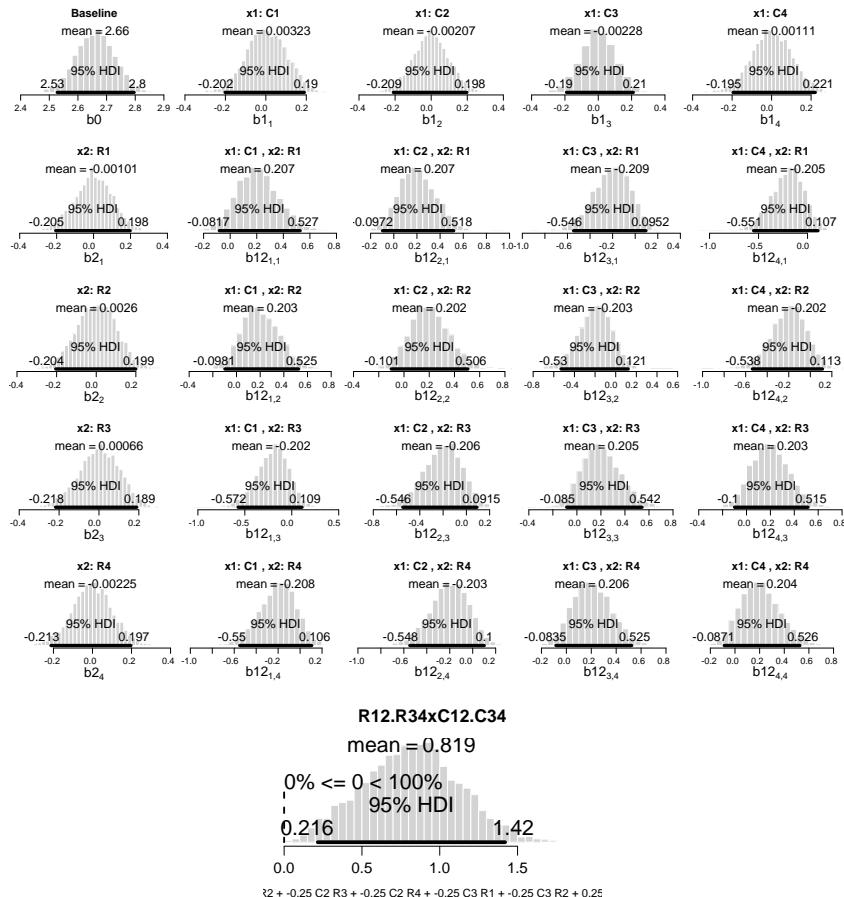


Figure 22.4: Posterior of Poisson exponential model applied to data of Table 22.3. The interaction contrast excludes zero, but no single main coefficient does.

cell probabilities for the halleye color data from Table 22.1.

22.3 Log linear models for contingency tables

This chapter only scratches the surface of methods for analyzing count data from nominal predictors. As mentioned earlier, these sorts of data are displayed as tables, and the counts in each cell are thought of as contingent upon the levels of the nominal predictor. Therefore the data are referred to as “contingency tables”. There can be more than two predictors, and models are generalized in the same way as ANOVA is generalized to more than two predictors. The formulation presented here has a simple link expressed as $P_{ij} = \exp(\alpha_0 + \alpha_r r_i + \alpha_c c_j + \alpha_{rc} r_i c_j)$, but this equation can also be written $\log(P_{ij}) = \alpha_0 + \alpha_r r_i + \alpha_c c_j + \alpha_{rc} r_i c_j$. This latter form lends the usual name for these models: log linear models for contingency tables. This is the terminology to use when you want to explore these models more deeply. Agresti and Hitchcock (2005) provide a brief review of Bayesian log linear models for contingency tables, but the method used in this chapter is not included because the hierarchical ANOVA model was only popularized later (Gelman, 2005, 2006). For a description of Bayesian inference regarding contingency tables with a Poisson ANOVA model, without the Gelman-style hyperprior, see Marin and Robert (2007, pp. 109–110), Congdon (2005, p. 134).

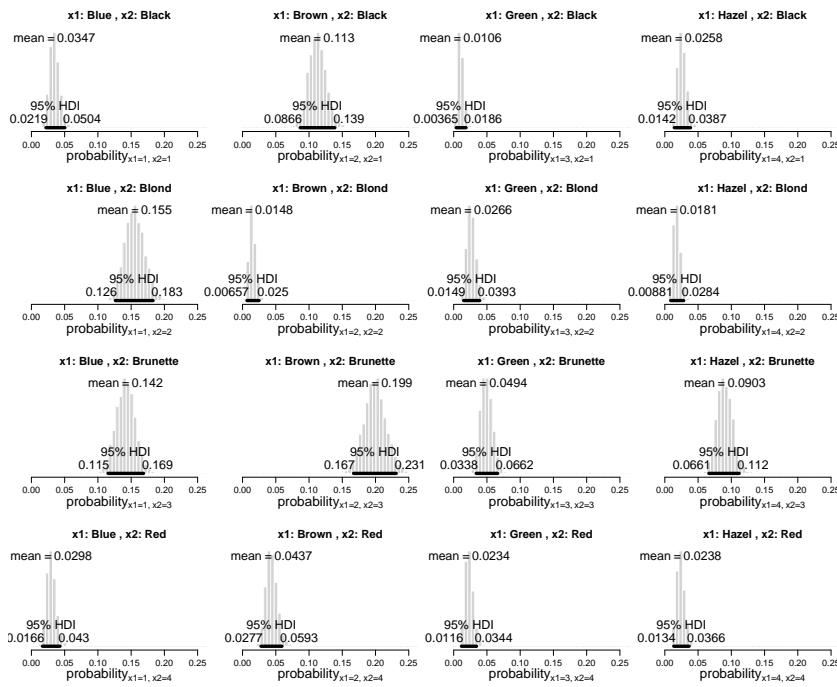


Figure 22.5: Posterior distribution on estimated cell abilities, for the data of Table 22.1, p. 490. See Section 22.2.1 for details.

and p. 202).

22.4 R code for Poisson exponential model

This program constitutes a small adaptation of the two-way ANOVA program listed in Section 19.3.1 (NOVAtwowayBRugs.R). The only notable changes are that (i) the normal likelihood is changed to a Poisson, (ii) all references to the *sigma* parameter of the normal likelihood are removed, (iii) the values are not standardized, because they must be non-negative integers, (iv) the initialization of the parameters is based on *log*, not *y*, because the parameters are exponentiated to map to

At the end of the program, some code for conducting a chi-square test is included. The purpose of this is merely for comparison of Bayesian results with NHST.

(PoissonExponentialBRugs.R)

```

1 graphics.off()
2 rm(list=ls(all=TRUE))
3 fnroot = "PoissonExponentialBRugs"
4 library(BRugs)          # Kruschke, J. K. (2010). Doing Bayesian data analysis:
5                         # A Tutorial with R and BUGS. Academic Press / Elsevier.
6 #-----#
7 # THE MODEL.
8
9 modelstring =
10 # BUGS model specification begins here...
11 model {
12   for ( i in 1:Ncells ) {
13     y[i] ~ dpois( lambda[i] )

```

```

14 lambda[i] <- exp( a0 + a1[x1[i]] + a2[x2[i]] + a1a2[x1[i],x2
15 ] )
16 #
17 a0 ~ dnorm(10,1.0E-6)
18 #
19 for ( j1 in 1:Nx1Lvl ) { a1[j1] ~ dnorm( 0.0 , a1tau ) }
20 a1tau <- 1 / pow( a1SD , 2 )
21 a1SD <- abs( a1SDunabs ) + .1
22 a1SDunabs ~ dt( 0 , 0.001 , 2 )
23 #
24 for ( j2 in 1:Nx2Lvl ) { a2[j2] ~ dnorm( 0.0 , a2tau ) }
25 a2tau <- 1 / pow( a2SD , 2 )
26 a2SD <- abs( a2SDunabs ) + .1
27 a2SDunabs ~ dt( 0 , 0.001 , 2 )
28 #
29 for ( j1 in 1:Nx1Lvl ) { for ( j2 in 1:Nx2Lvl ) {
30   a1a2[j1,j2] ~ dnorm( 0.0 , a1a2tau ) }
31 } }
32 a1a2tau <- 1 / pow( a1a2SD , 2 )
33 a1a2SD <- abs( a1a2SDunabs ) + .1
34 a1a2SDunabs ~ dt( 0 , 0.001 , 2 )
35 }
36 # ... end BUGS model specification
37 " # close quote for modelstring
38 # Write model to a file, and send to BUGS:
39 writeLines(modelstring,con="model.txt")
40 modelCheck( "model.txt" )

41 #
42 #-----#
43 # THE DATA.
44 # Specify data source:
45 dataSource = c( "HairEye" , "CrimeDrink" , "Toy" )[1]

46 # Load the data:
47 if ( dataSource == "HairEye" ) {
48   fnroot = paste( fnroot , dataSource , sep="" )
49   dataFrame = data.frame( # from Snee (1974)
50     Freq = c(68,119,26,7,20,84,17,94,15,54,14,10,5,29,14,
51     Eye = c("Brown","Brown","Brown","Brown","Blue","Blue"
52     Hair = c("Black","Brunette","Red","Blond","Black","Br
53 y = as.numeric(dataFrame$Freq)
54 x1 = as.numeric(dataFrame$Eye)
55 x1names = levels(dataFrame$Eye)
56 x2 = as.numeric(dataFrame$Hair)
57 x2names = levels(dataFrame$Hair)
58 Ncells = length(y)
59 Nx1Lvl = length(unique(x1))
60 Nx2Lvl = length(unique(x2))
61 x1contrastList = list( GREENvHAZEL = c(0,0,1,-1) )
62 x2contrastList = list( BLONDvRED = c(0,1,0,-1) )
63 x1x2contrastList = list( BLUE.BROWNxBLACK.BLOND
64                               = outer(c(-1,1,0,0),c(-1,1,0,0)) )
65 }
66 #

67 if ( dataSource == "CrimeDrink" ) {
68   fnroot = paste( fnroot , dataSource , sep="" )
69   dataFrame = data.frame( # from Kendall (1943) via Snee (1974 )
70     Freq = c(50,88,155,379,18,63,43,62,110,300,14,144) ,
71     Crime = c("Arson","Rape","Violence","Theft","Coining"
72                               , "Fraud","Arson","Rape","Violence","The

```

```

73     Drink = c("Drinker","Drinker","Drinker","Drinker","Drinker","Drinker",
74     y = as.numeric(dataFrame$Freq)
75     x1 = as.numeric(dataFrame$Crime)
76     x1names = levels(dataFrame$Crime)
77     x2 = as.numeric(dataFrame$Drink)
78     x2names = levels(dataFrame$Drink)
79     Ncells = length(y)
80     Nx1Lvl = length(unique(x1))
81     Nx2Lvl = length(unique(x2))
82     x1contrastList = list( FRAUDvOTHER = c(-1/5,-1/5,1,-1/5,      -1/5,-1/5) ,
83                           FRAUDvRAPE = c(0,0,1,-1,0,0) )
84     x2contrastList = list( DRINKERvNON = c(1,-1) )
85     x1x2contrastList = list( FRAUD.OTHERxDRINKER.NON
86                             = outer(c(-1/5,-1/5,1,-1/5,-1/5),c(-1,1)) ,
87                             FRAUD.RAPExDRINKER.NON
88                             = outer(c(0,0,1,-1,0,0),c(-1,1)) )
89   }
90
91 if ( dataSource == "Toy" ) {
92   dataMultiplier = 2 # Try 2 (chi-sq warns) , 6 (p>.05) , 7 (p<.05      ) , 10
93   fnroot = paste( fnroot , dataSource , dataMultiplier , sep="      " )
94   dataFrame = data.frame(
95     Freq = c( 2,2,1,1, 2,2,1,1, 1,1,2,2, 1,1,2,2 ) * dataMultipli
96     Col = c("C1","C2","C3","C4", "C1","C2","C3","C4", "C1",
97             Row = c("R1","R1","R1","R1", "R2","R2","R2","R2", "R3",
98     y = as.numeric(dataFrame$Freq)
99     x1 = as.numeric(dataFrame$Col)
100    x1names = levels(dataFrame$Col)
101    x2 = as.numeric(dataFrame$Row)
102    x2names = levels(dataFrame$Row)
103    Ncells = length(y)
104    Nx1Lvl = length(unique(x1))
105    Nx2Lvl = length(unique(x2))
106    x1contrastList = NULL
107    x2contrastList = NULL
108    x1x2contrastList = list( R12.R34xC12.C34 = outer(c(1,1,-
109   }                                         1,-1)/2,c(1,1,-1,-1)/2) )
110
111 # Specify the data in a form that is compatible with BRugs mode I, as a list:
112 dataList = list(
113   y = y ,
114   x1 = x1 ,
115   x2 = x2 ,
116   Ncells = Ncells ,
117   Nx1Lvl = Nx1Lvl ,
118   Nx2Lvl = Nx2Lvl
119 )
120 # Get the data into BRugs:
121 modelData( bugsData( dataList ) )
122
123 #-----
124 # INTIALIZE THE CHAINS.
125
126 nchain = 5
127 modelCompile( numChains = nchain )
128
129 if ( F ) {
130   modelGenInits() # often won't work for diffuse prior
131 } else {

```

```

132 # initialization based on data
133 theData = data.frame( y=log(y) , x1=factor(x1,labels=x1names) ,
134                               x2=factor(x2,labels=x2names) )
135 a0 = mean( theData$y )
136 a1 = aggregate( theData$y , list( theData$x1 ) , mean )[2] - a0      0
137 a2 = aggregate( theData$y , list( theData$x2 ) , mean )[2] - a0      0
138 linpred = as.vector( outer( a1 , a2 , "+" ) + a0 )
139 a1a2 = aggregate( theData$y , list(theData$x1,theData$x2      ), mean)[,3] - linpred
140 genInitList <- function() {
141   return(
142     list(
143       a0 = a0 ,
144       a1 = a1 ,
145       a2 = a2 ,
146       a1a2 = matrix( a1a2 , nrow=Nx1Lvl , ncol=Nx2Lvl ) ,
147       a1SDunabs = sd(a1) ,
148       a2SDunabs = sd(a2) ,
149       a1a2SDunabs = sd(a1a2)
150     )
151   )
152 }
153 for ( chainIdx in 1 : nchain ) {
154   modellnits( bugslnits( genInitList ) )
155 }
156 }

157 #
158 #-----#
159 # RUN THE CHAINS
160
161 # burn in
162 BurnInSteps = 1000
163 modelUpdate( BurnInSteps )
164 # actual samples
165 samplesSet( c( "a0" , "a1" , "a2" , "a1a2" , "a1SD" , "a2SD" , "a1a2SD" ) )
166 stepsPerChain = ceiling(5000/nchain)
167 thinStep = 500
168 modelUpdate( stepsPerChain , thin=thinStep )
169

170 #
171 # EXAMINE THE RESULTS
172
173 source("plotChains.R")
174 source("plotPost.R")
175
176 checkConvergence = F
177 if ( checkConvergence ) {
178   sumInfo = plotChains( "a0" , saveplots=F , filenameroot=fn      root )
179   sumInfo = plotChains( "a1" , saveplots=F , filenameroot=fn      root )
180   sumInfo = plotChains( "a2" , saveplots=F , filenameroot=fn      root )
181   sumInfo = plotChains( "a1a2" , saveplots=F , filenameroot=fnroot )
182   readline("Press any key to clear graphics and continue")
183   graphics.off()
184   sumInfo = plotChains( "a1SD" , saveplots=F , filenameroot=fnroot )
185   sumInfo = plotChains( "a2SD" , saveplots=F , filenameroot=fnroot )
186   sumInfo = plotChains( "a1a2SD" , saveplots=F , filenameroot=t=fnroot )
187   readline("Press any key to clear graphics and continue")
188   graphics.off()
189 }
190

```

```

191 # Extract and plot the SDs:
192 a1SDSample = samplesSample("a1SD")
193 a2SDSample = samplesSample("a2SD")
194 a1a2SDSample = samplesSample("a1a2SD")
195 windows(10,3)
196 layout( matrix(1:3,nrow=1) )
197 par( mar=c(3,1,2.5,0) , mgp=c(2,0.7,0) )
198 histInfo = plotPost( a1SDSample , xlab="a1SD" , main="a1 SD " , breaks=30 )
199 histInfo = plotPost( a2SDSample , xlab="a2SD" , main="a2 SD " , breaks=30 )
200 histInfo = plotPost( a1a2SDSample , xlab="a1a2SD" , main=" Interaction SD" ,
201                      breaks=30 )
202 dev.copy2eps(file=paste(fnroot,"SD.eps",sep=""))
203
204 # Extract a values:
205 a0Sample = samplesSample( "a0" )
206 chainLength = length(a0Sample)
207 a1Sample = array( 0 , dim=c( datalist$Nx1Lvl , chainLength ) )
208 for ( x1idx in 1:datalist$Nx1Lvl ) {
209   a1Sample[x1idx,] = samplesSample( paste("a1[",x1idx,"]",sep="") )
210 }
211 a2Sample = array( 0 , dim=c( datalist$Nx2Lvl , chainLength ) )
212 for ( x2idx in 1:datalist$Nx2Lvl ) {
213   a2Sample[x2idx,] = samplesSample( paste("a2[",x2idx,"]",sep="") )
214 }
215 a1a2Sample = array(0, dim=c( datalist$Nx1Lvl , datalist$N      x2Lvl , chainLength ) )
216 for ( x1idx in 1:datalist$Nx1Lvl ) {
217   for ( x2idx in 1:datalist$Nx2Lvl ) {
218     a1a2Sample[x1idx,x2idx,] = samplesSample( paste( "a1a2[      ,x1idx,",",x2idx,"]",sep="" ) )
219   }
220 }
221
222 # Convert to zero-centered b values:
223 m12Sample = array( 0 , dim=c( datalist$Nx1Lvl , datalist$Nx      2Lvl , chainLength ) )
224 for ( stepIdx in 1:chainLength ) {
225   m12Sample[,stepIdx] = ( a0Sample[stepIdx]
226                           + outer( a1Sample[,stepIdx] ,
227                                     a2Sample[,stepIdx] , "+" )
228                           + a1a2Sample[,stepIdx] )
229 }
230 b0Sample = apply( m12Sample , 3 , mean )
231 b1Sample = ( apply( m12Sample , c(1,3) , mean )
232               - matrix(rep( b0Sample ,Nx1Lvl),nrow=Nx1Lvl,byrow=T) )
233 b2Sample = ( apply( m12Sample , c(2,3) , mean )
234               - matrix(rep( b0Sample ,Nx2Lvl),nrow=Nx2Lvl,byrow=T) )
235 linpredSample = array(0,dim=c(datalist$Nx1Lvl,datalis      t$Nx2Lvl,chainLength))
236 for ( stepIdx in 1:chainLength ) {
237   linpredSample[,stepIdx] = ( b0Sample[stepIdx]
238                             + outer( b1Sample[,stepIdx] ,
239                                       b2Sample[,stepIdx] , "+" ) )
240 }
241 b1b2Sample = m12Sample - linpredSample
242
243 # Plot b values:
244 windows((datalist$Nx1Lvl+1)*2.75,(datalist$Nx2Lvl+1)      *2.25)
245 layoutMat = matrix( 0 , nrow=(datalist$Nx2Lvl+1) , ncol=(d      atalist$Nx1Lvl+1) )
246 layoutMat[1,1] = 1
247 layoutMat[1,2:(datalist$Nx1Lvl+1)] = 1:datalist$Nx1Lv      l + 1
248 layoutMat[2:(datalist$Nx2Lvl+1),1] = 1:datalist$Nx2Lv      l + (datalist$Nx1Lvl + 1)

```

```

250 layoutMat[2:(datalist$Nx2Lvl+1),2:(datalist$Nx1Lvl+      1)] = matrix(
251   1:(datalist$Nx1Lvl*datalist$Nx2Lvl) + (datalist$Nx2Lvl      l+datalist$Nx1Lvl+1) ,
252   ncol=datalist$Nx1Lvl , byrow=T )
253 layout( layoutMat )
254 par( mar=c(4,0.5,2.5,0.5) , mgp=c(2,0.7,0) )
255 histinfo = plotPost( b0Sample , xlab=expression(beta * 0) ,      main="Baseline" ,
256                      breaks=30 )
257 for ( x1idx in 1:datalist$Nx1Lvl ) {
258   histinfo = plotPost( b1Sample[x1idx,] , xlab=bquote(beta      *1[(x1idx)]) ,
259                         main=paste("x1:",x1names[x1idx]) , breaks=30 )
260 }
261 for ( x2idx in 1:datalist$Nx2Lvl ) {
262   histinfo = plotPost( b2Sample[x2idx,] , xlab=bquote(beta      *2[(x2idx)]) ,
263                         main=paste("x2:",x2names[x2idx]) , breaks=30 )
264 }
265 for ( x2idx in 1:datalist$Nx2Lvl ) {
266   for ( x1idx in 1:datalist$Nx1Lvl ) {
267     hdiLim = HDIofMCMC(b1b2Sample[x1idx,x2idx,])
268     if ( hdiLim[1]>0 | hdiLim[2]<0 ) { compVal=0 } else { compVal=      NULL }
269     histinfo = plotPost( b1b2Sample[x1idx,x2idx,] , breaks=3      0 , compVal=compVal ,
270                           xlab=bquote(beta*12[list(x1==.(x1idx),x2==.(x2idx))      ]) ,
271                           main=paste("x1:",x1names[x1idx]," x2:",x2names[x2id      x]) )
272   }
273 }
274 dev.copy2eps(file=paste(fnroot,"b.eps",sep=""))
275
276 # Display contrast analyses
277 nContrasts = length( x1contrastList )
278 if ( nContrasts > 0 ) {
279   nPlotPerRow = 5
280   nPlotRow = ceiling(nContrasts/nPlotPerRow)
281   nPlotCol = ceiling(nContrasts/nPlotRow)
282   windows(3.75*nPlotCol,2.5*nPlotRow)
283   layout( matrix(1:(nPlotRow*nPlotCol),nrow=nPlotRow,n      col=nPlotCol,byrow=T) )
284   par( mar=c(4,0.5,2.5,0.5) , mgp=c(2,0.7,0) )
285   for ( cldx in 1:nContrasts ) {
286     contrast = matrix( x1contrastList[[cldx]],nrow=1 ) # make      it a row matrix
287     incldx = contrast!=0
288     histInfo = plotPost( contrast %*% b1Sample , compVal=0 , bre      aks=30 ,
289                           xlab=paste( round(contrast[incldx],2) , x1names[incldx      ] ,
290                           c(rep("+",sum(incldx)-1),"") , collapse=" " ) ,
291                           cex.lab = 1.0 ,
292                           main=paste( "X1 Contrast:", names(x1contrastList)[cldx      ] ) )
293   }
294   dev.copy2eps(file=paste(fnroot,"x1Contrasts.eps",se      p=""))
295 }
296 #
297 nContrasts = length( x2contrastList )
298 if ( nContrasts > 0 ) {
299   nPlotPerRow = 5
300   nPlotRow = ceiling(nContrasts/nPlotPerRow)
301   nPlotCol = ceiling(nContrasts/nPlotRow)
302   windows(3.75*nPlotCol,2.5*nPlotRow)
303   layout( matrix(1:(nPlotRow*nPlotCol),nrow=nPlotRow,n      col=nPlotCol,byrow=T) )
304   par( mar=c(4,0.5,2.5,0.5) , mgp=c(2,0.7,0) )
305   for ( cldx in 1:nContrasts ) {
306     contrast = matrix( x2contrastList[[cldx]],nrow=1 ) # make      it a row matrix
307     incldx = contrast!=0
308     histInfo = plotPost( contrast %*% b2Sample , compVal=0 , bre      aks=30 ,

```

```

309         xlab=paste( round(contrast[inclIdx],2) , x2names[inclIdx]      ] ,
310                     c(rep("+",sum(inclIdx)-1),"") , collapse=" " ) ,
311                     cex.lab = 1.0 ,
312                     main=paste( "X2 Contrast:", names(x2contrastList)[cldx]      ] ) )
313     }
314   dev.copy2eps(file=paste(fnroot,"x2Contrasts.eps",se           p=""))
315 }
316 #
317 nContrasts = length( x1x2contrastList )
318 if ( nContrasts > 0 ) {
319   nPlotPerRow = 5
320   nPlotRow = ceiling(nContrasts/nPlotPerRow)
321   nPlotCol = ceiling(nContrasts/nPlotRow)
322   windows(3.75*nPlotCol,2.5*nPlotRow)
323   layout( matrix(1:(nPlotRow*nPlotCol),nrow=nPlotRow,n       col=nPlotCol,byrow=T) )
324   par( mar=c(4,0.5,2.5,0.5) , mgp=c(2,0.7,0) )
325   for ( cldx in 1:nContrasts ) {
326     contrast = x1x2contrastList[[cldx]]
327     contrastArr = array( rep(contrast,chainLength) ,
328                           dim=c(NROW(contrast),NCOL(contrast),chainLength) )
329     contrastLab = ""
330     for ( x1idx in 1:Nx1Lvl ) {
331       for ( x2idx in 1:Nx2Lvl ) {
332         if ( contrast[x1idx,x2idx] != 0 ) {
333           contrastLab = paste( contrastLab , "+" ,
334                               signif(contrast[x1idx,x2idx],2) ,
335                               x1names[x1idx] , x2names[x2idx] )
336         }
337       }
338     }
339     histInfo = plotPost( apply( contrastArr * b1b2Sample , 3 , su      m ) ,
340                         compVal=0 , breaks=30 , xlab=contrastLab , cex.lab = 0.75 ,
341                         main=paste( names(x1x2contrastList)[cldx] ) )
342   }
343   dev.copy2eps(file=paste(fnroot,"x1x2Contrasts.eps",           sep=""))
344 }
345 #
346 # Compute credible cell probability at each step in the MCMC chain
347 lambda12Sample = 0 * b1b2Sample
348 for ( chainIdx in 1:chainLength ) {
349   lambda12Sample[,chainIdx] = exp(
350     b0Sample[chainIdx]
351     + outer( b1Sample[,chainIdx] , b2Sample[,chainIdx] , "+" )
352     + b1b2Sample[,chainIdx] )
353 }
354 cellp = 0 * lambda12Sample
355 for ( chainIdx in 1:chainLength ) {
356   cellp[,chainIdx] = ( lambda12Sample[,chainIdx]
357                         / sum( lambda12Sample[,chainIdx] ) )
358 }
359 #
360 # Display credible cell probabilities
361 windows((datalist$Nx1Lvl)*2.75,(datalist$Nx2Lvl)*2.      25)
362 layoutMat = matrix( 1:(datalist$Nx2Lvl*datalist$Nx1Lvl)      ) ,
363                     nrow=(datalist$Nx2Lvl) , ncol=(datalist$Nx1Lvl) , byrow      =T )
364 layout( layoutMat )
365 par( mar=c(4,1.5,2.5,0.5) , mgp=c(2,0.7,0) )
366 maxp = max( cellp )
367 for ( x2idx in 1:datalist$Nx2Lvl ) {
368   for ( x1idx in 1:datalist$Nx1Lvl ) {

```

```

368 histinfo = plotPost( cellp[x1idx,x2idx,] ,
369   breaks=seq(0,maxp,length=51) , xlim=c(0,maxp) ,
370   xlab=bquote(probability[list(x1==.(x1idx),x2==.(x2i
371   main=paste("x1:",x1names[x1idx]," x2:",x2names[x2id
372   HDItextPlace=0.95 )
373 }
374 }
375 dev.copy2eps(file=paste(fnroot,"CellP.eps",sep=""))
376
377
378 #=====
379 # Conduct NHST Pearson chi-square test of independence.
380
381 # Convert DataFrame to frequency table:
382 obsFreq = matrix( 0 , nrow=Nx1Lvl , ncol=Nx2Lvl )
383 for ( x1idx in 1:Nx1Lvl ) {
384   for ( x2idx in 1:Nx2Lvl ) {
385     obsFreq[x1idx,x2idx] = y[ DataFrame[,2]==x1names[x1idx]
386                               & DataFrame[,3]==x2names[x2idx] ]
387   }
388 }
389 obsFreq = t(obsFreq) # merely to match orientation of histogram display
390 chisqtest = chisq.test( obsFreq )
391 print( "obs :" )
392 print( chisqtest$observed )
393 print( "( obs - exp )^2 / exp :" )
394 print( ( chisqtest$observed - chisqtest$expected )^2 / chisqtest$expected )
395 print( chisqtest )
396
397 #=====

```

22.5 Exercises

Exercise 22.1.[Purpose: Sample size and precision of estimate.] Consider the data regarding hair color and eye color in Table 22.1, which is available in the code of Section 22.4 (PoissonExponentialBrugs.R).

(A) Divide the original data frequencies by 2 and take the log of the result. Thus, the modified data are ceiling(c(68,...,16) / 2). This modification (nearly) preserves the relative proportions in each cell but halves the sample size. Run the program with these modified data, show the histograms, and compute the width of the 95% HDI's for the contrasts.

(B) Multiply the original data vector by 5, so that all the cell frequencies are five times larger than the original. This modification preserves the proportions in each cell but quintuples the sample size. Run the program again. Show histograms. Compute the width of the 95% HDI's for the contrasts. How do the precision of the contrasts differ for the reduced and the enlarged data?

Exercise 22.2.[Purpose: Explore “main effects” and interactions in a contingency table.] The data section of the program in Section 22.4 (PoissonExponentialBrugs.R) includes data regarding the type of crime committed by convicted criminals and whether or not the criminal is a regular drinker of alcoholic beverages (Snee, 1974).

(A) Run the program with these data selected. Is there a credible independence of the attributes? Describe the interaction in terms of the levels of the attributes (i.e.,

type of crime and drinker or non-drinker).

(B) Is there a (“simple”) main effect of fraud versus rape? That is, marginalizing across drinking, is there a credible difference in the proportion of criminals who committed fraud or committed rape? Show the results of an appropriate contrast.

(C) Is the difference from the previous part the same among drinkers and non-drinkers? In other words, is the effect of criminal type the same at all levels of drinking? Show the results of an appropriate interaction contrast.

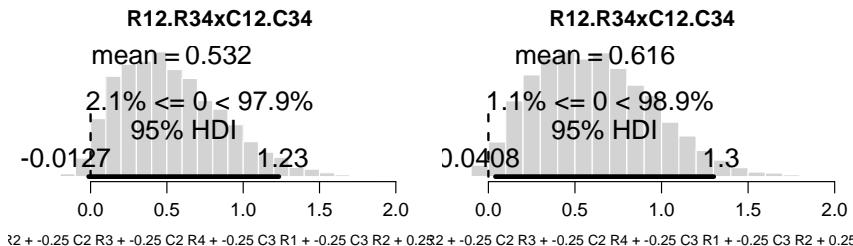


Figure 22.6: For Exercise 22.3. Left histogram shows posterior interaction contrast for the “toy” data when `dataMultiplier = 6`. A chi-square test (not shown) indicates that $p > .05$. Right histogram shows posterior interaction contrast for the “toy” data when `dataMultiplier = 7`. A chi-square test (not shown) indicates that $p < .05$.

Exercise 22.3. [Purpose: Compare Bayesian analysis with chi-square test of independence.] This exercise assumes that you have some familiarity with `fitdistr` chi-square tests of independence. In the code of Section 22.4 (`PoissonExponentialBrugs.R`), the data section includes some “toy” data, which we can manipulate and compare results of Bayesian and chi-square analyses.

(A) Set the `dataMultiplier` to 6, and run the program. Include the histograms of the parameters and of the interaction contrast. Does the 95% HDI contrast include zero? The end of the program runs a chi-square test. What is the value from the test? Hint: See Figure 22.6.

(B) Set the `dataMultiplier` to 7, and run the program. Include the histograms of the parameters and of the interaction contrast. Does the 95% HDI contrast include zero? The end of the program runs a chi-square test. What is the value from the test? Hint: See Figure 22.6.

(C) Typical chi-square tests rely on approximating the sampling distribution of discrete Pearson chi-square values (i.e., $\sum_c (f_{rc} - \hat{f}_{rc})^2 / f_{rc}$) with the continuous chi-square distribution (which derives from the sum of standardized normal variables). When the expected frequencies, \hat{f}_{rc} , are too small, then the approximation is not very good, and the estimated p value may be wrong. A usual heuristic for declaring the chi-square test to be suspect is when (at least 10% of the) expected frequencies are less than 5. This is why computer packages will issue warnings when expected values are too small. Bayesian analysis has no such problems. Set the `dataMultiplier` to 2, and run the program. Does the Bayesian analysis complain or do anything wrong? (No.) The end of the program runs a chi-square test. Is there a warning message? (Yes; report what it is.)

Chapter 23

Tools in the Trunk

Contents

23.1	Reporting a Bayesian analysis	508
23.1.1	Essential points	508
23.1.2	Optional points	509
23.1.3	Helpful points	509
23.2	MCMC burn-in and thinning	510
23.3	Functions for approximating highest density intervals	513
23.3.1	R code for computing HDI of a grid approximation	513
23.3.2	R code for computing HDI of a MCMC sample	513
23.3.3	R code for computing HDI of a function	155
23.4	Reparameterization of probability distributions	516
23.4.1	Examples	516
23.4.2	Reparameterization of two parameters	517

She changes her hair, and he changes his style,
She paints on her face, and he wears a fake smile,
She shrink wraps her head, and he stretches the truth;
But they'll always be stuck with their done wasted youth.

This chapter includes a few important topics that carry over many different models throughout the book. The first topic is how to report a Bayesian analysis in a scientific journal. The second topic is the details behind computing an HDI.

The third issue is something that has not come up explicitly often in the book, but lurks in the shadows all the time. This is the topic of reparameterization. For example, the beta distribution has shape parameters α and β , but we regularly transformed them into parameters μ and σ . The complication of reparameterization is that a probability distribution on a parameter is not the same on the transformed parameter. In particular, if we specify a prior on a parameter, but then transform the parameter, the prior on the transformed parameter is different. (The quatrain at the beginning of the chapter has reparameterization in mind.)

Kruschke, J. K. (2010) Doing Bayesian Data Analysis: A Tutorial with R and BUGS. Academic Press
Elsevier. Copyright 2010 by John K. Kruschke. Draft of May 11, 2010. Please do not distribute this
preliminary draft. If you report Bayesian analyses based on this book, please do cite it! “

23.1 Reporting a Bayesian analysis

Bayesian data analyses are not yet standard procedure in fields of research, and no conventional format for reporting them has been established. Therefore the researcher who reports a Bayesian analysis must be sensitive to the background knowledge of his or her specific audience, and frame the description accordingly. I once assigned an exercise to students in which they had to write up the results of a Bayesian analysis as it would appear in a research journal. Because I am a psychologist, a student earnestly asked me, "Do we have to write it as if it were for a psychology journal or for a science journal?" After swallowing my feeling of injury at the implication that psychology is not science, I politely asked the student to clarify the question. The student said, "Well, a psychology journal you have to explain what you did, but for a science journal there has to puzzle it out."

23.1.1 Essential points

When reporting a Bayesian analysis, the writer must address the following essential points:

Motivate the use of Bayesian (non-NHST) analysis. Many audiences, including journal editors and reviewers, are used to NHST and unfamiliar with Bayesian methods. You may motivate your use of Bayesian data analysis on several grounds, depending in part on the particular application. For example, Bayesian models are designed to be appropriate to the data structure, without having to make approximation assumptions typical in NHST. The inferences from a Bayesian analysis are richer and more informative than NHST. And, of course, there is no reliance on p-values.

Clearly describe the model and its parameters. Because the posterior distribution is a distribution over parameter values, the parameters must be clearly defined. This task of describing the model can be arduous for complex hierarchical models, but it is necessary and crucial if your analysis is to mean anything to your audience. Because the model refers to the data and their structure, the model description requires that the structure of the data has already been explained.

Clearly describe and justify the prior. It is crucial to convince your audience that your prior is appropriate and does not predetermine the results. There is no escape from this requirement, even when using "objective" or "uninformative" priors, because there is no unique choice of such priors and even they can have major consequences for model comparison. The prior should be amenable to a skeptical audience. The prior should be at least mildly informed to match the scale of the data being modeled. If there is copious previous research using very similar models, it should not be ignored. Optionally, as mentioned again below, it may be useful to try different priors and report the robustness of the posterior.

Mention the MCMC details, especially evidence that the chains were converged (plenty of burn-in and no orphaned chains) and not clumpy (auto-correlation from sufficient thinning). Also state how many points there are in the final MCMC sample. Usually this section of the report can be brief, if your audience believes that you know what you're doing and therefore took adequate care to generate a trustworthy and representative sample from the posterior. See Section 23.2 for a reminder of how to do this.

Interpret the posterior. Many models have dozens or even hundreds of parameters, and therefore it is impossible to summarize all of them. The choice of which parameters or contrasts to report is driven by domain-specific theory and by the results themselves. You want to report the parameters and contrasts that are theoretically meaningful, and those that showed effects driven by data, whether expected or not. Reporting of HDI's can be done in text alone, to save space. (Diagrams of posteriors are useful for explanation in a textbook, but may be unnecessary in a concise report.) Include scatterplots of correlated important parameters if the scatter is not bivariate Gaussian. Be sure to describe effects of shrinkage if appropriate. If your model includes interactions of predictors, be careful how you interpret lower-order effects. Finally, if you are using a ROPE for posterior interpretation, justify its limits.

23.1.2 Optional points

The following points are not necessarily crucial to address every report, but certainly should be considered. Whether or not to include these points depends on the particulars of the application domain, the points the reporter wants to make to the audience to which the report is being addressed.

Robustness of the posterior for different priors. When there is real or imagined contention about the prior, it can be most convincing simply to conduct the analysis with different priors and demonstrate that the essential conclusion of the posterior do not change. Which priors should be used? Those that are appropriate for your audience, such as reviewers and editors of the journal to which the report is submitted. This uncertainty in the choice of priors may seem unappealing, but it accurately reflects the way science gets done, by incorporating previous research and addressing currently active researchers.

Posterior predictive check. By generating simulated data from the credible parameter values of the model, and examining the qualities of the simulated data, the veracity of the model can be further bolstered, if the simulated data resemble the actual data. On the other hand, if the simulated data are discrepant from the actual data in systematic and interpretable ways, then the posterior predictive check can inspire new research and new models.

Power analysis. For example, if there is only a weak effect in your results, what sample size would be needed to achieve some desired precise estimate? If you found a credible difference, what was the retrospective power of your experiment and what is its replication power? This sort of information can be useful not only for the researcher's own planning, but it can also be useful to the audience of the report to anticipate potential follow-up research and to assess the robustness of the currently reported results.

23.1.3 Helpful points

Finally, to help science be cumulative, make your results available on the web:

Post the raw data. There are two benefits of posting the original data. One benefit is that subsequent researchers can analyze the data with different models. New insights can be gained by alternative modeling interpretations. The gravity of the original

research is enhanced. A second benefit is that if an exact replication is conducted, the original data set can be concatenated with the new data set, to enhance sensitivity of the new data.

Post the MCMC sample of the posterior there are two benefits of making the posterior publicly available. One is that other researchers explore the posterior for effects and comparisons that were not in the report. Complex models have many parameters, and no single report can cover every possible aspect on the posterior distribution. The longevity and impact of the research is thereby enhanced. A second benefit is that if subsequent researchers do follow-up research with a similar design and model, then the posted posterior can inform the prior of subsequent analysis. Because the full posterior automatically incorporates the covariations between all the parameters, the full posterior can be more useful than summaries of marginal distributions in a report.

23.2 MCMC burn-in and thinning

When using a Markov chain to generate a Monte Carlo sample from a distribution, we want to be sure that the resulting chain is a truly representative sample from the distribution. There are several ways in which the chain might fail to be representative.

One problem is that the chains might start far from the true(s) of the distribution, and therefore the early steps in the chain might be unrepresentative of the distribution. The main way to mitigate this problem is to start the chains ~~randomly~~ lightly, instead of randomly, if you can. If you know a point that is likely to be in the midst ~~of~~ distribution, try starting the chains there. Many of the programs in this book use this. For example, in the multiple linear regression program of Section 17.5 (`MultipleLinearRegressionBrugs.R`), the chains were started at the maximum likelihood estimate because it was assumed that the posterior will be dominated by the data, not by the prior. Although a good starting point reduces the time it takes for the chains to find the bulk of the distribution, the early steps are still not necessarily representative because we still let the chains run a while to dilute the influence of the initial value; see Figure 7.2, p21. This initial set of steps is called the burn-in period.

Another problem is that the chains might not change values from step to step, and therefore take a long time to generate values from the breadth of the distribution. When consecutive steps have similar values, the chain ~~is~~ state highly autocorrelated. The problem is not merely taking a long time to generate the range of the distribution. The main problem is that an autocorrelated chain is “clumpy” over-represents some values and under-represents other values. To primary mitigate this problem is to thin the chain: Instead of using every step in the chain, we only every m^{th} step, where m could be 50 or 500 or larger, depending on the model and data.

Autocorrelation is measured simply as the correlation of ~~chain~~ values with the values L steps behind (or ahead), ~~which~~ called the “lag”. Expressed in terms of R code, if the chain is the vector `v` and it has `n` components, then the autocorrelation at lag $ACF(L) = \text{cor}(v[1:(W-L)], v[(L+1):W])$. At $L=0$, the $ACF(0)$ is 1.0, of course. If consecutive steps in the chain are very similar, then $ACF(1)$ will be close to 1.0. If consecutive steps in the chain are uncorrelated, then $ACF(1)$ will be close to zero.

Another problem is that it is possible for chains to get stuck in nonrepresentative regions of parameter space if they are initialized poorly or take a jump. As a check that the

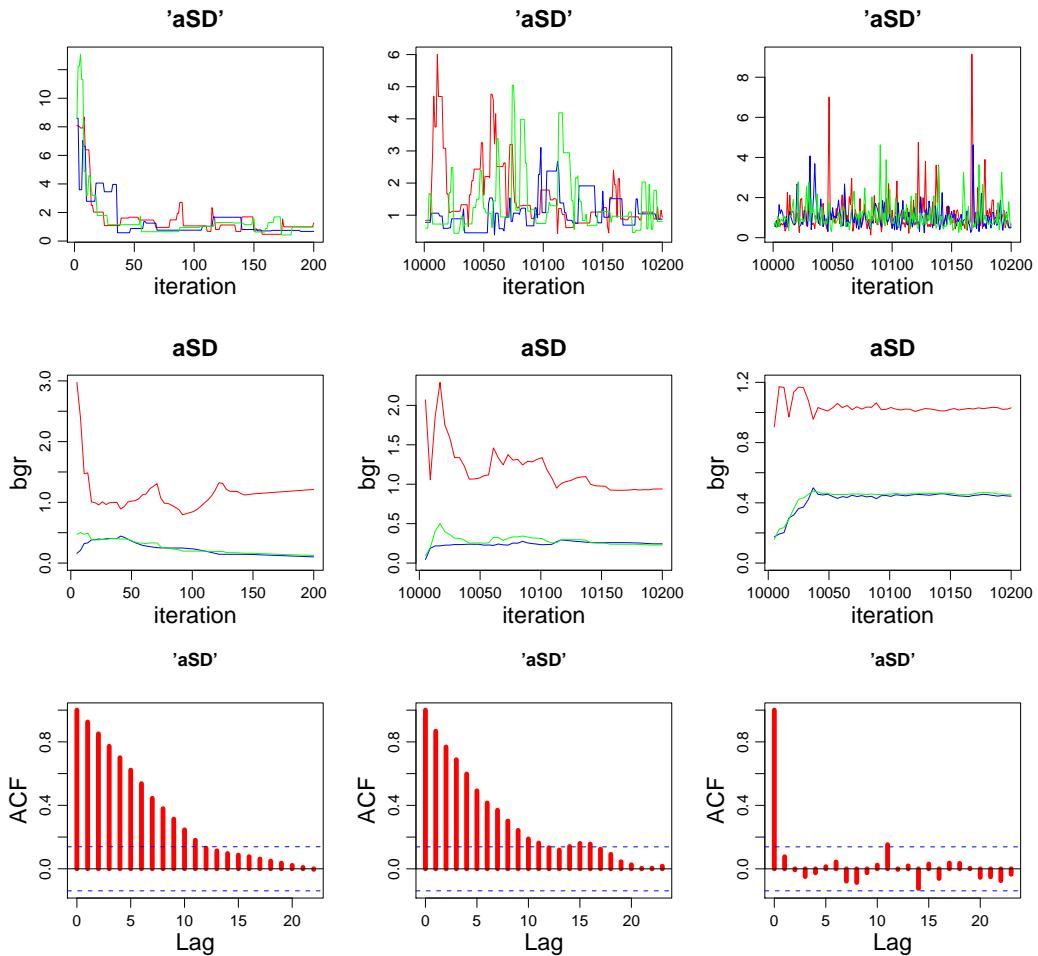


Figure 23.1: Effects of burn-in and thinning on MCMC chains. Left column: Zero burn-in steps and no thinning. Middle column: 10,000 burn-in steps and no thinning. Right column: 10,000 burn-in steps and thinning to 1 every 500 steps.

chains are actually exploring the bulk of the distribution and are not stuck in some unrepresentative zone, we run several chains at once and examine whether they are well mixed, without any outlying chain.

One measure of mixing considers the variability between chains relative to the variability within chains. If the chains are well mixed, then this ratio should be about 1, i.e., there should be no more variability between chains than within chains. But if the chains are not well mixed, and one or more chains is lingering far from the rest, then the variance between chains will be larger than the variance within chains. This ratio is merely a F ratio for the chain values over a specified window of steps. In BUGS it is called the bgr statistic, after Brooks, Gelman and Rubin (see Brooks & Gelman, 1998).

To illustrate these ideas, consider the parameter from oneway ANOVA, near the top of the hierarchical diagram of Figure 18.1, p. 403. The standard deviation expresses the estimated variation across groups. In the BUGS program of Section 18.4.1 (ANOVAonewayBRugs), this parameter is denoted σ . The top-left panel of Figure 23.1 shows the initial steps of three chains when they are intentionally started far from truly representa-

tive values. You can see that it takes about 50 steps (labeled “iterations” in the graph) for the chains to descend to representative values. The top panel of Figure 23.1 shows the chains after 10,000 burn-in steps (which is overkill for this particular application). The chains appear to be converged, without systematically ~~earring~~ or decreasing.

The bottom row of panels in Figure 23.1 shows the autocorrelation function, ACF. The left and middle columns show results when there is no thinning of the chains, and you can see that the ACF remains high for lags up to 15 or 20 steps. The correlation is visible in the chains themselves as sustained plateaus during which values from step to step barely change. The right column shows results when the chains were thinned, keeping only 1 step in every 500. You can see that after even one of the thinnings, $\text{ACF}(L=1)$ is nearly zero.

The middle row of Figure 23.1 shows the bgr statistic, which measures how well mixed the chains are, over a limited length of the chains. As mentioned above, the bgr should be around 1.0. The plots show the bgr as a curve (plotted as red by BUGS) hovering near a value of 1.0. The two other curves, lower in graphs (and colored green and blue by BUGS), show the between-chain and within-chain variance. In this particular application, the chains converge quickly and well.

The graphs in Figure 23.1 were produced by the following R program, which can be called after running any BUGS model. The argument `nodename` is a string that is name of any variable monitored by BUGS. (`plotChains.R`)

```

1 plotChains = function( nodename , saveplots=F , filenamero      ot="DeleteMe" ) {
2   summarytable = samplesStats(nodename)
3   show( summarytable )
4   nCompon = NROW(summarytable)
5   nPlotPerRow = 5
6   nPlotRow = ceiling(nCompon/nPlotPerRow)
7   nPlotCol = ceiling(nCompon/nPlotRow)
8   windows(3.75*nPlotCol,3.5*nPlotRow)
9   par( mar=c(4,4,3,1) , mgp=c(2,0.7,0) )
10  samplesHistory( nodename , ask=F , mfrow=c(nPlotRow,nPlotCol) ,
11                  cex.lab=1.5 , cex.main=1.5 )
12  if ( saveplots ) {
13    dev.copy2eps( file=paste( filenamero      me) ,
14                  "history.eps" , sep="" ) }
15  windows(3.75*nPlotCol,3.5*nPlotRow)
16  par( mar=c(4,4,3,1) , mgp=c(2,0.7,0) )
17  samplesAutoC( nodename , chain=1 , ask=F , mfrow=c(nPlotRo w,nPlotCol) ,
18                  cex.lab=1.5 , cex.main=1.5 )
19  if ( saveplots ) {
20    dev.copy2eps( file=paste( filenamero      me) ,
21                  "autocorr.eps" , sep="" ) }
22  windows(3.75*nPlotCol,3.5*nPlotRow)
23  par( mar=c(4,4,3,1) , mgp=c(2,0.7,0) )
24  samplesBgr( nodename , ask=F , mfrow=c(nPlotRow,nPlotCol) ,
25                  cex.lab=1.5 , cex.main=1.5 )
26  if ( saveplots ) {
27    dev.copy2eps( file=paste( filenamero      me) ,
28                  "bgr.eps" , sep="" ) }
29  return( summarytable )
30 }
```

23.3 Functions for approximating highest density intervals

HDI's have been used routinely throughout the book to describe distributions. This section provides details regarding how the HDI's are computed. The algorithm for computing an HDI on a grid approximation applied to any dimensionality or any shape distribution. The algorithms for computing an HDI of an MCMC sample or for a function apply only to single parameters with unimodal distributions.

23.3.1 R code for computing HDI of a grid approximation

This function was first used in the R code of Section 6.7 `BernGrid.R`, and again in Section 8.8.1 `BernTwoGrid.R`.

We can imagine the grid approximation of a distribution as a landscape of poles sticking up from each point on the parameter grid, with the height of each pole indicating the probability mass at that discrete point. We can imagine the highest density region by visualizing a rising tide: We gradually flood the landscape, monitoring the total mass of the poles that protrude above water, stopping the flood when 95% (say) of the poles remains protruding. The waterline at that moment defines the highest density region.

The program, listed below, finds the approximate highest density region in a somewhat analogous way. It uses one extra trick at the beginning, however: first sorts all the poles in order of height, from tallest to shortest. The idea is to move down the sorted queue of poles until the cumulative probability has just barely exceeded 95% (or whatever). The resulting height is the “waterline” that defines all points inside the highest density. See the comments in the top of the code for details of how to use the function.

(`HDIofGrid.R`)

```

1 HDIofGrid = function( probMassVec , credMass=0.95 ) {
2   # Arguments:
3   #   probMassVec is a vector of probability masses at each grid point.
4   #   credMass is the desired mass of the HDI region.
5   # Return value:
6   #   A list with components:
7   #     indices is a vector of indices that are in the HDI
8   #     mass is the total mass of the included indices
9   #     height is the smallest component probability mass in the HDI
10  # Example of use: For determining HDI of a beta(30,12) distribution
11  # approximated on a grid:
12  # > probDensityVec = dbeta( seq(0,1,length=201) , 30 , 12 )
13  # > probMassVec = probDensityVec / sum( probDensityVec )
14  # > HDIinfo = HDIofGrid( probMassVec )
15  # > show( HDIinfo )
16  sortedProbMass = sort( probMassVec , decreasing=T )
17  HDIheightIdx = min( which( cumsum( sortedProbMass ) >= credMass ) )
18  HDIheight = sortedProbMass[ HDIheightIdx ]
19  HDImass = sum( probMassVec[ probMassVec >= HDIheight ] )
20  return( list( indices = which( probMassVec >= HDIheight ) ,
21             mass = HDImass , height = HDIheight ) )
22 }
```

23.3.2 R code for computing HDI of a MCMC sample

The algorithms for computing the HDI for an MCMC sample or for a function rely on a crucial property: For a unimodal probability distribution on a single variable, the HDI of

massM is the narrowest possible interval of that mass. Figure 23.2 illustrates this is true. Consider the 90% HDI as shown. We construct another interval of 90% mass by moving the limits of the HDI to right, such that each limit is moved to a point that covers 4%, as marked in grey in Figure 23.2. The new interval does not cover 90%, because the 4% lost on the left is replaced by the 4% gained on the right.

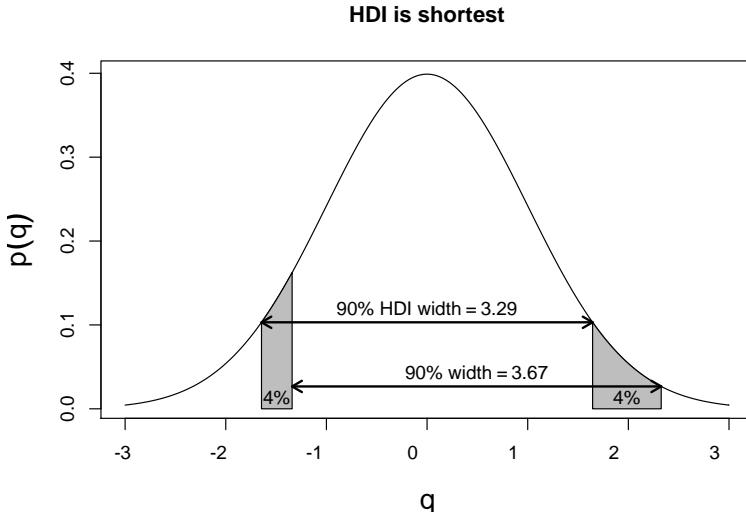


Figure 23.2: For a unimodal distribution, the HDI is the narrowest interval of that mass. This figure shows the 90% HDI and another interval that covers 90% mass.

Consider the grey regions in Figure 23.2. Their left edges are at the same height, because the left edges are defined by the HDI. Their areas are the same, because, by definition, the areas are both 4%. Notice, however, that the left grey area is taller than the right grey area, because the left edge falls at a point where the distribution is increasing in density, but the right edge falls at a point where the distribution is decreasing in density. Therefore, because the areas are the same but the left is taller than the right, the width of the left grey zone must be less than the width of the right grey zone. Consequently, the distance between right edges of the two grey zones must be greater than the HDI width. These widths are marked in Figure 23.2.

This argument applies for any size of grey zone, and for any HDI. The argument relies on unimodality, however. Given the argument and diagram in Figure 23.2, it is not too hard to believe the converse: For a unimodal distribution of one variable, for any mass M , the interval containing mass M that has the narrowest width is the HDI for that mass. The algorithms described below are based on this property of the HDI. The algorithms find the HDI by searching among candidate intervals of mass. The shortest one found is declared to be the approximate HDI. It is an approximation, of course. (See Chen and Shao (1999) for more details, and Chen, He, Shao, and Xu (2003) for dealing with the unusual situation of multimodal distributions.)

Below is the code for finding the HDI of an MCMC sample. It is brief and hopefully self-explanatory after the discussion above.

(HDIofMCMC.R)

```

1 HDIofMCMC = function( sampleVec , credMass=0.95 ) {
2   # Computes highest density interval from a sample of representative values,
```

```

3   # estimated as shortest credible interval.
4   # Arguments:
5   # sampleVec
6   # is a vector of representative values from a probability distribution.
7   # credMass
8   # is a scalar between 0 and 1, indicating the mass within the credible
9   # interval that is to be estimated.
10  # Value:
11  # HDIlim is a vector containing the limits of the HDI
12  sortedPts = sort( sampleVec )
13  cidxInc = floor( credMass * length( sortedPts ) )
14  nCIs = length( sortedPts ) - cidxInc
15  ciWidth = rep( 0 , nCIs )
16  for ( i in 1:nCIs ) {
17    ciWidth[ i ] = sortedPts[ i + cidxInc ] - sortedPts[ i ]
18  }
19  HDImin = sortedPts[ which.min( ciWidth ) ]
20  HDImax = sortedPts[ which.min( ciWidth ) + cidxInc ]
21  HDIlim = c( HDImin , HDImax )
22  return( HDIlim )
23 }
```

23.3.3 R code for computing HDI of a function

This program finds the HDI of a probability density function that is specified mathematically in R. For example, it can find HDI's of normal densities or beta densities or gamma densities, because those densities are specified as functions in R.

What the program accomplishes is just a search of HDI's, trying to the shortest one, but it does this by using some commands and R functions that have not been used much, or at all, elsewhere in the book. One function that the program uses is the inverse cumulative density function (ICDF) for whatever probability distribution is being tested. We have seen one case of an ICDF previously, namely the `pnorm`, which is the inverse of the cumulative-density function for the normal distribution. In R, the ICDF of the normal is the `qnorm(x)` function, where the argument is a value between zero and one. The program for finding the HDI takes, as one of its arguments, name for the ICDF of the function. For example, if we want to find an HDI of a normal density, we pass in `ICDFname="qnorm"`.

The crucial function called by the program is the `optimize` routine. The `optimize` routine searches the minimum of a specified function over a specified domain. In the program below, we define a function called `intervalWidth` that returns the width of the interval that starts at `lowTailPr` and has 95% mass. This `intervalWidth` function is repeatedly called from the `optimize` routine until it converges to a minimum.

(HDIofICDF.R)

```

1 HDIofICDF = function( ICDFname , credMass=0.95 , tol=1e-8 , ... ) {
2   # Arguments:
3   # ICDFname is R's name for the inverse cumulative density function
4   # of the distribution.
5   # credMass is the desired mass of the HDI region.
6   # tol is passed to R's optimize function.
7   # Return value:
8   # Highest density interval (HDI) limits in a vector.
9   # Example of use: For determining HDI of a beta(30,12) distribution, type
10  # HDIofICDF( qbta , shape1 = 30 , shape2 = 12 )
```

```

11 # Notice that the parameters of the ICDFname must be explicitly named;
12 # e.g., HDIofICDF( qbta , 30 , 12 ) does not work.
13 # Adapted and corrected from Greg Snow's TeachingDemos package.
14 incredMass = 1.0 - credMass
15 intervalWidth = function( lowTailPr , ICDFname , credMass , ... ) {
16   ICDFname( credMass + lowTailPr , ... ) - ICDFname( lowTailPr , ... )
17 }
18 optInfo = optimize( intervalWidth , c( 0 , incredMass ) , ICDFname ,
19                     credMass=credMass , tol=tol , ... )
20 HDIlowTailPr = optInfo$minimum
21 return( c( ICDFname( HDIlowTailPr , ... ) ,
22           ICDFname( credMass + HDIlowTailPr , ... ) ) )
23 }
```

23.4 Reparameterization of probability distributions

There are situations in which it is natural to express a distribution on one scale, but parameterize it mathematically on a different scale. For example, we may think intuitively of the standard deviation of a normal distribution, but have to parameterize it in terms of the precision (i.e., reciprocal of the variance). As another example, we may think intuitively of an underlying bias on an infinite scale, but have to parameterize it on a zero-to-one interval. The question is, if we express a probability distribution on one scale, what is the equivalent distribution on a transformed scale?

The answer is not difficult to figure out, especially for single parameters. Let the “destination” parameter be denoted θ and suppose that $\theta = f(x)$ for the “source” parameter x , with a monotonic and differentiable function f . Let the probability distribution on x be denoted $p(x)$. Then the corresponding probability distribution on θ is $p(\theta) = p(f^{-1}(\theta)) f'(f^{-1}(\theta))$, where $f'(x)$ is the derivative of f with respect to x .

Here's why. Consider a small (actually infinitesimal) interval under the distribution $p(x)$, at a particular value x . Call the width of the interval dx . The probability mass in that interval is the product of the density and the width, $p(x)dx$. We want to construct a density on θ , which we denote $p(\theta) = p(f(x))$, that has the same probability mass in the corresponding interval at $\theta = f(x)$. The width of the corresponding interval on θ is, by definition of derivative, $d\theta = dx/f'(x)$. So the probability mass in that interval is $p(\theta)d\theta = p(f(x))d\theta/f'(x)$. Therefore, to equate the probability masses in the corresponding intervals, we require that $p(f(x))d\theta/f'(x) = p(x)dx$, which, when re-arranged, yields $p(f(x)) = p(x)/f'(x)$.

23.4.1 Examples

We can apply the general formula to the case in which $f(x) = \text{sig}(x) = 1/(1 + \exp(-x))$, and the distribution on x is $p(x) = (f(x))^a (1 - f(x))^b = B(a; b) = a!(b!)^{-1} (1 - \exp(-x))^{a+b-1}$. Notice that the derivative of $f(x)$ is $f'(x) = \exp(-x)/(1 + \exp(-x))^2 = \text{sig}(x)(1 - \text{sig}(x)) = 1/(1 + \exp(-x))^2$. Therefore, the equivalent probability density at $\theta = f(x)$ is $p(\theta) = p(x) = p(x)/f'(x) = a!(b!)^{-1} (1 - \exp(-x))^{a+b-1} / [1/(1 + \exp(-x))^2] = B(a; b) = \text{beta}(a; b)$. The upper row of Figure 23.3 shows this situation when $a = b = 1$. An intuitive way to think of this situation is that the probability is dense near $\theta = 0$, but that is exactly where the sigmoidal transformation stretches the distribution. On the other hand, the probability on θ is sparse at large positive or large negative values, but it is exactly where the sigmoidal

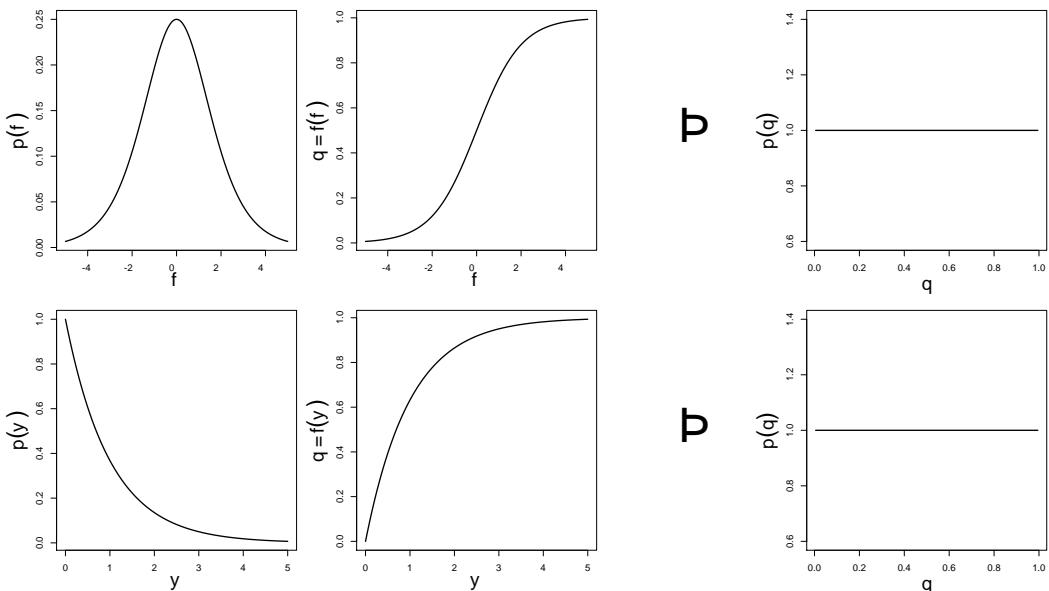


Figure 23.3: Top row shows a reparameterization that maps $[-1, 1]$ to $[0, 1]$. Bottom row shows a reparameterization that maps $[0, \infty)$ to $[0, 1]$.

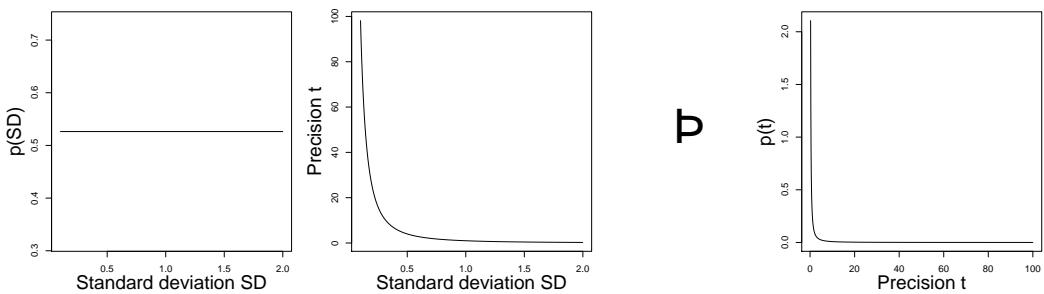


Figure 23.4: A uniform distribution on standard deviation mapped to the corresponding distribution on precision.

transformation compresses the distribution.

As another example, consider a case in which $f(\theta) = 1/\exp(\theta)$, with the probability density $p(\theta) = \exp(\theta)$. Notice that the derivative of the transformation $f(\theta) = \exp(\theta)$, and therefore the equivalent density at $f(\theta)$ is $p(f(\theta)) = p(\theta) = f'(\theta) = 1$. In other words, the equivalent density is the uniform density, as shown in the lower panel of Figure 23.3.

As a final example, Figure 23.4 shows a uniform distribution on standard deviation transformed to the corresponding distribution on precision. By definition, precision is the reciprocal of squared standard deviation.

23.4.2 Reparameterization of two parameters

When there is more than one parameter being transformed, it becomes a little more involved. Suppose we have a probability density on a parameter space $(\theta_1; \theta_2)$. Let $\theta_1 = f_1(\theta_1; \theta_2)$ and $\theta_2 = f_2(\theta_1; \theta_2)$. Our goal is to find the probability density $p(\theta_1; \theta_2)$

that corresponds to $p(x_1; x_2)$. We do this by considering infinitesimal corresponding regions. Consider a point $x_1; x_2$. The probability mass of a small region near that point is the density at that point times the area of the small region $d_1 d_2$. The corresponding region in the transformed parameters should have the same mass. That mass is the density at the transformed point times the area of the region mapped to from the originating region. In vector calculus textbooks, you can find discussions demonstrating that the area of the mapped-to region is $\det(J) d_1 d_2$ where J is the Jacobian matrix: $J_{rc} = df_r(x_1; x_2) = d_c$ and $\det(J)$ is the determinant of the Jacobian matrix. Setting the two masses equal and re-arranging yields $p(x_1; x_2) = p(r_1; r_2) \cdot \det(J)$.

As we have had no occasions to apply this transformation, examples will be provided here. But the method has been mentioned for those intrepid who may wish to venture into the wilderness of multivariate probability distributions with nothing more than pen and paper.

References

- Adcock, C. J. (1997). Sample size determination: a review. *The Statistician*, 46, 261–283.
- Agresti, A., & Hitchcock, D. B. (2005). Bayesian inference for categorical data analysis. *Statistical Methods & Applications*, 14(3), 297–330.
- Albert, J. H., & Rossman, A. J. (2001). *Workshop statistics: Discovery with data, a Bayesian approach*. Emeryville, CA: Key College Publishing.
- Berger, J. O. (1985). *Statistical decision theory and Bayesian analysis*, 2nd edn. New York: Springer.
- Berger, J. O., & Berry, D. A. (1988). Statistical analysis: the illusion of objectivity. *American Scientist*, 76(2), 159–165.
- Berger, R. L., Boos, D. D., & Guess, F. M. (1988). Tests and confidence sets for comparing two mean residual life functions. *Biometrics*, 44(1), 103–115.
- Berry, D. A. (1996). *Statistics: A Bayesian perspective*. Belmont, CA: Duxbury Press/Wadsworth.
- Berry, D. A., & Hochberg, Y. (1999). Bayesian perspectives on multiple comparisons. *Journal of Statistical Planning and Inference*, 82(1-2), 215–227.
- Bishop, C. M. (2006). *Pattern recognition and machine learning*. New York: Springer.
- Bliss, C. I. (1934). The method of probits. *Science*, 79(2037), 38–39.
- Bolstad, W. M. (2007). *Introduction to Bayesian statistics*, 2nd ed.). Hoboken, NJ: Wiley.
- Brambor, T., Clark, W. R., & Golder, M. (2006). Understanding interaction models: Improving empirical analyses. *Political Analysis*, 14, 63–82.
- Braumoeller, B. F. (2004). Hypothesis testing and multiple interaction terms. *International Organization*, 58(04), 807–820.
- Brehmer, B. (1974). Hypotheses about relations between variables in the learning of probabilistic inference tasks. *Organizational Behavior and Human Performance*, 1, 1–27.
- Brooks, S. P., & Gelman, A. (1998). Alternative methods for monitoring convergence of iterative simulations. *Journal of Computational and Graphical Statistics*, 7(3), 434–455.
- Carlin, B. P., & Chib, S. (1995). Bayesian model choice via *Markov chain Monte Carlo* methods. *Journal of the Royal Statistical Society, Series B*, 57(3), 473–484.

- Carlin, B. P., & Louis, T. A. (2000). *Bayes and empirical Bayes methods for data analysis* (2nd ed.). Boca Raton, FL: Chapman & Hall/CRC.
- Carlin, B. P., & Louis, T. A. (2009). *Bayesian methods for data analysis* (3rd ed.). Boca Raton, FL: CRC Press.
- Casella, G., & Moreno, E. (2006). Objective Bayesian variable selection. *Journal of the American Statistical Association*, 101(473), 157–167.
- Chaloner, K., & Verdinelli, I. (1995). Bayesian experimental design: A review. *Statistical Science*, 10(3), 273–304.
- Chen, M.-H., He, X., Shao, Q.-M., & Xu, H. (2003). A Monte Carlo test in computing HPD regions. *Development of Modern Statistics and Related Topics: In Celebration of Professor Yaoting Zhang's 70th Birthday*, 38–52.
- Chen, M. H., & Shao, Q. M. (1999). Monte Carlo estimation of Bayesian credible and HPD intervals. *Journal of Computational and Graphical Statistics*, 8, 69–92.
- Clyde, M., & George, E. I. (2004). Model uncertainty. *Statistical Science*, 21, 81–94.
- Congdon, P. (2005). *Bayesian models for categorical data*. West Sussex, England: Wiley.
- Damgaard, L. H. (2007). Technical note: How to use WinBUGS to draw inferences in animal models. *Journal of Animal Sciences*, 85, 1363–1368.
- Dawes, J. (2008). Do data characteristics change according to the number of scale points used? an experiment using 5-point, 7-point and 10-point scales. *International Journal of Market Research*, 50(1), 61–77.
- de Saint-Exupéry, A. (1943). *The little prince*. San Diego: Harcourt.
- De Santis, F. (2004). Statistical evidence and sample size coordination for Bayesian hypothesis testing. *Journal of Statistical Planning and Inference*, 124, 121–144.
- De Santis, F. (2007). Using historical data for Bayesian sample size determination. *Journal of the Royal Statistical Society: Series B*, 70, 95–113.
- DeGroot, M. H. (2004). *Optimal statistical decisions*. New York: Wiley Interscience.
- Edwards, W., Lindman, H., & Savage, L. J. (1963). Bayesian statistical inference for psychological research. *Psychological Review*, 70, 193–242.
- Feldman, H. A. (1988). Families of lines: random effects in linear regression analysis. *Journal of Applied Physiology*, 64(4), 1721–1732.
- Freedman, L. S., Lowe, D., & Macaskill, P. (1984). Stopping rules for clinical trials incorporating clinical opinion. *Biometrics*, 40, 575–586.
- Gallistel, C. R. (2009). The importance of proving the null. *Psychological Review*, 116(2), 439–453.
- Gelfand, A. E., & Dey, D. K. (1994). Bayesian model choice via posterior odds ratios and exact calculations. *Journal of the Royal Statistical Society, Series B*, 50, 501–514.

- Gelman, A. (2005). Analysis of variance — why it is more important than ever. *The Annals of Statistics* 33(1), 1–53.
- Gelman, A. (2006). Prior distributions for variance parameters in hierarchical models. *Bayesian Analysis* 1(3), 515–533.
- Gelman, A., Carlin, J. B., Stern, H. S., & Rubin, D. B. (2004). Bayesian data analysis (2nd ed.). Boca Raton, Florida: CRC Press.
- Gelman, A., & Hill, J. (2007). Data analysis using regression and multilevel/hierarchical models New York: Cambridge University Press.
- Gelman, A., Hill, J., & Yajima, M. (2009). Why we (usually) don't have to worry about multiple comparisons. Available from <http://www.stat.columbia.edu/~gelman/research/unpublished/multiple2.pdf>
- Geman, S., & Geman, D. (1984). Stochastic relaxation, Gibbs distributions, and the Bayesian restoration of images. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 6, 721–741.
- George, D. N., & Pearce, J. M. (1999). Acquired distinctiveness is controlled by stimulus relevance not correlation with reward. *Journal of Experimental Psychology: Animal Behavior Processes* 25(3), 363–373.
- George, E. I. (2000). The variable selection problem. *Journal of the American Statistical Association* 95(452).
- Gigerenzer, G., & Hrage, U. (1995). How to improve Bayesian reasoning without instruction: Frequency formats. *Psychological Review* 102, 684–704.
- Gigerenzer, G., Krauss, S., & Vitouch, O. (2004). The null hypothesis: What you always wanted to know about significance testing but were afraid to ask. In K. A. Karpman (Ed.), *The Sage handbook of quantitative methodology for the social sciences* (pp. 391–408). Thousand Oaks, CA: Sage.
- Gilks, W. R., Thomas, A., & Spiegelhalter, D. J. (1994). A dage and program for complex Bayesian modelling. *The Statistician* 43(1), 169–177.
- Gill, J. (2002). Bayesian methods for the social and behavioral sciences. Boca Raton, Florida: CRC Press.
- Gilovich, T., Vallone, R., & Tversky, A. (1985). The hot hand in basketball: On the misperception of random sequences. *Cognitive Psychology* 17(3), 295–314.
- Gopalan, R., & Berry, D. (1998). Bayesian multiple comparisons using Dirichlet process priors. *Journal of the American Statistical Association* 93(431), 1130–1139.
- Gosset, W. S. (1908). The probable error of a mean. *Biometrika* 6, 1–25.
- Greenland, S. (2008). Invited commentary: variable selection versus shrinkage in the control of multiple confounders. *American Journal of Epidemiology* 167(5), 523–529.

- Guber, D. L. (1999). Getting what you pay for: The debate over equity in public school expenditures. *Journal of Statistics Education* 7(2). Available from <http://www.amstat.org/publications/JSE/secure/v7n2/datasets.guber.cfm>
- Hahn, U., Chater, N., & Richardson, L. B. (2003). Similarity transformationCognition 87(1), 1–32.
- Han, C., & Carlin, B. P. (2001). Markov chain Monte Carlo methods for computing Bayes factors: A comparative reviewJournal of the American Statistical Association 96(455), 1122–1132.
- Hand, D. J., Daly, F., Lunn, A. D., McConway, K. J., & Ostrowski, S. (1994).A handbook of small data setsLondon: Chapman & Hall.
- Hobbs, B. P., & Carlin, B. P. (2008, January). Practical Bayesian design and analysis for drug and device clinical trialsJournal of Biopharmaceutical Statistics 18(1), 54–80.
- Homan, P. J., Earle, T. C., & Slovic, P. (1981). Multidimensional functional learning (MFL) and some new conceptions of feedbackOrganizational Behavior and Human Performance 27(1), 75–102.
- Holcomb, J., & Spalsbury, A. (2005). Teaching students to use summary statistics and graphics to clean and analyze dataJournal of Statistics Education 13(3). Available from http://www.amstat.org/publications/jse/v13n3/database_ts.holcomb.html
- Hyndman, R. J. (1996). Computing and graphing highest density regions.The American statistician 50(2), 120–126.
- Joseph, L., Wolfson, D. B., & du Berger, R. (1995a). Sample calculations for binomial proportions via highest posterior density intervalsThe Statistician 44, 143–154.
- Joseph, L., Wolfson, D. B., & du Berger, R. (1995b). Some comments on Bayesian sample size determinationThe Statistician 44, 167–171.
- Kalish, M. L., Griffiths, T. L., & Lewandowsky, S. (2007). Iterated learning of generational knowledge transmission reveals inductive biasesPsychonomic Bulletin & Review 14(2), 288.
- Kass, R. E., & Raftery, A. E. (1995). Bayes factorsJournal of the American Statistical Association 90, 773–795.
- Keith, T. (2005).Multiple regression and beyondColumbus, OH: Allyn & Bacon.
- Kolmogorov, A. N. (1956).Foundations of the theory of probabilityNew York: Chelsea.
- Krauss, S., Martignon, L., & Horage, U. (1999). Simplifying Bayesian inference: The general case. In L. Magnani, N. J. Nersessian, & P. Thagard (Eds.), Model-based reasoning in scientific discovery(pp. 165–180). New York: Springer.
- Kruschke, J. K. (1993). Human category learning: Implications for backpropagation models.Connection Science 5, 3–36.

- Kruschke, J. K. (1996). Dimensional relevance shifts in category learning. *Connection Science* 8, 201–223.
- Kruschke, J. K. (2008). Bayesian approaches to association learning: From passive to active learning. *Learning & Behavior* 36(3), 210–226.
- Kruschke, J. K. (2009). Highlighting: A canonical experiment. In B. Ross (Ed.), *The psychology of learning and motivation* (Vol. 51, pp. 153–185). Elsevier/Academic Press.
- Kruschke, J. K. (2010). Bayesian data analysis. *Wiley Interdisciplinary Reviews: Cognitive Science* ** (**), **-**.
- Learner, E. E. (1978). Specification searches. New York: Wiley.
- Lee, M. D., & Webb, M. R. (2005). Modeling individual differences in cognition. *Psychonomic Bulletin & Review* 12(4), 605–621.
- Liang, F., Paulo, R., Molina, G., Clyde, M. A., & Berger, J. O. (2008). Mixtures of g priors for Bayesian variable selection. *Journal of the American Statistical Association* 103, 410–423.
- Likert, R. (1932). A technique for the measurement of attitudes. *Archives of Psychology* 14, 1–55.
- Lindley, D. V. (1997). The choice of sample size. *The Statistician* 46, 129–138.
- Lindley, D. V., & Phillips, L. D. (1976). Inference for a Bernoulli process (a Bayesian view). *The American Statistician* 30(3), 112–119.
- Lindquist, M. A., & Gelman, A. (2009). Correlations and multiple comparisons in functional imaging – a statistical perspective. *Perspectives in Psychological Science* 4(3), 310–313.
- Liu, C. C., & Aitkin, M. (2008). Bayes factors: Prior sensitivity and model generalizability. *Journal of Mathematical Psychology* 52, 362–375.
- Lynch, S. M. (2007). *Introduction to applied Bayesian statistics and estimation for social scientists*. New York: Springer.
- MacKay, D. J. C. (2003). *Information theory, inference & learning algorithms*. Cambridge, UK: Cambridge University Press.
- Marin, J.-M., & Robert, C. P. (2007). *Bayesian core: A practical approach to computational Bayesian statistics*. New York: Springer.
- Maxwell, S. E., & Delaney, H. D. (2004). *Designing experiments and analyzing data: a model comparison perspective* (2nd ed.). Mahwah, NJ: Erlbaum.
- McCullagh, P., & Nelder, J. (1989). *Generalized linear models*, 2nd ed. Boca Raton, FL: Chapman and Hall/CRC.
- McDonald, J. H. (2009). *Handbook of biological statistics* (2nd ed.). Baltimore, Maryland: Sparky House Publishing.

- McDonald, J. H., Seed, R., & Koehn, R. K. (1991). Allozymes and morphometric characters of three species of *Mytilus* in the Northern and Southern Hemispheres. *Marine Biology*, 111(3), 323–333.
- McIntyre, L. (1994). Using cigarette data for an introduction to multiple regression. *Journal of Statistics Education* 2(1). Available from http://www.amstat.org/publications/jse/v2n1/dataset_s.mcintyre.html
- Meng, C. Y. K., & Dempster, A. P. (1987). A Bayesian approach to the multiplicity problem for significance testing with binomial data. *Biometrics* 43(2), 301–311.
- Metropolis, N., Rosenbluth, A. W., Rosenbluth, M. N., Teller, H., & Teller, E. (1953). Equations of state calculations by fast computing machines. *Journal of Chemical Physics* 21, 1087–1091.
- Meyer, R., & Yu, J. (2000). BUGS for a Bayesian analysis of stochastic volatility models. *Econometrics Journal* 3(2), 198–215.
- Miller, J. (2009). What is the probability of replicating statistically significant effects? *Psychonomic Bulletin & Review* 16(4), 617–640.
- Moore, T. L. (2006). Paradoxes in film ratings. *Journal of Statistics Education* 14(1). Available from www.amstat.org/publications/jse/v14n1/datasets.moore_e.html
- Mueller, P., Parmigiani, G., & Rice, K. (2007). FDR and Bayes multiple comparisons rules. In J. M. Bernardo et al. (Eds) *Bayesian statistics 8* Oxford, UK: Oxford University Press.
- Navarro, D. J., Griffiths, T. L., Steyvers, M., & Lee, M. D. (2006). Modeling individual differences using Dirichlet processes. *Journal of Mathematical Psychology* 50, 101–122.
- Nelder, J. A., & Wedderburn, R. W. M. (1972). Generalized linear models. *Journal of the Royal Statistical Society. Series A (General)* 125(3), 370–384.
- Nietzsche, F. (1967). *The will to power* New York: Random House. (Translated by W. Kaufmann and R. J. Hollingdale)
- Ntzoufras, I. (2009). Bayesian modeling using WinBUGS. Hoboken, NJ: Wiley.
- Oswald, C. J. P., Yee, B. K., Rawlins, J. N. P., Bannerman, D., Goddard, M., & Honey, R. C. (2001). Involvement of the entorhinal cortex in a process of attentional modulation: Evidence from a novel variant of an IEDS procedure. *Behavioral Neuroscience* 115(4), 841–849.
- Pham-Gia, T., & Turkkan, N. (1992). Sample size determination in Bayesian analysis. *The Statistician* 41, 389–392.
- Poldrack, R. A. (2006). Can cognitive processes be inferred from neuroimaging data? *Trends in Cognitive Sciences* 10(2), 59–63.
- Proschan, F. (1963). Theoretical explanation of observed decreasing failure rate. *Technometrics* 375–383.

- Qian, S. S., & Shen, Z. (2007). Ecological applications of hierarchical analysis of variance. *Ecology* 88(10), 2489–2495.
- Robert, C. P., & Casella, G. (2004). *Monte Carlo statistical methods* (2nd ed.). New York: Springer.
- Rosa, L., Rosa, E., Sarner, L., & Barrett, S. (1998). A cross-cultural therapeutic touch. *Journal of the American Medical Association* 279(13), 1005–1010.
- Rouder, J. N., & Lu, J. (2005). An introduction to Bayesian hierarchical models with an application in the theory of signal detection. *Psychonomic Bulletin & Review* 12(4), 573–604.
- Rouder, J. N., Lu, J., Speckman, P., Sun, D., & Jiang, Y. (2005). A hierarchical model for estimating response time distributions. *Psychonomic Bulletin & Review* 12(2), 195–223.
- Rouder, J. N., Speckman, P. L., Sun, D., Morey, R. D., & Iverson, G. (2009). Bayesian t-tests for accepting and rejecting the null hypothesis. *Psychonomic Bulletin & Review* 16, 225–237.
- Roy, A., Ghosal, S., & Rosenberger, W. F. (2009). Convergence properties of sequential Bayesian D-optimal designs. *Journal of Statistical Planning and Inference* 139, 425–440.
- Sadiku, M. N. O., & Toghi, M. R. (1999). A tutorial on simulation of queueing models. *International Journal of Electrical Engineering Education* 36, 102–120.
- Scott, J. G., & Berger, J. O. (2006). An exploration of aspects of Bayesian multiple testing. *Journal of statistical planning and inference* 136(7), 2144–2162.
- Snee, R. D. (1974). Graphical display of two-way contingency tables. *The American Statistician* 28(1), 9–12.
- Solari, F., Liseo, B., & Sun, D. (2008). Some remarks on Bayesian inference for one-way ANOVA models. *Annals of the Institute of Statistical Mathematics* 60, 483–498.
- Spiegelhalter, D. J., Freedman, L. S., & Parmar, M. K. B. (1999). Bayesian approaches to randomized trials. *Journal of the Royal Statistical Society. Series A* 157, 357–416.
- Stevens, S. S. (1946). On the theory of scales of measurement. *Science* 103(2684), 677–680.
- Thomas, A. (2004). BRugs user manual (the R interface to BUGS). Available from <http://mathstat.helsinki.fi/openbugs/data/Docu/BRugsManual.html>
- Thomas, A., O'Hara, B., Ligges, U., & Sturtz, S. (2006, March). Making BUGS openR News 6(1), 12–17.
- Tsionas, E. G. (2002). Bayesian inference in the noncentral Student-t model. *Journal of Computational and Graphical Statistics* 11(1), 208–221.
- Wagenmakers, E. J. (2007). A practical solution to the pervasive problems of p values. *Psychonomic Bulletin & Review* 14(5), 779–804.

- Walker, L. J., Gustafson, P., & Frimer, J. A. (2007). The application of Bayesian analysis to issues in developmental research. *International Journal of Behavioral Development*, 31(4), 366.
- Wang, F., & Gelfand, A. E. (2002). A simulation-based approach to Bayesian sample size determination for performance under a given model and for assessing models. *Statistical Science*, 17, 193–208.
- Weiss, R. (1997). Bayesian sample size calculations for hypothesis testing. *The Statistician*, 46, 185–191.
- Werner, M., Stabenau, J. R., & Pollin, W. (1970). The thematic perception test method for the differentiation of families of schizophrenics, delinquents and "normals". *Journal of Abnormal Psychology*, 75(2), 139–145.
- Western, B., & Jackman, S. (1994). Bayesian inference for comparative research. *The American Political Science Review*, 88(2), 412–423.
- Winer, B. J., Brown, D. R., & Michels, K. M. (1991). *Statistical principles in experimental design*, 3rd ed. New York: McGraw-Hill.

The Index on the following pages is incomplete. It includes ~~only~~ a few items to test the indexing facility. The index will be ~~expanded~~ at a later date.

Index

- aggregate , 188
- Bernoulli distribution, 66
- Bernoulli versus binomial, 67
- binomial probability distribution, 217
- Binomial versus Bernoulli, 67
- Bonferroni correction, 228
- BUGS, 115
 - categorical density function, 198
 - censoring in BUGS, 171, 188
 - dcat , 198
 - dev.copy2eps , 17, 20
 - dgamma170
- encapsulated PostScript, 20
- entropy, 264
- EPS format, 20
- exchangeability, 165
- experimentwise false alarm rate, 228
- Iteration and condensation, 179
- gamma170
- gamma distribution, 170
- gamma function, 170
- Grinch, 252
- help in R, 17
- highest density interval, 34
- I(lower,upper) in BUGS, 171
- logistic, 305
- logit, 306
- MCMC: Markov chain Monte Carlo, 109
- natural frequencies, 60
- Markov representation, 60
- negative binomial distribution, 220
 - in R, 233
- nested indexing in BUGS, 172
- OpenBUGS, 115
 - per comparison false alarm rate, 228
 - planned comparison, 229
 - Poisson distribution, 235, 492, 493
 - post-hoc comparison, 229
 - posterior predictive check, 81, 232
 - precision of normal distribution, 320
 - probit, 308
 - product space, 208
 - pseudoprior, 202
- R programs
 - ANOVAonewayBRugs.R, 405, 413
 - ANOVAonewayNonhomogvarBRugs.R, 418, 419
 - ANOVAtwowayBRugs.R, 436
 - ANOVAtwowayBRugsWithinSubj.R, 446
 - BayesUpdate.R, 57
 - BernBeta.R, 77
 - BernBetaBugsFull.R, 116–118
 - BernBetaModelCompBRugs.R, 198
 - BernBetaMuKappaBugs.R, 171, 185
 - BernGrid.R, 90
 - BernMetropolisTemplate.R, 121
 - BernTwoBugs.R, 140, 142, 149
 - BernTwoBugsPriorOnly.R, 141
 - BernTwoFurrowsBugs.R, 154
 - BernTwoGrid.R, 144
 - BernTwoMetropolis.R, 146
 - BetaPosteriorPredictions.R, 82
 - BinomNHSTpoissonrate.R, 235
 - FilconBRugs.R, 180, 188
 - FilconBRugsPower.R, 276

- FilconCoKappaBrugs.R, 193
FilconModelCompBrugs.R, 200
FilconModelCompPseudoPrior-
 Brugs.R, 202
HDlofGrid.R, 513
HDlofICDF.R, 515
HDlofMCMC.R, 514
HtWtDataGenerator.R, 358
IntegralOfDensity.R, 41
Kruschke1996CSbugs.R, 283, 284, 286
LogisticOneWayAnovaBrugs.R, 463
LogisticOneWayAnovaHeteroVar-
 Brugs.R, 470
minNforHDIpower.R, 273
MultiLinRegressHyperBrugs.R, 394
MultiLinRegressInterBrugs.R, 384
MultipleLinearRegressionBrugs.R,
 375, 378, 390
MultipleLogisticRegressionBrugs.R,
 460
NHSTTwoTierStoppingExercise.R, 237
OneOddGroupModelComp.R, 252
OneOddGroupModelCompEx12.1.R,
 255
OrdinalProbitRegressionBrugs.R, 482
plotChains.R, 512
plotPost.R, 151
PoissonExponentialBrugs.R, 497
RunningProportion.R, 40
SimpleGraph.R, 17
SimpleLinearRegressionBrugs.R, 347
SimpleLinearRegressionRepeated-
 Brugs.R, 355, 362
SimpleRobustLinearRegression-
 Brugs.R, 353, 359
SystemsBrugs.R, 331, 335
ToyModelComp.R, 211
YmetricXsingleBrugs.R, 323, 333
R: help, 17
R: Tinn-R editor, 18
region of practical equivalence (ROPE), 72

sampling distribution, 218
Santa Claus, 252
sigmoid, 305

t distribution, 323, 324
 folded, 403–405, 418
for outliers, 325, 331, 352, 353, 368,
 379
thematic apperception test, 256
Tinn-R (R editor), 18
transdimensional MCMC, 197