

# Non-Linear Regression Analysis

By Chanaka Kaluarachchi

Otago : Unibersity



# Presentation outline

- Linear regression
- Checking linear Assumptions
- Linear vs non-linear
- Non linear regression analysis

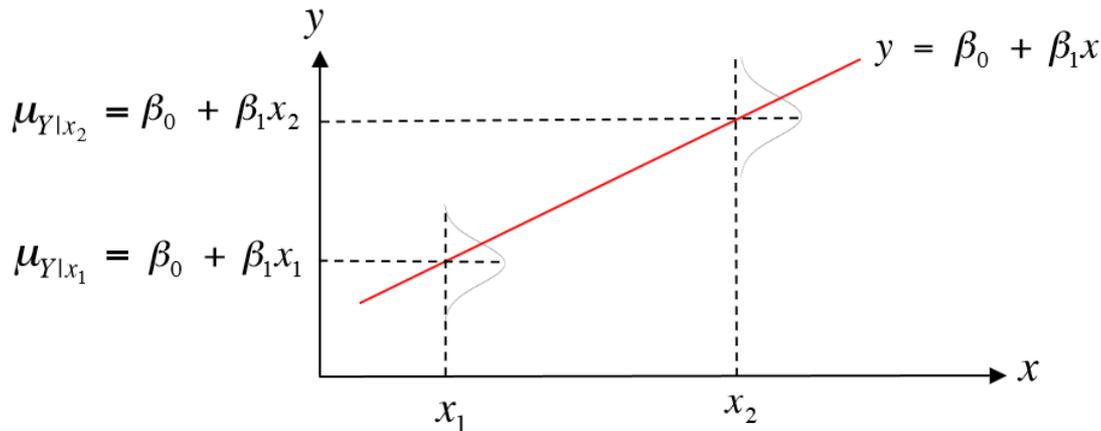
# Linear regression (reminder)

- Linear regression is an approach for modelling dependent variable( $y$ ) and one or more explanatory variables ( $x$ ).

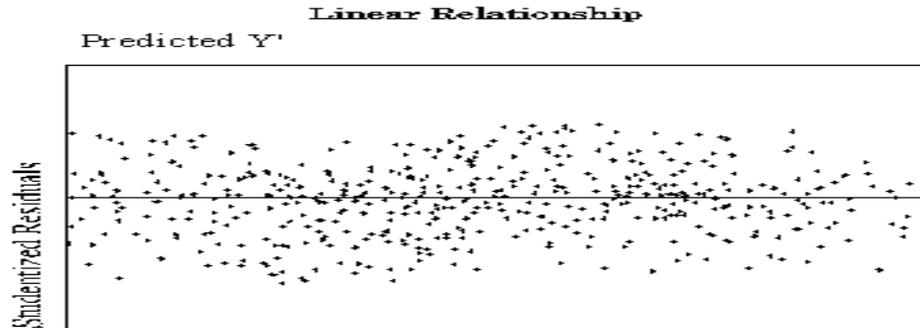
$$y = \beta_0 + \beta_1 x + \varepsilon$$

Assumptions:

$\varepsilon \sim N(0, \sigma^2)$  – iid (independently identically distributed)



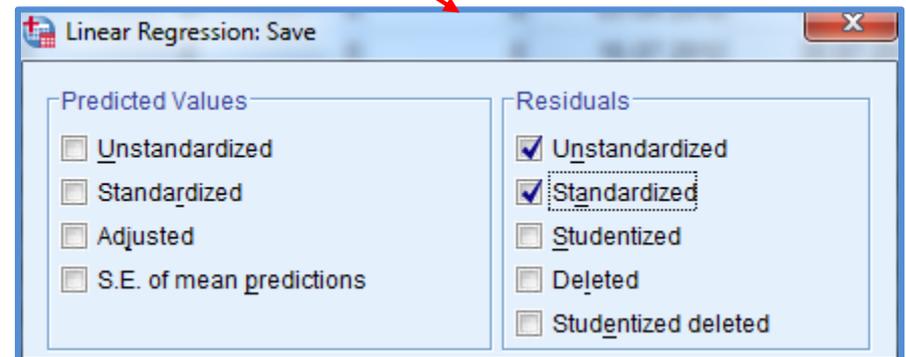
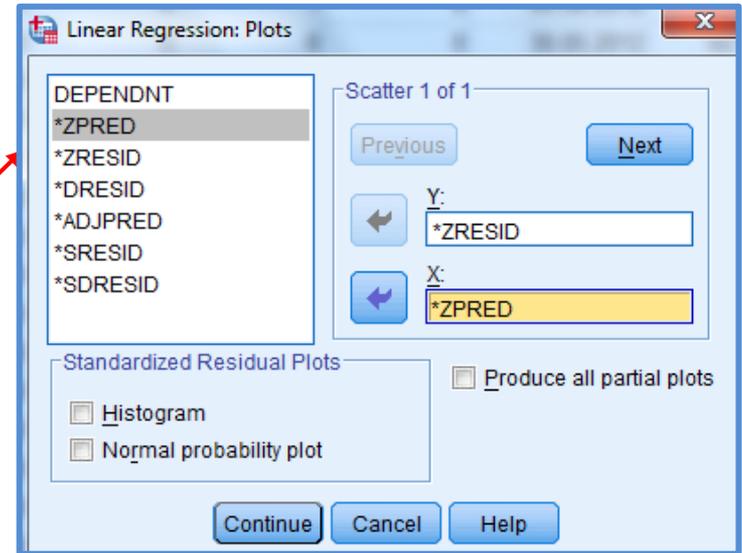
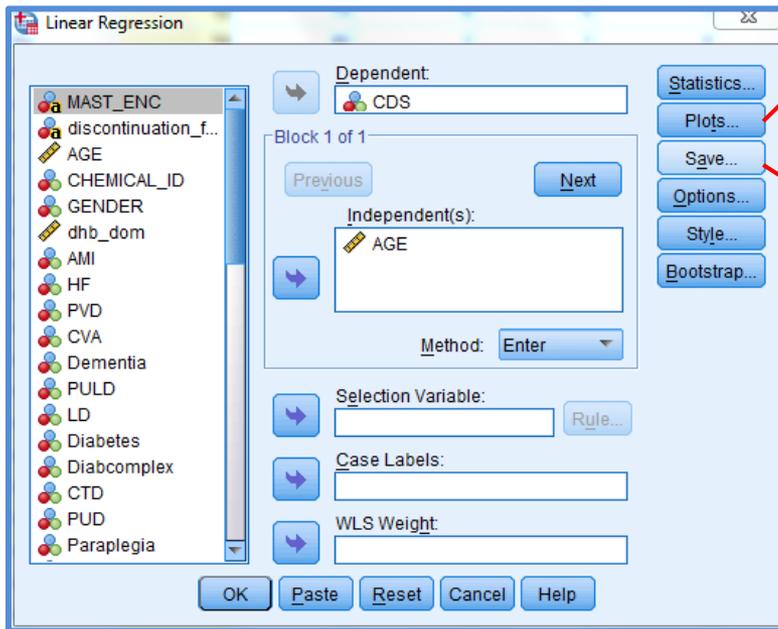
# Checking linear Assumptions



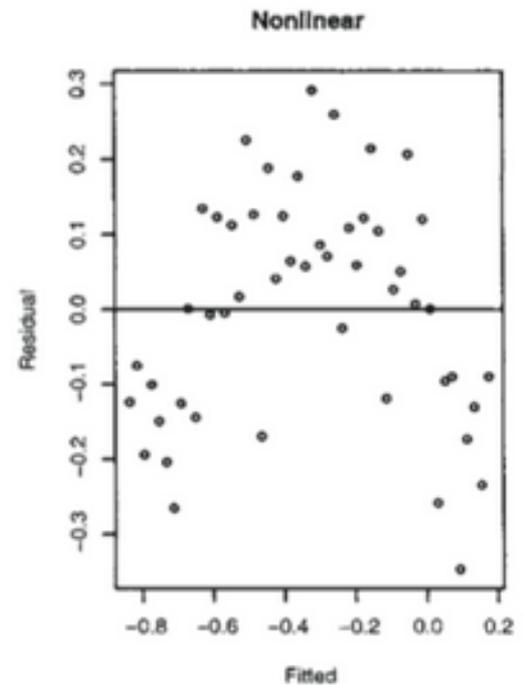
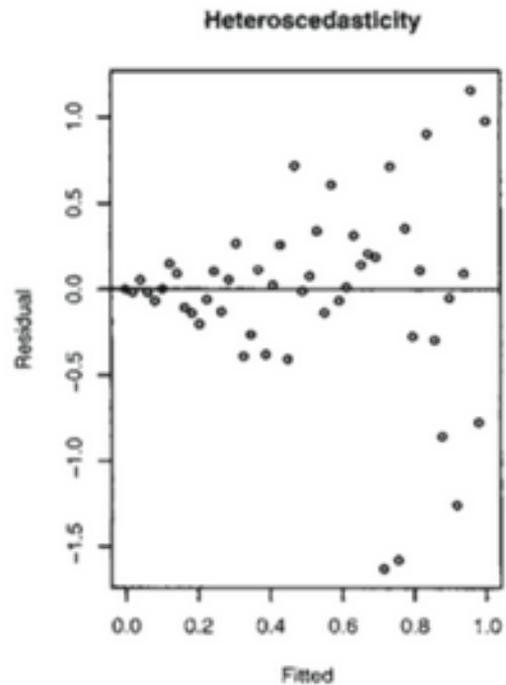
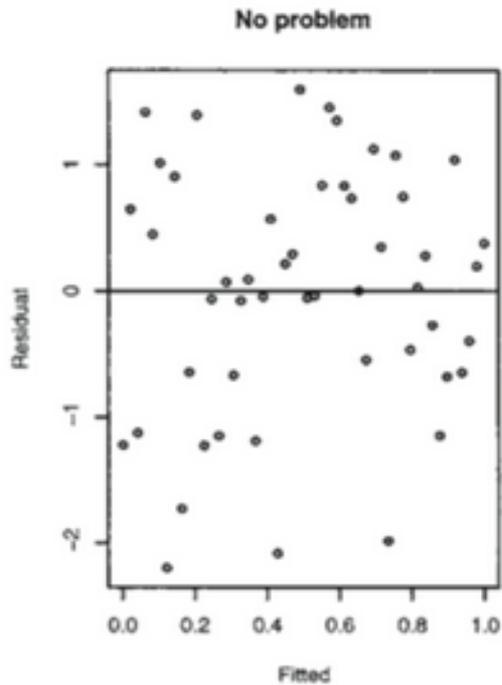
iid- residual plot ( $\varepsilon$  vs  $\hat{y}$ ) can be inspect to check that assumptions are met.

- Constant variance- Scattering is a constant magnitude
- Normal data- few outliers, systematic spared above and below the axis
- Liner relationship- No curve in the residual plot

# Residual plot in SPSS



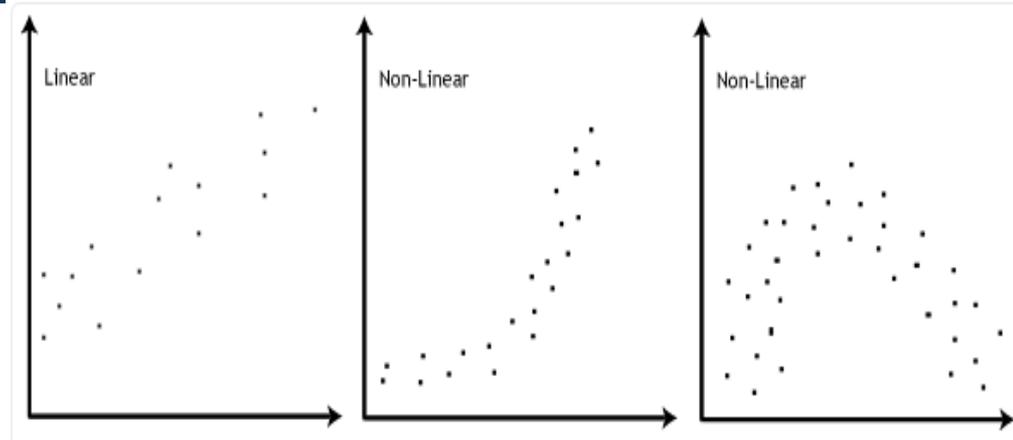
# Residual plots in SPSS



# Linear vs non linear

## Linear

- Linear scatter plot
- No curves in residual plot
- Correlation between variable is significant



## Non-linear

- Curves in scatter plot
- Curves in residual plot
- No significant correlation between variables

# Non linear regression

- Non linear regression arises when predictors and response follows particular function form.

$$y = f(\beta, x) + \varepsilon$$

## Examples

$$y = \beta^2 x + \varepsilon \quad \text{- non linear}$$

$$y = \beta x^2 + \varepsilon \quad \text{- linear}$$

$$y = \frac{1}{\beta} x + \varepsilon \quad \text{- non linear}$$

$$y = \beta \frac{1}{x} + \varepsilon \quad \text{- linear}$$

$$y = e^{\beta x} + \varepsilon \quad \text{- non linear}$$

$$y = \beta \ln x + \varepsilon \quad \text{- linear}$$

$$y = \frac{1}{1+\beta x} + \varepsilon \quad \text{- non linear}$$

# Transformation

- Some nonlinear regression problems can be moved to a linear domain by a suitable transformation of the model formulation.
- Four common transformations to induce linearity are: logarithmic transformation, square root transformation, inverse transformation and the square transformation

## Examples

$$\bullet \quad y = e^{\beta x} \quad \longrightarrow \quad \ln y = \beta x \quad \text{if } y \geq 0$$

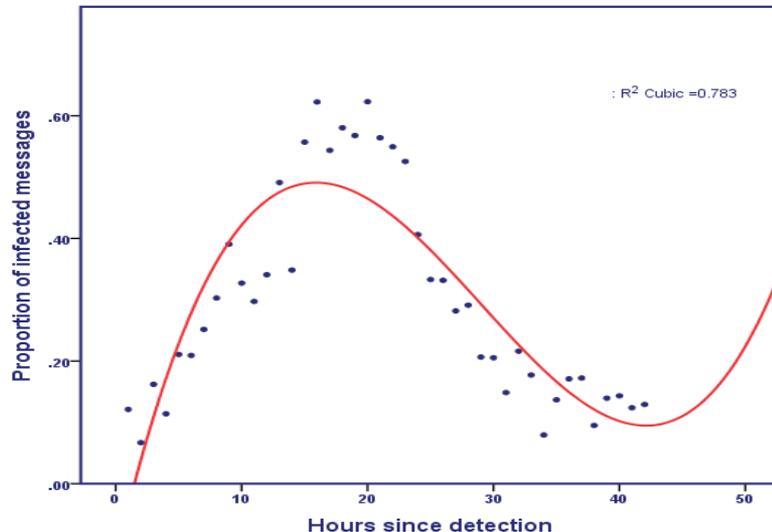
$$\bullet \quad y = \frac{1}{1+\beta x} \quad \longrightarrow \quad \frac{1}{y} - 1 = \beta x \quad \text{if } y \neq 0$$

# Curve Estimation

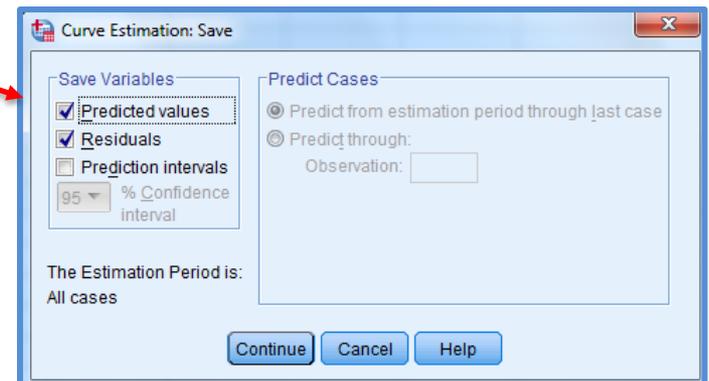
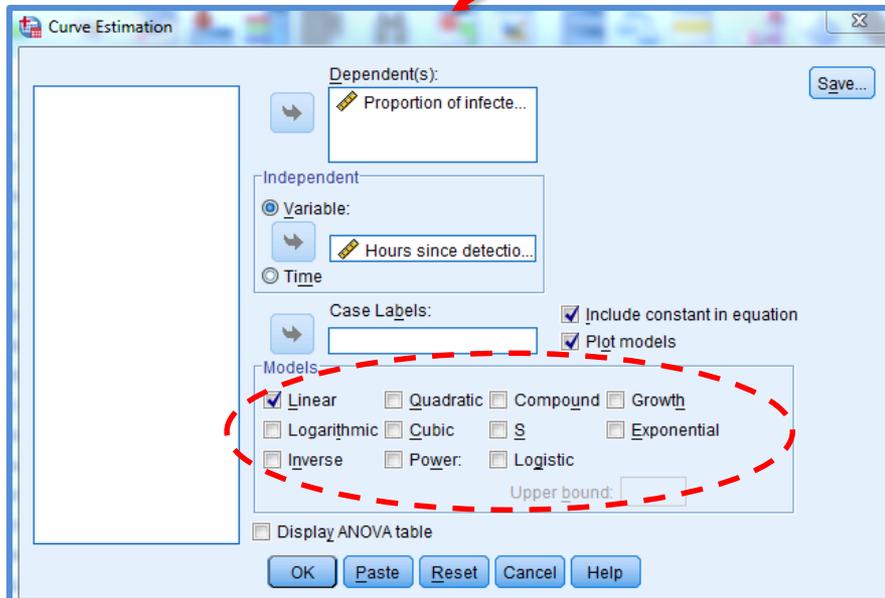
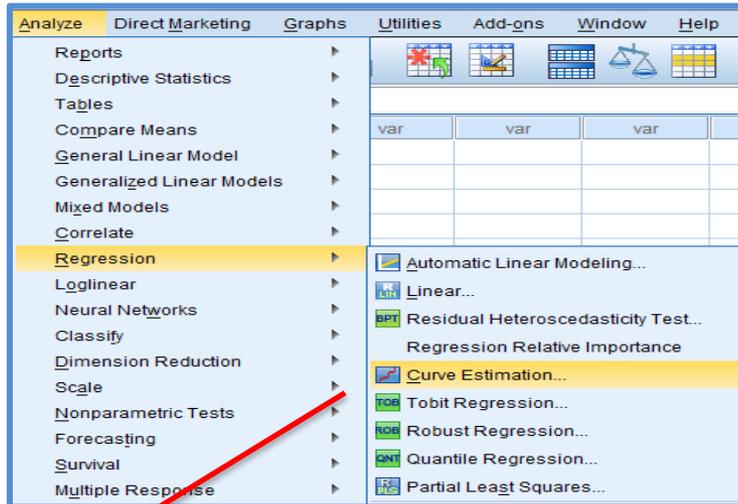
Curve fitting is the process of constructing a curve, or mathematical function, that has the best fit to a series of data points.

## Example – Viral growth model

- An internet service provider (ISP) is determining the effects of a virus on its networks. As part of this effort, they have tracked the (approximate) percentage of e-mail traffic on its networks over time, from the moment of discovery until the threat was contained.



# Curve Estimation- Cont.



# Output

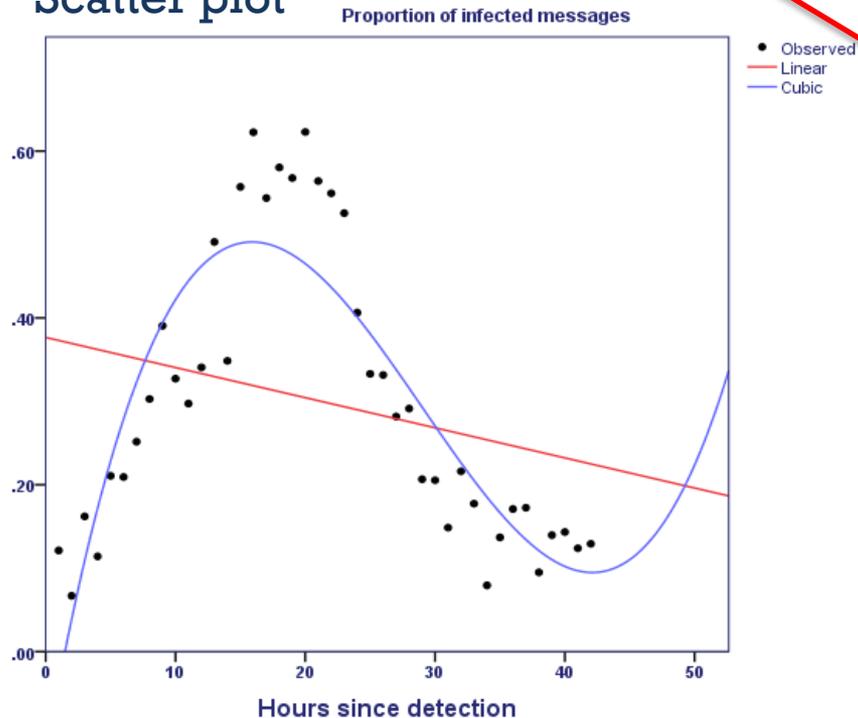
## Model Summary and Parameter Estimates

Dependent Variable: Proportion of infected messages

Equation	Model Summary					Parameter Estimates			
	R Square	F	df1	df2	Sig.	Constant	b1	b2	b3
Linear	.066	2.844	1	40	.100	.377	-.004		
Cubic	.783	45.736	3	38	.000	-.123	.088	-.004	4.399E-5

The independent variable is Hours since detection.

## Scatter plot



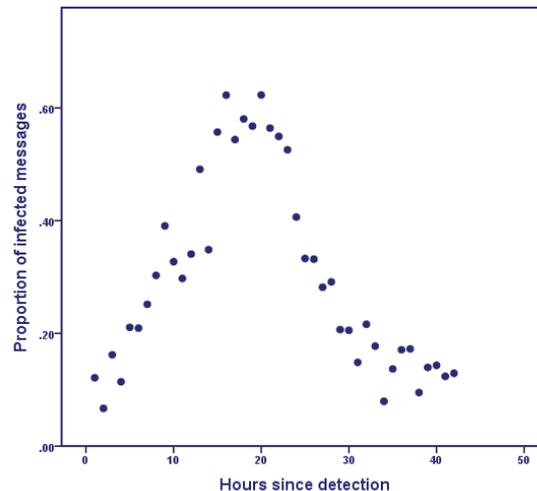
Higher the  $R^2$  better the model fit

P value < 0.05 means model is significant

# Segmentation

We can split the graph in to segments and fit a segmented model.

## Example – Viral growth model



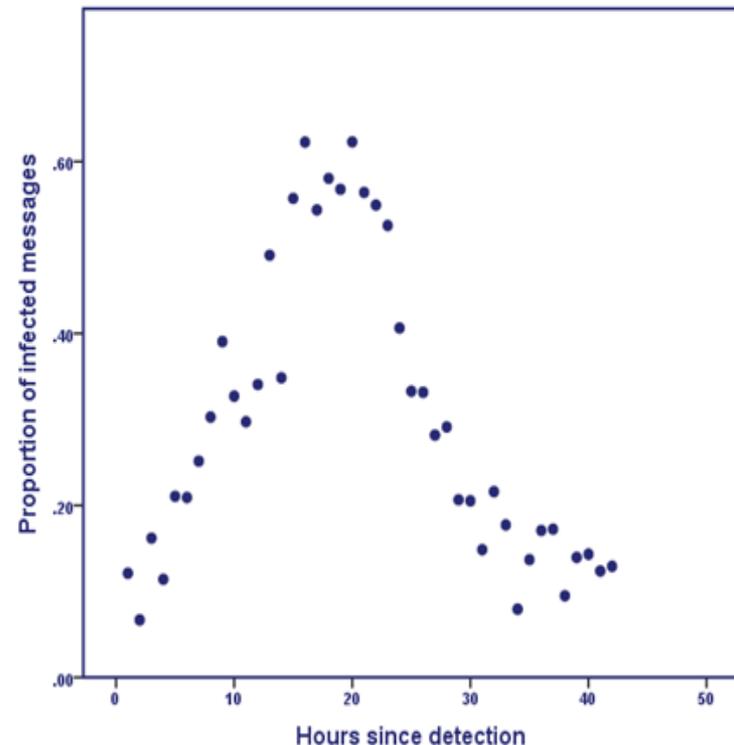
We can fit a logistic equation for the first 19 hours and an asymptotic regression for the remaining hours should provide a good fit and interpretability over the entire time period.

# Logistic model and choosing starting values

$$y = \frac{\beta_1}{1 + \beta_2 e^{-\beta_3 x}}$$

## Starting values

- $\beta_1$  - upper value of growth (0.65)
- $\beta_2$  - ratio upper value and lowest value ( $0.65/0.13=5$ )
- $\beta_3$  - estimated slope between points in plot.  
( $0.6-0.12/19-3$ )= $0.03$

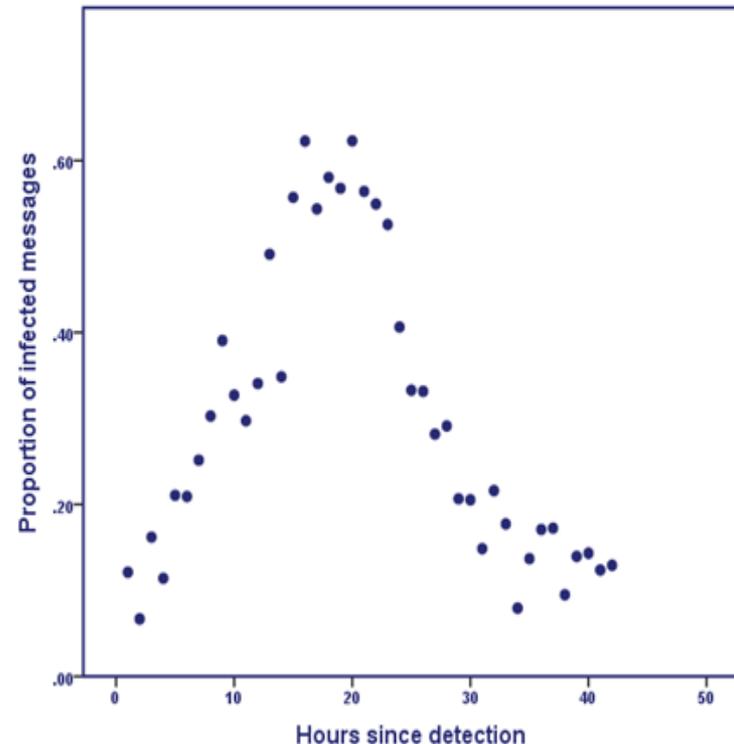


# Asymptotic regression model

$$y = \theta_1 + \theta_2 e^{\theta_3 x}$$

## Starting values

- $\theta_1$  - lowest value (0)
- $\theta_2$  - difference upper value and lowest value (0.6)
- $\theta_3$  - estimated slope between points in plot.  
(0.6-0.1/20-40)=-0.025



Analyze Direct Marketing Graphs Utilities Add-ons Window Help

Reports  
Descriptive Statistics  
Tables  
Compare Means  
General Linear Model  
Generalized Linear Models  
Mixed Models  
Correlate  
Regression  
Loglinear  
Neural Networks  
Classify  
Dimension Reduction  
Scale  
Nonparametric Tests  
Forecasting  
Survival  
Multiple Response  
PS Matching  
Missing Value Analysis...  
Multiple Imputation  
Complex Samples  
Simulation...  
Quality Control  
ROC Curve...

Utilities  
Add-ons  
Window  
Help

var var var

Automatic Linear Mo  
Linear...  
Residual Heterosce  
Regression Relative  
Curve Estimation...  
Tobit Regression...  
Robust Regression...  
Quantile Regression  
Partial Least Square  
Heckman Regressio  
Binary Logistic...  
Multinomial Logistic  
Ordinal...  
Probit...  
Nonlinear...  
Weight Estimation...

### Nonlinear Regression

Dependent: Proportion of infected messages [infected]

Model Expression:  

$$(time < 20) * b1 / (1 + b2 * \exp(-b3 * time)) + (time \geq 20) * (a1 + a2 * \exp(a3 * (time - 19)))$$

Hours since detecti...  
Proportion of infecte...

Function group:  
All  
Arithmetic  
CDF & Noncentral CDF  
Conversion  
Current Date/Time  
Date Arithmetic  
Date Creation  
Date Extraction

Functions and Special Variables:

Parameters...  
Parameters

Nonlinear Regression: Parameters

Name:

Starting Value:

Add  
Change  
Remove

b1(0.65)  
b2(5)  
b3(0.03)  
a1(0)  
a2(0.05)  
a3(-0.025)

Use starting values from previous analysis

Continue Cancel Help

Paste Reset Cancel Help

Loss...  
Constraints...  
Save...  
Options...

SPSS Statistics Data Editor window showing a dataset with columns 'time' and 'infected'. The 'infected' column contains values ranging from 0.07 to 0.62. A 'Nonlinear Regression' dialog box is open, with 'Proportion of infected messages [infected]' selected as the dependent variable. The model expression is:  $(time < 20) * b1 / (1 + b2 * \exp(-b3 * time)) + (time \geq 20) * (a1 + a2 * \exp(a3 * (time - 19)))$ . The 'Nonlinear Regression: Parameter Constraints' dialog box is also open, showing parameters b1(0.65), b2(5), b3(0.03), a1(0.0), a2(0.5), and a3(-0.025). The 'Define parameter constraint' option is selected, and constraints are listed as b1 >= 0, b2 >= 0, b3 >= 0, a1 >= 0, a2 >= 0, and a3 <= 0.

'Nonlinear Regression: Save...' dialog box with options checked for 'Predicted values', 'Residuals', and 'Loss function values'. Buttons for 'Continue', 'Cancel', and 'Help' are visible.

# Output

Parameter Estimates

Parameter	Estimate	Std. Error	95% Confidence Interval	
			Lower Bound	Upper Bound
b1	.734	.127	.477	.991
b2	7.428	1.375	4.638	10.217
b3	.184	.040	.103	.265
a1	.091	.030	.030	.153
a2	.661	.044	.572	.750
a3	-.150	.027	-.205	-.095

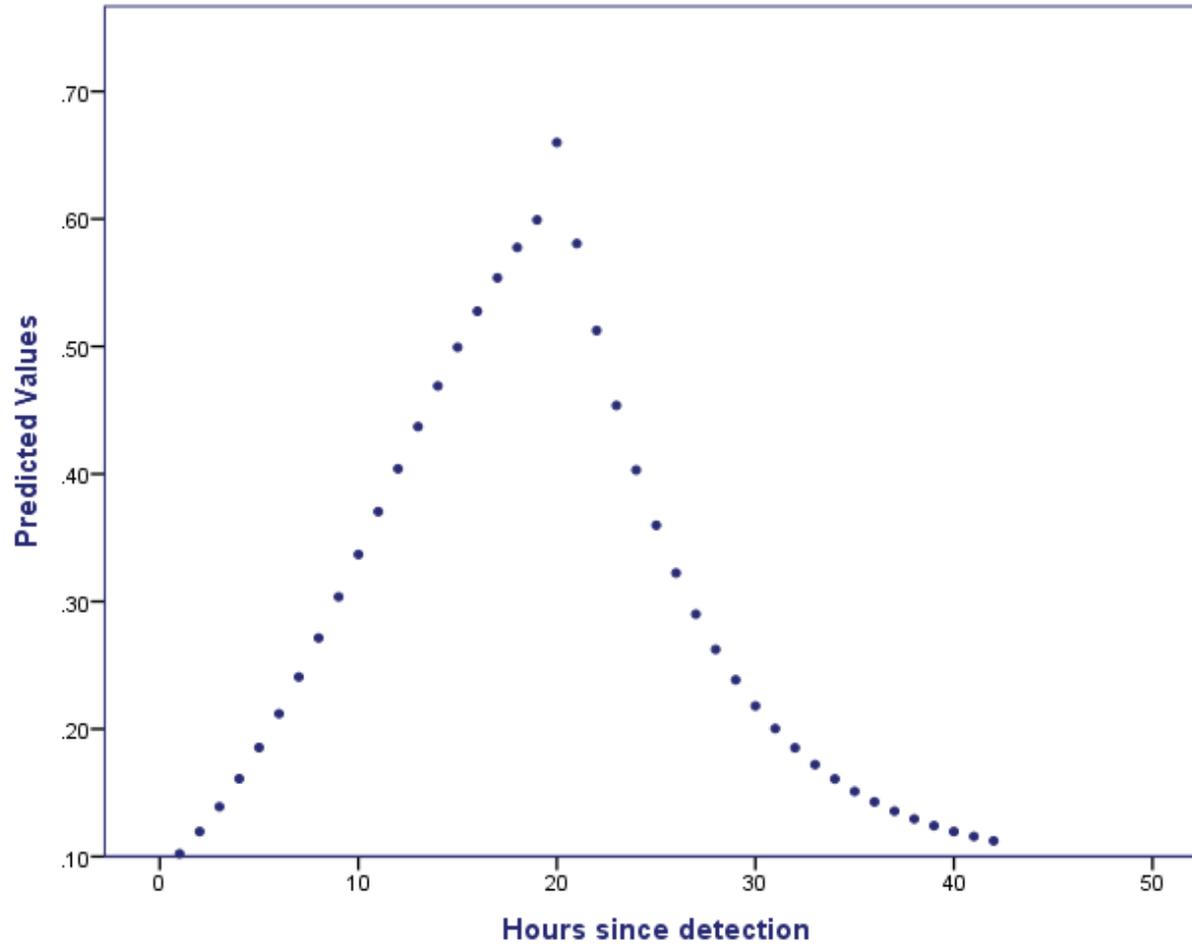
ANOVA<sup>a</sup>

Source	Sum of Squares	df	Mean Squares
Regression	4.884	6	.814
Residual	.082	36	.002
Uncorrected Total	4.966	42	
Corrected Total	1.212	41	

Dependent variable: Proportion of infected messages

a. R squared =  $1 - (\text{Residual Sum of Squares}) / (\text{Corrected Sum of Squares}) = .933$ .

# Output





Thanks for listening

Otago : University

