# DeepTriangle: A Deep Learning Approach to Loss Reserving

Kevin Kuo[a]

[a]*RStudio, 250 Northern Ave, Boston, MA 02210, United States*

**Abstract**

We propose a novel approach for loss reserving based on deep neural networks. The approach allows for joint modeling of paid losses and claims outstanding, and incorporation of heterogeneous inputs. We validate the models on loss reserving data across lines of business, and show that they improve on the predictive accuracy of existing stochastic methods. The models require minimal feature engineering and expert input, and can be automated to produce forecasts more frequently than manual workflows.

*Keywords:* loss reserving, machine learning, neural networks
*JEL:* G22

## 1. Introduction

In the loss reserving exercise for property and casualty insurers, actuaries are concerned with forecasting future payments due to claims. Accurately estimating these payments is important from the perspectives of various stakeholders in the insurance industry. For the management of the insurer, the estimates of unpaid claims inform decisions in underwriting, pricing, and strategy. For the investors, loss reserves, and transactions related to them, are essential components in the balance sheet and income statement of the insurer. And, for the regulators, accurate loss reserves are needed to appropriately understand the financial soundness of the insurer.

There can be time lags both for reporting of claims, where the insurer is not notified of a loss until long after it has occurred, and for final development of claims, where payments continue long after the loss has been reported. Also, the amounts of claims are uncertain before they have fully developed. These factors contribute to the difficulty of the loss reserving problem, for which extensive literature exists and active research is being done. We refer the reader to England and Verrall (2002) for a survey of the problem and existing techniques.

Deep learning has garnered increasing interest in recent years due to successful applications in many fields (LeCun, Bengio, and Hinton 2015) and has recently made its way into the loss reserving literature. Wüthrich (2018b) augments the traditional chain ladder method with neural networks to incorporate claims features, and Gabrielli, Richman, and Wuthrich (2018) embeds the over-dispersed Poisson (ODP) model into a neural network.
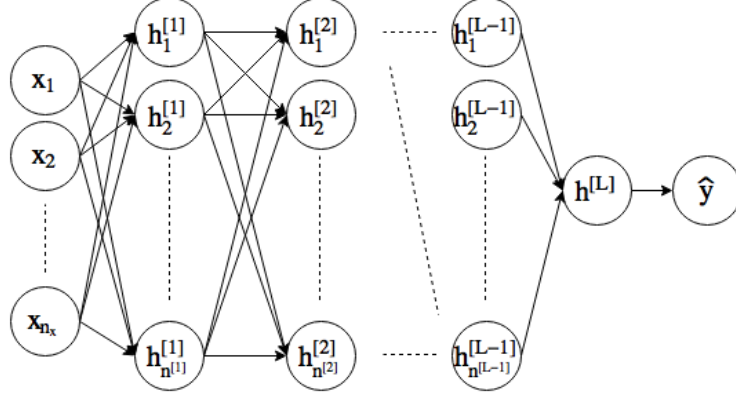
---

Figure 1: Feedforward neural network.

In developing our framework, which we call DeepTriangle, we also draw inspiration from the existing stochastic reserving literature: Quarg and Mack (2004) utilize incurred loss along with paid loss data, Miranda, Nielsen, and Verrall (2012) incorporate claim count information in addition to paid losses, and Avanzi et al. (2016) consider the dependence between lines of business within an insurer's portfolio.

The approach that we develop differs from existing works in many ways, and has the following advantages. First, it enables joint modeling of paid losses and claims outstanding for multiple companies simultaneously in a single model. In fact, the architecture can also accommodate arbitrary additional inputs, such as claim count data and economic indicators, should they be available to the modeler. Second, it requires no manual input during model updates or forecasting, which means that predictions can be generated more frequently than traditional processes, and, in turn, allows management to react to changes in the portfolio sooner.

## 2. Neural network preliminaries

For comprehensive treatments of neural network mechanics and implementation, we refer the reader to Goodfellow, Bengio, and Courville (2016) and Chollet and Allaire (2018). In order to establish common terminology used in this paper, we present a brief overview in this section.

We motivate the discussion by considering an example feedforward network with fully connected layers represented in Figure 1, where the goal is to predict an output $y$ from input $x = (x_1, x_2, \ldots, x_{n_x})$, where $n_x$ is the number of elements of $x$. The intermediate values, $h_j^{[l]}$, known as hidden units, are organized into layers, which try to transform the input data into representations that successively become more useful at predicting the output. The nodes in the figure are computed, for each layer $l = 1, \ldots, L$, as

$$h_j^{[l]} = g^{[l]}(z_j^{[l]}), \tag{1}$$

where

$$z_j^{[l]} = \sum_k w_{jk}^{[l]} h_k^{[l-1]} + b_j^{[l]}. \tag{2}$$

Here, the index $j$ spans $\{1, \ldots, n^{[l]}\}$, where $n^{[l]}$ denotes the number of units in layer $l$. The functions $g^{[l]}$ are called activation functions, whose values $h_j^{[l]}$ are known as activations. The $w^{[l]}$ values come from matrices $W^{[l]}$, with dimensions $n^{[l]} \times n^{[l-1]}$. Together with the biases $b_j^{[l]}$, they represent the weights, which are learned during training, for layer $l$.

For $l = 1$, we define the previous layer activations as the input, so that $n^{[0]} = n_x$. Hence, the calculation for the first hidden layer becomes

$$h_j^{[1]} = g^{[1]} \left( \sum_k w_{jk}^{[1]} x_k + b_j^{[1]} \right). \tag{3}$$

Also, for the output layer $l = L$, we compute the prediction

$$\hat{y} = h_j^{[L]} = g^{[L]} \left( \sum_k w_{jk}^{[L]} h_k^{[L-1]} + b_j^{[L]} \right). \tag{4}$$

We can then think of a neural network as a sequence of function compositions $f = f_L \circ f_{L-1} \circ \cdots \circ f_1$ parameterized as $f(x; W^{[1]}, b^{[1]}, \ldots, W^{[L]}, b^{[L]})$.

Each neural network model is specified with a specific loss function, which is used to measure how close the model predictions are to the actual values. During model training, the parameters discussed above are iteratively updated in order to minimize the loss function. Each update of the parameters typically involves only a subset, or mini-batch, of the training data, and one complete pass through the training data, which includes many updates, is known as an epoch. Training a neural network often requires many passes through the data.

## 3. Neural architecture for loss reserving

As shown in Figure 2, DeepTriangle is a multi-task network with two prediction goals: claims outstanding and paid loss. We construct one model for each line of business and each model is trained on data from multiple companies.

### 3.1. Training Data

Let indices $1 \leq i \leq I$ denote accident years and $1 \leq j \leq J$ denote development years under consideration. Also, let $\{P_{ij}\}$ and $\{OS_{ij}\}$ denote the *incremental* paid losses and the *total* claims outstanding, respectively.

Then, at the end of calendar year $I$, we have access to the observed data

$$\{P_{ij} : i = 1, \ldots, I; j = 1, \ldots, I - i + 1\} \tag{5}$$

and

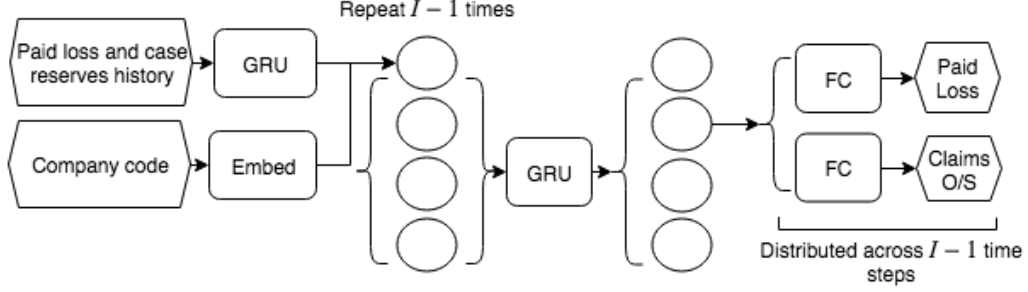$$\{OS_{ij} : i = 1, \ldots, I; j = 1, \ldots, I - i + 1\}. \tag{6}$$

Figure 2: DeepTriangle architecture. *Embed* denotes embedding layer, *GRU* denotes gated recurrent unit, *FC* denotes fully connected layer.

Assume that we are interested in development through the $I$th development year; in other words, we only forecast through the eldest maturity in the available data. The goal then is to obtain predictions for future values $\{\widehat{P}_{ij} : i = 2, \ldots, I; j = i + 1, \ldots, I\}$ and $\{\widehat{OS}_{ij} : i = 2, \ldots, I; j = i + 1, \ldots, I\}$. We can then determine ultimate losses for each accident year $i \in 1, \ldots, I$ by calculating

$$\widehat{UL}_i = \left( \sum_{j=1}^{I-i+1} P_{ij} \right) + \left( \sum_{j=I-i+2}^{I} \widehat{P}_{ij} \right). \tag{7}$$

*3.2. Response and predictor variables*

In DeepTriangle, each training sample is associated with an accident year-development year pair, which we refer to thereinafter as a *cell*. The response for the sample associated with accident year $i$ and development year $j$ is the sequence

$$(Y_{i,j}, Y_{i,j+1}, \ldots, Y_{i,I-i+1}), \tag{8}$$

where each $Y_{ij} = (P_{ij}, OS_{ij})/NPE_i$, where $NPE_i$ denotes the net earned premium for accident year $i$. Working with loss ratios makes training more tractable by normalizing values into a similar scale.

The predictor for the sample contains two components. The first component is the observed history as of the end of the calendar year associated with the cell:

$$(Y_{i,1}, Y_{i,2}, \ldots, Y_{i,j-1}). \tag{9}$$

In other words, for each accident year and at each evaluation date for which we have data, we attempt to predict future development of the accident year's paid losses and claims outstanding based on the observed history as of that date. While we are ultimately interested in $P_{ij}$, the paid losses, we include claims outstanding as an auxiliary output of the model. Since the two quantities are related, we expect to obtain better performance by jointly training than predicting each quantity independently (Collobert and Weston 2008).

The second component of the predictor is the company identifier associated with the experience. Because we include experience from multiple companies in each training

iteration, we need a way to differentiate the data from different companies. We discuss handling of the company identifier in more detail in the next section.

### 3.3. Model Architecture

DeepTriangle utilizes a sequence-to-sequence architecture inspired by Sutskever, Vinyals, and Le (2014) and Srivastava, Mansimov, and Salakhutdinov (2015).

We utilize gated recurrent units (GRU) (Chung et al. 2014), which is a type of recurrent neural network (RNN) building block that is appropriate for sequential data. A graphical representation of a GRU is shown in Figure 3, and the associated equations are as follows:

$$\tilde{h}^{<t>} = \tanh(W_h[\Gamma_r h^{<t-1>}, x^{<t>}] + b_h) \tag{10}$$

$$\Gamma_r^{<t>} = \sigma(W_r[h^{<t-1>}, x^{<t>}] + b_r) \tag{11}$$

$$\Gamma_u^{<t>} = \sigma(W_u[h^{<t-1>}, x^{<t>}] + b_u) \tag{12}$$

$$h^{<t>} = \Gamma_u^{<t>} \tilde{h}^{<t>} + (1 - \Gamma_u^{<t>})h^{<t-1>}. \tag{13}$$

Here, $h^{<t>}$ and $x^{<t>}$ represent the activation and input values, respectively, at time $t$, and $\sigma$ denotes the logistic sigmoid function defined as

$$\sigma(x) = \frac{1}{1 + \exp(-x)}. \tag{14}$$

$W_h$, $W_r$, $W_u$, $b_h$, $b_r$, and $b_u$ are the appropriately sized weight matrices and biases to be learned.

We first encode the sequential predictor with a GRU to obtain a summary of the historical values. We then repeat the output $I - 1$ times before passing them to a decoder GRU. The factor $I - 1$ is chosen here because for the $I$th accident year, we need to forecast $I - 1$ timesteps into the future. Each timestep of the decoded sequence is then concatenated with the company embedding before being passed to two subnetworks, corresponding to the two prediction outputs, of fully connected layers, each of which shares weights across the timesteps.

The company code input is first passed to an embedding layer. In this process, each company is mapped to a fixed length vector in $\mathbb{R}^k$, where $k$ is a hyperparameter. The mapping is learned during the training of the entire network instead of a separate data preprocessing step. Companies that are similar in the context of our claims forecasting problem are mapped to vectors that are close to each other in terms of Euclidean distance. Intuitively, one can think of this representation as a proxy for characteristics of the companies, such as size of book and case reserving philosophy. Categorical embedding is a common technique in deep learning that has been successfully applied to recommendation systems (Cheng et al. 2016) and retail sales prediction (Guo and Berkhahn 2016). In the actuarial science literature, Richman and Wuthrich (2018) utilize embedding layers to capture characteristics of regions in mortality forecasting, while Gabrielli, Richman, and Wuthrich (2018) apply them to lines of business factors in loss reserving.

Rectified linear unit (ReLU) (Nair and Hinton 2010), defined as

$$x \mapsto \max(0, x), \tag{15}$$

is used as the activation function for the fully connected layers, including both of the output layers.
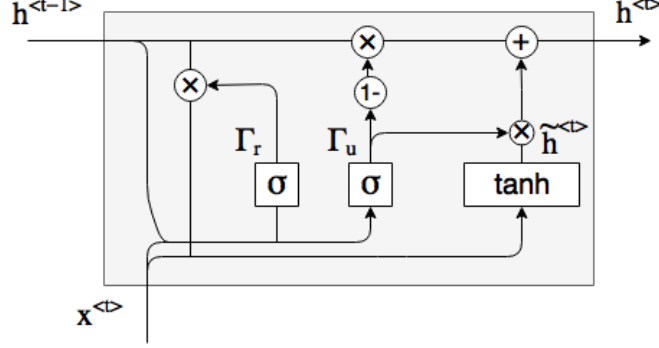
Figure 3: Gated recurrent unit.

### 3.4. Deployment considerations

While one may not have access to the latest experience data of competitors, the company code predictor can be utilized to incorporate data from companies within a group insurer. During training, the relationships among the companies are inferred based on historical development behavior. This approach provides an automated and objective alternative to manually aggregating, or clustering, the data based on knowledge of the degree of homogeneity among the companies.

If new companies join the portfolio, or if the companies and associated claims are reorganized, one would modify the embedding input size to accommodate the new codes, leaving the rest of the architecture unchanged, then refit the model. The network would then assign embedding vectors to the new companies.

Since the model outputs predictions for each triangle cell, one can calculate the traditional age-to-age, or loss development, factors (LDF) using the model forecasts. Having a familiar output may enable easier integration of DeepTriangle into existing actuarial workflows.

Insurers often have access to richer information than is available in regulatory filings, which underlies the experiments in this paper. For example, in addition to paid and incurred losses, one may include claim count triangles so that the model can also learn from, and predict, frequency information.

## 4. Experiments

### 4.1. Data

We validate the modeling approach on data from National Association of Insurance Commissioners (NAIC) Schedule P triangles (Meyers and Shi 2011). The dataset corresponds to claims from accident years 1988-1997, with development experience of 10 years for each accident year.

Following Meyers (2015), we restrict ourselves to a subset of the data which covers four lines of business (commercial auto, private personal auto, workers' compensation, and other liability) and 50 companies in each line of business. This is done to facilitate comparison to existing results.

We use the following variables from the dataset in our study: line of business, company code, accident year, development lag, incurred loss, cumulative paid loss, and net earned premium. Claims outstanding, for the purpose of this study, is derived as incurred loss less cumulative paid loss.

We use data as of year end 1997 for training, and evaluate predictive performance on the development year 10 ultimates.

## 4.2. Evaluation metrics

We aim to produce scalar metrics to evaluate the performance of the model on each line of business. To this end, for each company and each line of business, we calculate the actual and predicted ultimate losses as of development year 10, for all accident years combined, then compute the root mean squared percentage error (RMSPE) and mean absolute percentage error (MAPE) over companies in each line of business. Percentage errors are used in order to have unit-free measures for comparing across companies with vastly different sizes of portfolios. Formally, if $\mathcal{C}_l$ is the set of companies in line of business $l$,

$$MAPE_l = \frac{1}{|\mathcal{C}_l|} \sum_{C \in \mathcal{C}_l} \left| \frac{\widehat{UL}_C - UL_C}{UL_C} \right|, \tag{16}$$

and

$$RMSPE_l = \sqrt{\frac{1}{|\mathcal{C}_l|} \sum_{C \in \mathcal{C}_l} \left( \frac{\widehat{UL}_C - UL_C)}{UL_C} \right)^2} \tag{17}$$

where $\widehat{UL}_C$ and $UL_C$ are the predicted and actual cumulative ultimate losses, respectively, for company $C$.

An alternative approach for evaluation could involve weighting the company results by the associated earned premium or using dollar amounts. However, due to the distribution of company sizes in the dataset, the weights would concentrate on a handful of companies. Hence, to obtain a more balanced evaluation, we choose to report the unweighted percentage-based measures outlined above.

## 4.3. Implementation and training

The loss function for the each output is computed as the average over the forecasted time steps of the mean squared error of the predictions. The losses for the outputs are then averaged to obtain the network loss. Formally, for the sample associated with cell $(i, j)$, we can write the per-sample loss as

$$\frac{1}{I - i + 1 - (j - 1)} \sum_{k=j}^{I-i+1} \frac{(\widehat{P_{ik}} - P_{ik})^2 + (\widehat{OS_{ik}} - OS_{ik})^2}{2}. \tag{18}$$

For optimization, we use the AMSGRAD (Reddi, Kale, and Kumar 2018) variant of ADAM with a learning rate of 0.0005. We train each neural network for a maximum of 1000 epochs with the following early stopping scheme: if the loss on the validation set does not improve over a 200-epoch window, we terminate training and revert back to the

Table 1: Performance comparison of various models. DeepTriangle and AutoML are abbreviated do DT and ML, respectively.

| Line of Business | Mack | ODP | CIT | LIT | ML | **DT** |
|---|---|---|---|---|---|---|
| **MAPE** | | | | | | |
| Commercial Auto | 0.060 | 0.217 | 0.052 | 0.052 | 0.068 | **0.043** |
| Other Liability | 0.134 | 0.223 | 0.165 | 0.152 | 0.142 | **0.109** |
| Private Passenger Auto | 0.038 | 0.039 | 0.038 | 0.040 | 0.036 | **0.025** |
| Workers' Compensation | 0.053 | 0.105 | 0.054 | 0.054 | 0.067 | **0.046** |
| **RMSPE** | | | | | | |
| Commercial Auto | 0.080 | 0.822 | 0.076 | 0.074 | 0.096 | **0.057** |
| Other Liability | 0.202 | 0.477 | 0.220 | 0.209 | 0.181 | **0.150** |
| Private Passenger Auto | 0.061 | 0.063 | 0.057 | 0.060 | 0.059 | **0.039** |
| Workers' Compensation | 0.079 | 0.368 | 0.080 | 0.080 | 0.099 | **0.067** |

weights on the epoch with the lowest validation loss. The validation set used in the early stopping criterion is defined to be the subset of the training data that becomes available after calendar year 1995. For each line of business, we create an ensemble of 100 models, each trained with the same architecture but different random weight initialization. This is done to reduce the variance inherent in the randomness associated with neural networks.

We implement DeepTriangle using the keras R package (Chollet, Allaire, and others 2017) with the TensorFlow (Abadi et al. 2015) backend. Code for producing the experiment results is available online.[1]

### 4.4. Results and discussion

In Table 1 we tabulate the out-of-time performance of DeepTriangle against other models: the Mack chain-ladder model (Mack 1993), the bootstrap ODP model (England and Verrall 2002), an AutoML model, and a selection of Bayesian Markov chain Monte Carlo (MCMC) models from Meyers (2015) including the correlated incremental trend (CIT) and leveled incremental trend (LIT) models. For the stochastic models, we use the means of the predictive distributions as the point estimates to which we compare the actual outcomes. For DeepTriangle, we report the averaged predictions from the ensembles.

The AutoML model is developed by automatically searching over a set of common machine learning techniques. In the implementation we use, it trains and cross-validates a random forest, an extremely-randomized forest, a random grid of gradient boosting machines, a random grid of deep feedforward neural networks, and stacked ensembles thereof (The H2O.ai team 2018). Details of these algorithms can be found in Friedman, Hastie, and Tibshirani (2001). Because the machine learning techniques produce scalar outputs, we use an iterative forecasting scheme where the prediction for a timestep is used in the predictor for the next timestep.

We see that DeepTriangle improves on the performance of the popular chain ladder and ODP models, common machine learning models, and Bayesian stochastic models.

---

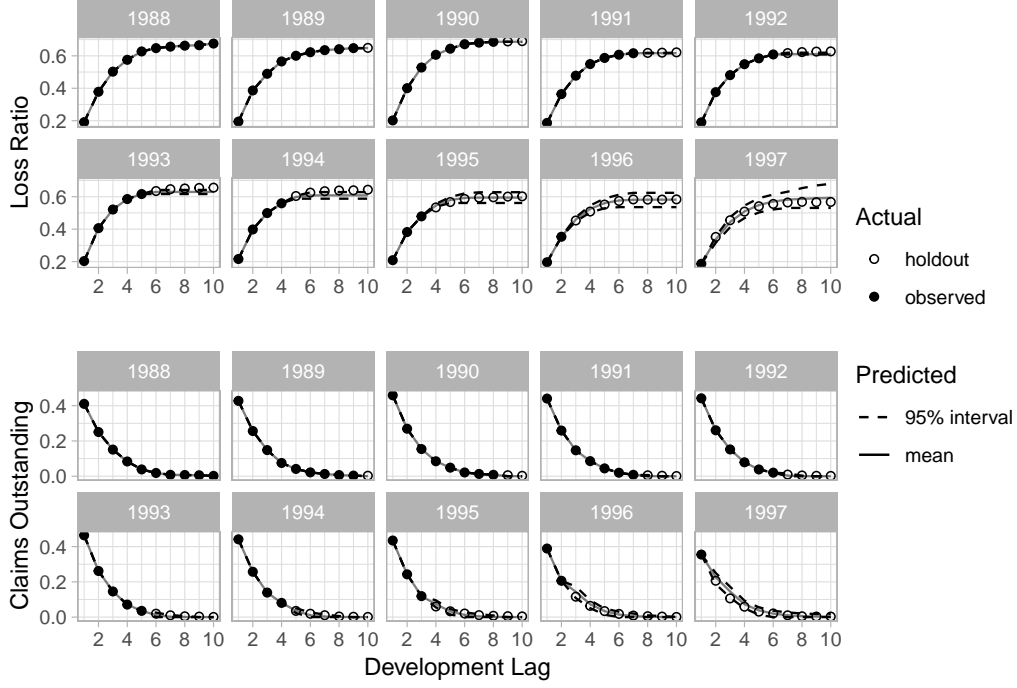[1] https://github.com/kevinykuo/deeptriangle

Figure 4: Development by accident year for Company 1767, commercial auto.

In addition to aggregated results for all companies, we also investigate qualitatively the ability of DeepTriangle to learn development patterns of individual companies. Figures 4 and 5 show the paid loss development and claims outstanding development for the commercial auto line of Company 1767 and the workers' compensation line of Company 337, respectively. We see that the model captures the development patterns for Company 1767 reasonably well. However, it is unsuccessful in forecasting the deteriorating loss ratios for Company 337's workers' compensation book.

We do not study uncertainty estimates in this paper nor interpret the forecasts as posterior predictive distributions; rather, they are included to reflect the stochastic nature of optimizing neural networks. We note that others have exploited randomness in weight initialization in producing predictive distributions (Lakshminarayanan, Pritzel, and Blundell 2017), and further research could study the applicability of these techniques to reserve variability.

## 5. Conclusion

We introduce DeepTriangle, a deep learning framework for forecasting paid losses. Our models are able to attain performance comparable, by our metrics, to modern stochastic reserving techniques without expert input. By utilizing neural networks, we can incorporate multiple heterogeneous inputs and train on multiple objectives simultaneously, and also allow customization of models based on available data.
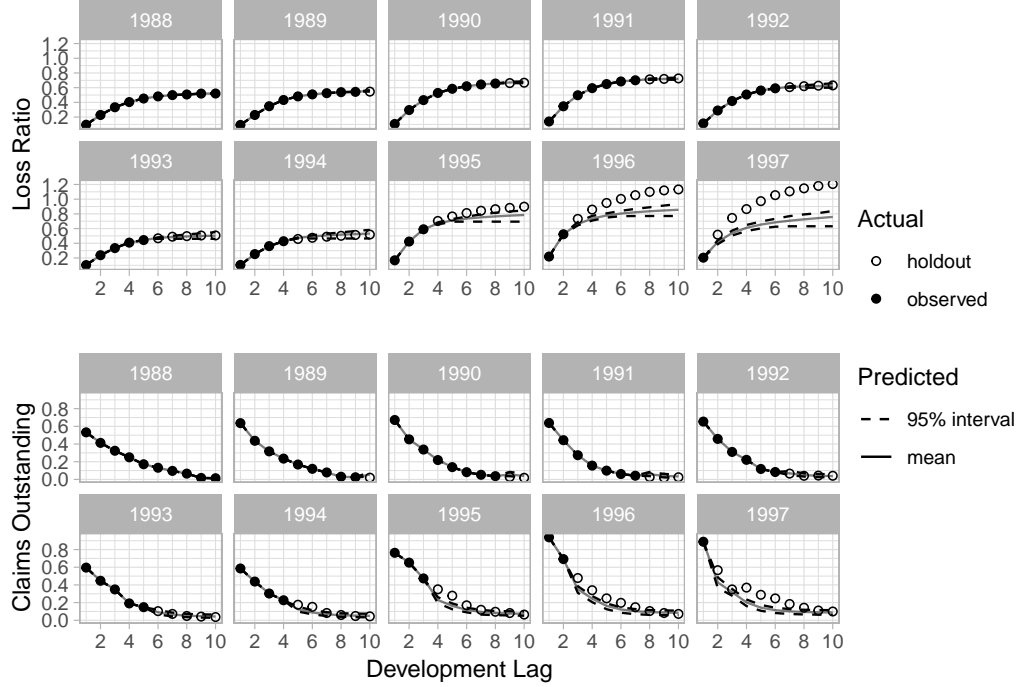
Figure 5: Development by accident year for Company 337, workers' compensation.

We analyze an aggregated dataset with limited features in this paper because it is publicly available and well studied, but one can extend DeepTriangle to incorporate additional data, such as claim counts.

Deep neural networks can be designed to extend recent efforts, such as Wüthrich (2018a), on applying machine learning to claims level reserving. They can also be designed to incorporate additional features that are not handled well by traditional machine learning algorithms, such as claims adjusters' notes from free text fields and images.

While this study focuses on prediction of point estimates, future extensions may include outputting distributions in order to address reserve variability.

## References

Abadi, Martı'n, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S. Corrado, et al. 2015. "TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems." https://www.tensorflow.org/.

Avanzi, Benjamin, Greg Taylor, Phuong Anh Vu, and Bernard Wong. 2016. "Stochastic Loss Reserving with Dependence: A Flexible Multivariate Tweedie Approach." *Insurance: Mathematics and Economics* 71. Elsevier: 63–78.

Cheng, Heng-Tze, Mustafa Ispir, Rohan Anil, Zakaria Haque, Lichan Hong, Vihan Jain, Xiaobing Liu, et al. 2016. "Wide & Deep Learning for Recommender Systems." *Proceedings of the 1st Workshop on Deep Learning for Recommender Systems - DLRS 2016.* ACM Press. https://doi.org/10.1145/2988450.2988454.

Chollet, Francois, and JJ Allaire. 2018. *Deep Learning with R.* Manning Publications.

Chollet, François, JJ Allaire, and others. 2017. "R Interface to Keras." https://github.com/rstudio/keras; GitHub.

Chung, Junyoung, Caglar Gulcehre, KyungHyun Cho, and Yoshua Bengio. 2014. "Empirical Evaluation of Gated Recurrent Neural Networks on Sequence Modeling."

Collobert, Ronan, and Jason Weston. 2008. "A Unified Architecture for Natural Language Processing: Deep Neural Networks with Multitask Learning." In *Proceedings of the 25th International Conference on Machine Learning*, 160–67. ACM.

England, Peter D, and Richard J Verrall. 2002. "Stochastic Claims Reserving in General Insurance." *British Actuarial Journal* 8 (3). Cambridge University Press: 443–518.

Friedman, Jerome, Trevor Hastie, and Robert Tibshirani. 2001. *The Elements of Statistical Learning.* Vol. 1. 10. Springer series in statistics New York, NY, USA:

Gabrielli, Andrea, Ronald Richman, and Mario V Wuthrich. 2018. "Neural Network Embedding of the over-Dispersed Poisson Reserving Model." *Available at SSRN.*

Goodfellow, Ian, Yoshua Bengio, and Aaron Courville. 2016. *Deep Learning.* MIT Press Cambridge.

Guo, Cheng, and Felix Berkhahn. 2016. "Entity Embeddings of Categorical Variables." *CoRR* abs/1604.06737. http://arxiv.org/abs/1604.06737.

Lakshminarayanan, Balaji, Alexander Pritzel, and Charles Blundell. 2017. "Simple and Scalable Predictive Uncertainty Estimation Using Deep Ensembles." In *Advances in Neural Information Processing Systems 30*, edited by I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, 6402–13. Curran Associates, Inc. http://papers.nips.cc/paper/7219-simple-and-scalable-predictive-uncertainty-estimation-using-deep-ensembles.pdf.

LeCun, Yann, Yoshua Bengio, and Geoffrey Hinton. 2015. "Deep Learning." *Nature* 521 (7553). Nature Publishing Group: 436.

Mack, Thomas. 1993. "Distribution-Free Calculation of the Standard Error of Chain Ladder Reserve Estimates." *Astin Bulletin* 23 (2): 213–25.

Meyers, Glenn. 2015. *Stochastic Loss Reserving Using Bayesian MCMC Models.* Casualty Actuarial Society.

Meyers, Glenn, and Peng Shi. 2011. "Loss Reserving Data Pulled from NAIC Schedule P." http://www.casact.org/research/index.cfm?fa=loss_reserves_data.

Miranda, Marı'a Dolores Martı'nez, Jens Perch Nielsen, and Richard Verrall. 2012. "Double Chain Ladder." *ASTIN Bulletin: The Journal of the IAA* 42 (1). Cambridge

University Press: 59–76.

Nair, Vinod, and Geoffrey E Hinton. 2010. "Rectified Linear Units Improve Restricted Boltzmann Machines." In *Proceedings of the 27th International Conference on Machine Learning (Icml-10)*, 807–14.

Quarg, Gerhard, and Thomas Mack. 2004. "Munich Chain Ladder." *Blätter Der DGVFM* 26 (4). Springer: 597–630.

Reddi, Sashank J., Satyen Kale, and Sanjiv Kumar. 2018. "On the Convergence of Adam and Beyond." In *International Conference on Learning Representations*. https://openreview.net/forum?id=ryQu7f-RZ.

Richman, Ronald, and Mario V Wuthrich. 2018. "A Neural Network Extension of the Lee-Carter Model to Multiple Populations." *Available at SSRN 3270877*.

Srivastava, Nitish, Elman Mansimov, and Ruslan Salakhutdinov. 2015. "Unsupervised Learning of Video Representations Using Lstms." *CoRR* abs/1502.04681. http://arxiv.org/abs/1502.04681.

Sutskever, Ilya, Oriol Vinyals, and Quoc V Le. 2014. "Sequence to Sequence Learning with Neural Networks." In *Advances in Neural Information Processing Systems*, 3104–12.

The H2O.ai team. 2018. *H2o: R Interface for H2o*. http://www.h2o.ai.

Wüthrich, Mario V. 2018a. "Machine Learning in Individual Claims Reserving." *Scandinavian Actuarial Journal*. Taylor & Francis, 1–16.

———. 2018b. "Neural Networks Applied to Chain–Ladder Reserving." *European Actuarial Journal* 8 (2). Springer: 407–36.