

Article

DeepTriangle: A Deep Learning Approach to Loss Reserving

Kevin Kuo^{1,*} ¹ Kasa AI; kevin@kasa.ai* Correspondence: kevin@kasa.ai

Version September 16, 2019 submitted to Risks



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

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 et al. 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, [Gabrielli and Wüthrich \(2018\)](#) utilize neural networks to synthesize claims data, and [Gabrielli et al. \(2018\)](#) and [Gabrielli \(2019\)](#) embed classical parametric loss reserving models into neural networks. More specifically, the development in [Gabrielli et al. \(2018\)](#) and [Gabrielli \(2019\)](#) proposes initializing a neural network so that, before training, it corresponds exactly to a classical model, such as the over-dispersed Poisson model. The training iterations then adjust the weights of the neural network to minimize the prediction errors, which can be interpreted as a boosting procedure.

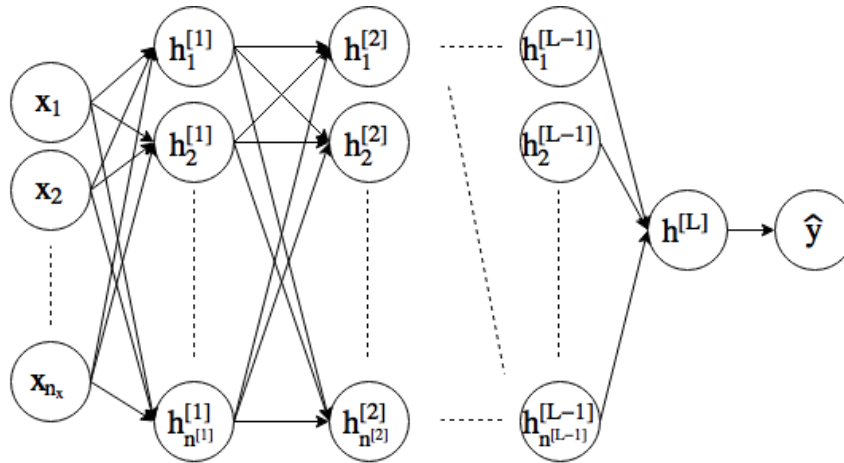


Figure 1. Feedforward neural network.

31 In developing our framework, which we call DeepTriangle¹, we also draw inspiration from
 32 the existing stochastic reserving literature. Works that propose utilizing data in addition to paid
 33 losses include [Quarg and Mack \(2004\)](#), which uses incurred losses, and [Miranda et al. \(2012\)](#), which
 34 incorporates claim count information. Moving beyond a single homogeneous portfolio, [Avanzi
 35 et al. \(2016\)](#) considers the dependencies among lines of business within an insurer's portfolio, while
 36 [Peremans et al. \(2018\)](#) proposes a robust general multivariate chain ladder approach to accommodate
 37 outliers. There is also a category of models, referred to as state space or adaptive models, that allow
 38 parameters to evolve recursively in time as more data is observed ([Chukhrova and Johannssen 2017](#)).
 39 This iterative updating mechanism is similar in spirit to the continuous updating of neural network
 40 weights during model deployment.

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

48 The rest of the paper is organized as follows: Section 2 provides a brief overview of neural
 49 network terminology, Section 3 discusses the dataset used and introduces the proposed neural network
 50 architecture, Section 4 defines the performance metrics we use to benchmark our models and discuss
 51 the results, and Section 5 concludes.

52 2. Neural Network Preliminaries

53 For comprehensive treatments of neural network mechanics and implementation, we refer the
 54 reader to [Goodfellow et al. \(2016\)](#) and [Chollet and Allaire \(2018\)](#). A more actuarially oriented discussion
 55 can be found in [Wuthrich and Buser \(2019\)](#). In order to establish common terminology used in this
 56 paper, we present a brief overview in this section.

57 We motivate the discussion by considering an example feedforward network with fully connected
 58 layers represented in Figure 1, where the goal is to predict an output y from input x . The intermediate
 59 values, known as hidden layers and represented by $h_j^{[l]}$, try to transform the input data into

¹ A portmanteau of *deep learning* and *loss development triangle*.

representations that successively become more useful at predicting the output. The nodes in the figure are computed, for each layer $l = 1, \dots, L$, as

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

where

$$z_j^{[l]} = w_j^{[l]T} h^{[l-1]} + b_j^{[l]}, \quad (2)$$

for $j = 1, \dots, n^{[l]}$. In these equations, a superscript $[l]$ denotes association with the layer l , a subscript j denotes association with the j -th component of the layer, of which there are $n^{[l]}$. The $g^{[l]}$ ($l = 1, \dots, L$) are called activation functions, whose values $h^{[l]}$ are known as activations. The vectors $w_j^{[l]}$ and scalars $b_j^{[l]}$ are known as weights and biases, respectively, and together represent the parameters of the neural network, which are learned during training.

For $l = 1$, we define the previous layer activations as the input, so that the calculation for the first hidden layer becomes

$$h_j^{[1]} = g^{[1]}(w_j^{[1]T} x + b_j^{[1]}). \quad (3)$$

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

$$\hat{y} = h_j^{[L]} = g^{[L]}(w_j^{[L]T} h^{[L-1]} + b_j^{[L]}). \quad (4)$$

We can then think of a neural network as a sequence of function compositions $f = f_L \circ f_{L-1} \circ \dots \circ f_1$ parameterized as $f(x; W^{[1]}, b^{[1]}, \dots, W^{[L]}, b^{[L]})$. Here, it should be mentioned that the $g^{[l]}$ ($l = 1, \dots, L$) are chosen to be nonlinear, except for possibly in the output layer. These nonlinearities are key to the success of neural networks, because otherwise we would have a trivial composition of linear models.

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. Data and Model Architecture

In this section, we discuss the dataset used for our experiments and the proposed model architecture.

3.1. Data Source

We use the 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. In Schedule P data, the data is aggregated into accident year-development year records. The procedure for constructing the dataset is detailed in Meyers (2015).

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. The company code is a categorical variable that denotes which insurer the records are associated with.

95 3.2. Training/Testing Setup

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

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

$$\{P_{i,j} : i = 1, \dots, I; j = 1, \dots, I - i + 1\} \quad (5)$$

100 and

$$\{OS_{i,j} : i = 1, \dots, I; j = 1, \dots, I - i + 1\}. \quad (6)$$

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

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

105 In our case, data as of year end 1997 is used for training. We then evaluate predictive performance
 106 on the development year 10 cumulative paid losses.

107 3.3. Response and Predictor Variables

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

$$(Y_{i,j}, Y_{i,j+1}, \dots, Y_{i,I-i+1}), \quad (8)$$

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

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

$$(Y_{i,1}, Y_{i,2}, \dots, Y_{i,j-1}). \quad (9)$$

116 In other words, for each accident year and at each evaluation date for which we have data, we attempt
 117 to predict future development of the accident year's paid losses and claims outstanding based on the
 118 observed history as of that date. While we are ultimately interested in $P_{i,j}$, the paid losses, we include
 119 claims outstanding as an auxiliary output of the model. We elaborate on the reasoning behind this
 120 approach in the next section.

121 The second component of the predictor is the company identifier associated with the experience.
 122 Because we include experience from multiple companies in each training iteration, we need a way
 123 to differentiate the data from different companies. We discuss handling of the company identifier in
 124 more detail in the next section.

125 3.4. Model Architecture

126 As shown in Figure 2, DeepTriangle is a multi-task network (Caruana 1997) utilizing a
 127 sequence-to-sequence architecture (Srivastava et al. 2015; Sutskever et al. 2014) with two prediction

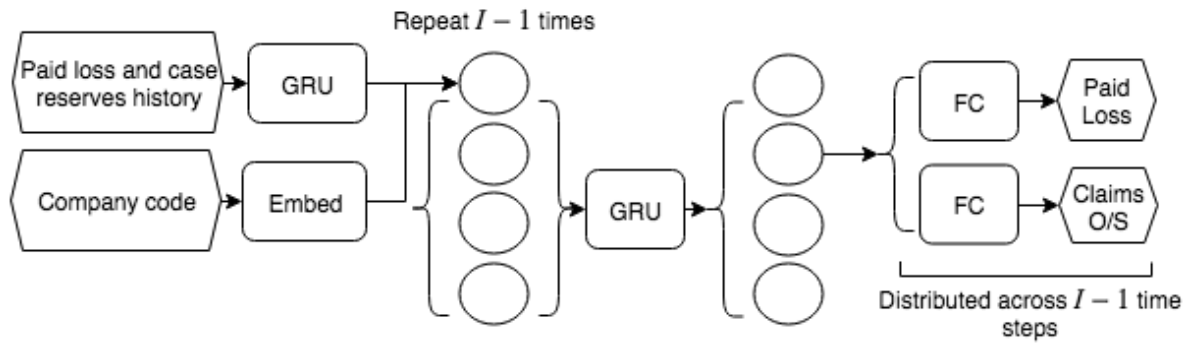


Figure 2. DeepTriangle architecture. *Embed* denotes embedding layer, *GRU* denotes gated recurrent unit, *FC* denotes fully connected layer.

128 goals: paid loss and claims outstanding. We construct one model for each line of business and each
 129 model is trained on data from multiple companies.

130 3.4.1. Multi-Task Learning

131 Since the two target quantities, paid loss and claims outstanding, are related, we expect to obtain
 132 better performance by jointly training than predicting each quantity independently. While [Caruana](#)
 133 [\(1997\)](#) contains detailed discourse on the specific mechanisms of multi-task learning, we provide some
 134 heuristics on why it may improve predictions: by utilizing the response data for claims outstanding, we
 135 are effectively increasing the training data size since we are providing more signals to the learning
 136 algorithm; there may be hidden features, useful for predicting paid losses, that are more easily learned
 137 by trying to predict claims outstanding; also, by trying to predict claims outstanding during training,
 138 we are imposing a bias towards neural network weight configurations which perform that task well,
 139 which lessens the likelihood of arriving at a model that overfits to random noise.

140 3.4.2. Sequential Input Processing

141 For handling the time series of paid losses and claims outstanding, we utilize gated recurrent
 142 units (GRU) ([Chung et al. 2014](#)), which is a type of recurrent neural network (RNN) building block
 143 that is appropriate for sequential data. A graphical representation of a GRU is shown in [Figure 3](#), and
 144 the associated equations are as follows²:

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

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

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

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

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

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

147 W_h , W_r , W_u , b_h , b_r , and b_u are the appropriately sized weight matrices and biases to be learned.
 148 Intuitively, the activations $h^{<t>}$ provide a way for the network to maintain state and “remember”

² Note the use of angle brackets to index position in a sequence rather than layers in a feedforward neural network as in [Section 2](#).

149 values from early values of the input sequence. The values $\tilde{h}^{<t>}$ can be thought of as candidates
 150 to replace the current state, and $\Gamma_u^{<t>}$ determines the weighting between the previous state and the
 151 candidate state. We remark that although the GRU (and RNN in general) may seem opaque at first,
 152 they contain sequential instructions for updating weights just like vanilla feedforward neural networks
 153 (and can in fact be interpreted as such (Goodfellow et al. 2016)).

154 We first encode the sequential predictor with a GRU to obtain a summary encoding of the historical
 155 values. We then repeat the output $I - 1$ times before passing them to a decoder GRU that outputs
 156 its hidden state for each time step. The factor $I - 1$ is chosen here because for the I th accident year,
 157 we need to forecast $I - 1$ timesteps into the future. For both the encoder and decoder GRU modules,
 158 we use 128 hidden units and a dropout rate of 0.2. Here, dropout refers to the regime where, during
 159 training, at each iteration, we randomly set the output of the hidden units to zero with a specified
 160 probability, in order to reduce overfitting (Srivastava et al. 2014). Intuitively, dropout accomplishes
 161 this by approximating an ensemble of sub-networks that can be constructed by removing some hidden
 162 units.

163 3.4.3. Company Code Embeddings

164 The company code input is first passed to an embedding layer. In this process, each company
 165 is mapped to a fixed length vector in \mathbb{R}^k , where k is a hyperparameter. In our case, we choose
 166 $k = \text{number of levels} - 1 = 49$, as recommended in Guo and Berkhahn (2016). In other words, each
 167 company is represented by a vector in \mathbb{R}^{49} . This mapping mechanism is part of the neural network
 168 and hence is learned during the training of the network, instead of in a separate data preprocessing
 169 step, so the learned numerical representations are optimized for predicted the future paid losses.
 170 Companies that are similar in the context of our claims forecasting problem are mapped to vectors that
 171 are close to each other in terms of Euclidean distance. Intuitively, one can think of this representation
 172 as a proxy for characteristics of the companies, such as size of book and case reserving philosophy.
 173 Categorical embedding is a common technique in deep learning that has been successfully applied
 174 to recommendation systems (Cheng et al. 2016) and retail sales prediction (Guo and Berkhahn 2016).
 175 In the actuarial science literature, Richman and Wuthrich (2018) utilize embedding layers to capture
 176 characteristics of regions in mortality forecasting, while Gabrielli et al. (2018) apply them to lines of
 177 business factors in loss reserving.

178 3.4.4. Fully Connected Layers and Outputs

179 Each timestep of the decoded sequence from the GRU decoder is then concatenated with the
 180 company embedding output. The concatenated values are then passed to two subnetworks of fully
 181 connected layers, each of which shares weights across the timesteps. The two subnetworks correspond
 182 to the paid loss and case outstanding predictions, respectively, and each consists of a hidden layer of
 183 64 units with a dropout rate of 0.2, followed by an output layer of 1 unit to represent the paid loss or
 184 claims outstanding at a time step.

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

$$x \mapsto \max(0, x), \quad (15)$$

186 is used as the activation function (which we denote by g in Section 2) for all fully connected layers,
 187 including both of the output layers. We remark that this choice of output activation implies we only
 188 predict nonnegative cash flows, i.e. no recoveries. This assumption is reasonable for the dataset we use
 189 in our experiments, but may be modified to accommodate other use cases.

190 3.5. Deployment Considerations

191 While one may not have access to the latest experience data of competitors, the company code
 192 predictor can be utilized to incorporate data from companies within a group insurer. During training,

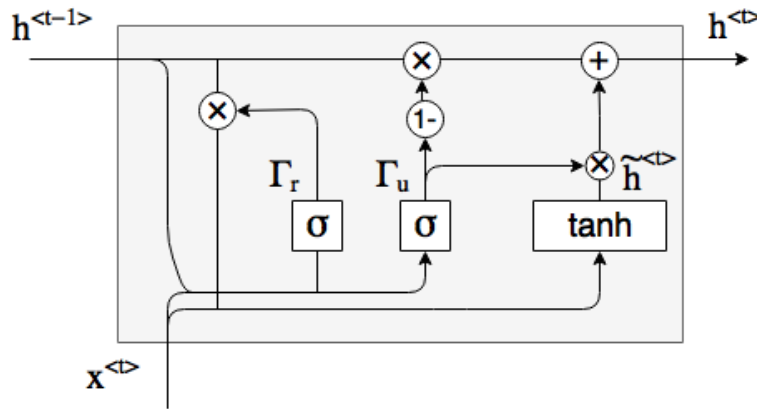


Figure 3. Gated recurrent unit.

193 the relationships among the companies are inferred based on historical development behavior. This
 194 approach provides an automated and objective alternative to manually aggregating, or clustering, the
 195 data based on knowledge of the degree of homogeneity among the companies.

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

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

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

206 4. Experiments

207 We now describe the performance metrics for benchmarking the models and training details, then
 208 discuss the results.

209 4.1. Evaluation Metrics

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

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

217 and

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

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

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

220 An alternative approach for evaluation could involve weighting the company results by the
 221 associated earned premium or using dollar amounts. However, due to the distribution of company
 222 sizes in the dataset, the weights would concentrate on a handful of companies. Hence, to obtain a
 223 more balanced evaluation, we choose to report the unweighted percentage-based measures outlined
 224 above. We note that the evaluation of reserving models is an ongoing area of research; and refer the
 225 reader to [Martinek \(2019\)](#) for a recent analysis.

226 4.2. Implementation and Training

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

$$\frac{1}{I - i + 1 - (j - 1)} \sum_{k=j}^{I-i+1} \frac{(\widehat{P}_{i,k} - P_{i,k})^2 + (\widehat{OS}_{i,k} - OS_{i,k})^2}{2}. \quad (18)$$

230 For optimization, we use the AMSGRAD ([Reddi et al. 2018](#)) variant of ADAM with a learning
 231 rate of 0.0005. We train each neural network for a maximum of 1000 epochs with the following early
 232 stopping scheme: if the loss on the validation set does not improve over a 200-epoch window, we
 233 terminate training and revert back to the weights on the epoch with the lowest validation loss. The
 234 validation set used in the early stopping criterion is defined to be the subset of the training data that
 235 becomes available after calendar year 1995. For each line of business, we create an ensemble of 100
 236 models, each trained with the same architecture but different random weight initialization. This is
 237 done to reduce the variance inherent in the randomness associated with neural networks.

238 We implement DeepTriangle using the keras R package ([Chollet et al. 2017](#)) and TensorFlow
 239 ([Abadi et al. 2015](#)), which are open source software for developing neural network models. Code for
 240 producing the experiment results is available online.³

241 4.3. Results and Discussion

242 In Table 1 we tabulate the out-of-time performance of DeepTriangle against other models: the
 243 Mack chain-ladder model ([Mack 1993](#)), the bootstrap ODP model ([England and Verrall 2002](#)), an
 244 AutoML model, and a selection of Bayesian Markov chain Monte Carlo (MCMC) models from [Meyers](#)

³ <https://github.com/kasaai/deeptriangle>.

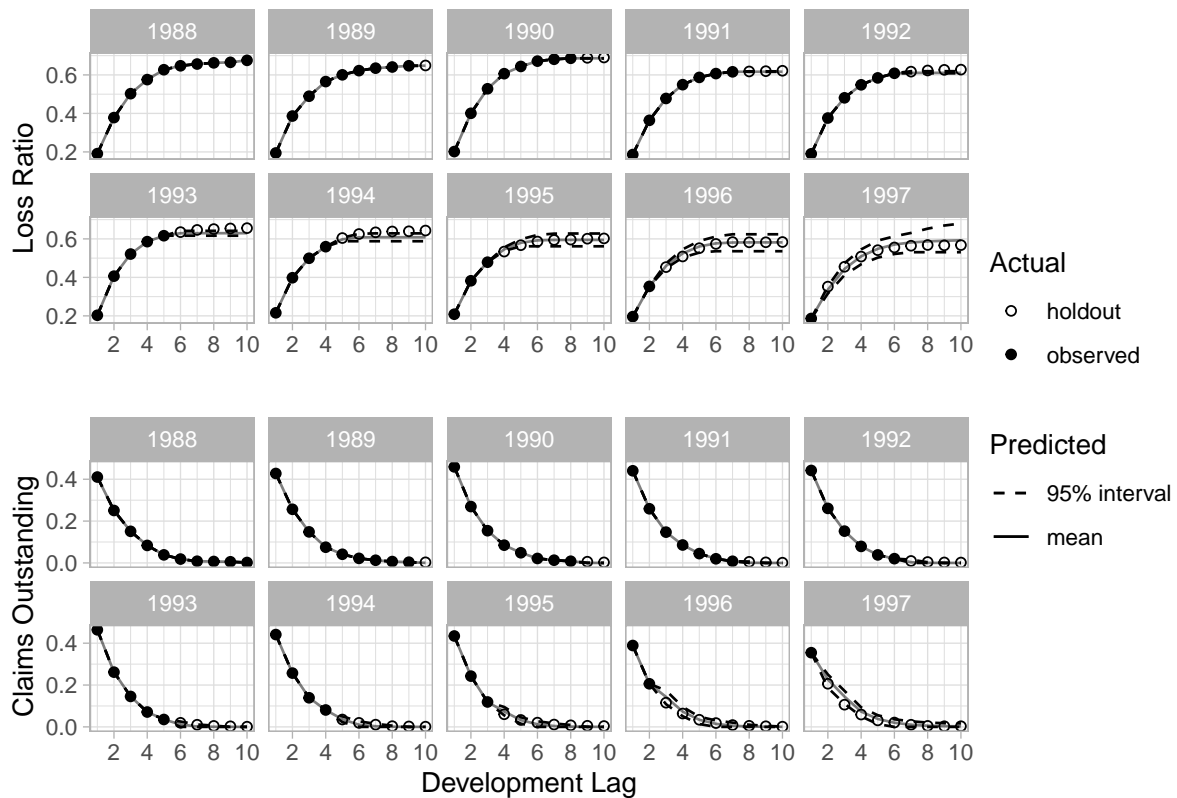


Figure 4. Development by accident year for Company 1767, commercial auto.

245 (2015) including the correlated incremental trend (CIT) and leveled incremental trend (LIT) models.
 246 For the stochastic models, we use the means of the predictive distributions as the point estimates to
 247 which we compare the actual outcomes. For DeepTriangle, we report the averaged predictions from
 248 the ensembles.

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

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

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

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

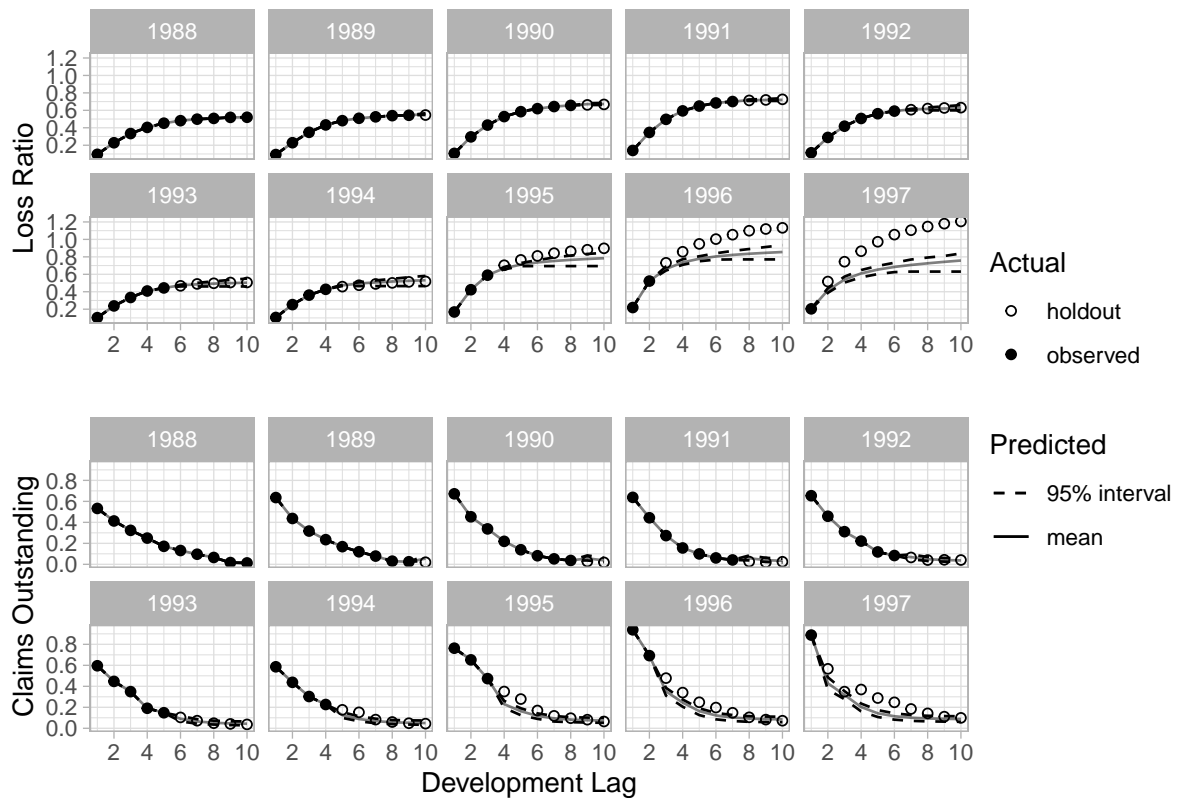


Figure 5. Development by accident year for Company 337, workers' compensation.

269 5. Conclusion

270 We introduce DeepTriangle, a deep learning framework for forecasting paid losses. Our models
 271 are able to attain performance comparable, by our metrics, to modern stochastic reserving techniques,
 272 without expert input. This means that one can automate model updating and report production
 273 at the desired frequency (although we note that, as with any automated machine learning system,
 274 a process involving expert review should be implemented). By utilizing neural networks, we can
 275 incorporate multiple heterogeneous inputs and train on multiple objectives simultaneously, and also
 276 allow customization of models based on available data. To summarize, this framework maintains
 277 accuracy while providing automatability and extensibility.

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

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

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

287 **Acknowledgments:** We thank Sigrid Keydana, Ronald Richman, the anonymous reviewers, and the volunteers
 288 on the Casualty Actuarial Society Committee on Reserves (CASCOR) who helped to improve the paper through
 289 helpful comments and discussions.

290 **Conflicts of Interest:** The author declares no conflict of interest.

291 **References**

- 292 Abadi, Martín, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S. Corrado,
293 Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfellow, Andrew Harp, Geoffrey Irving,
294 Michael Isard, Yangqing Jia, Rafal Jozefowicz, Lukasz Kaiser, Manjunath Kudlur, Josh Levenberg, Dandelion
295 Mané, Rajat Monga, Sherry Moore, Derek Murray, Chris Olah, Mike Schuster, Jonathon Shlens, Benoit Steiner,
296 Ilya Sutskever, Kunal Talwar, Paul Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda Viégas, Oriol Vinyals,
297 Pete Warden, Martin Wattenberg, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. 2015. TensorFlow: Large-scale
298 machine learning on heterogeneous systems. Software available from tensorflow.org.
- 299 Avanzi, Benjamin, Greg Taylor, Phuong Anh Vu, and Bernard Wong. 2016. Stochastic loss reserving with
300 dependence: A flexible multivariate tweedie approach. *Insurance: Mathematics and Economics* 71, 63–78.
- 301 Caruana, Rich. 1997. Multitask learning. *Machine learning* 28(1), 41–75.
- 302 Cheng, Heng-Tze, Mustafa Ispir, Rohan Anil, Zakaria Haque, Lichan Hong, Vihan Jain, Xiaobing Liu, Hemal Shah,
303 Levent Koc, Jeremiah Harmsen, and et al.. 2016. Wide & deep learning for recommender systems. *Proceedings*
304 *of the 1st Workshop on Deep Learning for Recommender Systems - DLRS 2016*. doi:10.1145/2988450.2988454.
- 305 Chollet, Francois and JJ Allaire. 2018. *Deep learning with R*. Manning Publications.
- 306 Chollet, François, JJ Allaire, et al.. 2017. R interface to keras. <https://github.com/rstudio/keras>.
- 307 Chukhrova, Nataliya and Arne Johannssen. 2017. State space models and the kalman-filter in stochastic claims
308 reserving: forecasting, filtering and smoothing. *Risks* 5(2), 30.
- 309 Chung, Junyoung, Caglar Gulcehre, KyungHyun Cho, and Yoshua Bengio. 2014. Empirical evaluation of gated
310 recurrent neural networks on sequence modeling.
- 311 England, Peter D and Richard J Verrall. 2002. Stochastic claims reserving in general insurance. *British Actuarial*
312 *Journal* 8(3), 443–518.
- 313 Friedman, Jerome, Trevor Hastie, and Robert Tibshirani. 2001. *The elements of statistical learning*, Volume 1.
314 Springer series in statistics New York, NY, USA:.
- 315 Gabrielli, Andrea. 2019. A neural network boosted double over-dispersed poisson claims reserving model.
316 Available at SSRN 3365517.
- 317 Gabrielli, Andrea, Ronald Richman, and Mario V Wuthrich. 2018. Neural network embedding of the
318 over-dispersed poisson reserving model. Available at SSRN.
- 319 Gabrielli, Andrea and Mario V Wüthrich. 2018. An individual claims history simulation machine. *Risks* 6(2), 29.
- 320 Goodfellow, Ian, Yoshua Bengio, and Aaron Courville. 2016. *Deep learning*. MIT Press Cambridge.
- 321 Guo, Cheng and Felix Berkhahn. 2016. Entity embeddings of categorical variables. *CoRR abs/1604.06737*.
- 322 Lakshminarayanan, Balaji, Alexander Pritzel, and Charles Blundell. 2017. Simple and scalable predictive
323 uncertainty estimation using deep ensembles. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus,
324 S. Vishwanathan, and R. Garnett (Eds.), *Advances in Neural Information Processing Systems* 30, pp. 6402–6413.
325 Curran Associates, Inc.
- 326 LeCun, Yann, Yoshua Bengio, and Geoffrey Hinton. 2015. Deep learning. *nature* 521(7553), 436.
- 327 Mack, Thomas. 1993. Distribution-free calculation of the standard error of chain ladder reserve estimates. *Astin*
328 *bulletin* 23(2), 213–225.
- 329 Martinek, László. 2019. Analysis of stochastic reserving models by means of naic claims data. *Risks* 7(2), 62.
- 330 Meyers, Glenn. 2015. *Stochastic loss reserving using Bayesian MCMC models*. Casualty Actuarial Society.
- 331 Meyers, Glenn and Peng Shi. 2011. Loss reserving data pulled from NAIC schedule p. [http://www.casact.org/
332 research/index.cfm?fa=loss_reserves_data](http://www.casact.org/research/index.cfm?fa=loss_reserves_data).
- 333 Miranda, María Dolores Martínez, Jens Perch Nielsen, and Richard Verrall. 2012. Double chain ladder. *ASTIN*
334 *Bulletin: The Journal of the IAA* 42(1), 59–76.
- 335 Nair, Vinod and Geoffrey E Hinton. 2010. Rectified linear units improve restricted boltzmann machines. In
336 *Proceedings of the 27th international conference on machine learning (ICML-10)*, pp. 807–814.
- 337 Peremans, Kris, Stefan Van Aelst, and Tim Verdonck. 2018. A robust general multivariate chain ladder method.
338 *Risks* 6(4), 108.
- 339 Quarg, Gerhard and Thomas Mack. 2004. Munich chain ladder. *Blätter der DGVFM* 26(4), 597–630.
- 340 Reddi, Sashank J., Satyen Kale, and Sanjiv Kumar. 2018. On the convergence of adam and beyond. In *International*
341 *Conference on Learning Representations*.

- 342 Richman, Ronald and Mario V Wüthrich. 2018. A neural network extension of the lee-carter model to multiple
343 populations. *Available at SSRN 3270877*.
- 344 Srivastava, Nitish, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. 2014. Dropout: a
345 simple way to prevent neural networks from overfitting. *The journal of machine learning research* 15(1), 1929–1958.
- 346 Srivastava, Nitish, Elman Mansimov, and Ruslan Salakhutdinov. 2015. Unsupervised learning of video
347 representations using lstms. *CoRR abs/1502.04681*.
- 348 Sutskever, Ilya, Oriol Vinyals, and Quoc V Le. 2014. Sequence to sequence learning with neural networks. In
349 *Advances in neural information processing systems*, pp. 3104–3112.
- 350 The H2O.ai team. 2018. *h2o: R Interface for H2O*. R package version 3.20.0.8.
- 351 Wüthrich, Mario V. 2018a. Machine learning in individual claims reserving. *Scandinavian Actuarial Journal*, 1–16.
- 352 Wüthrich, Mario V. 2018b. Neural networks applied to chain–ladder reserving. *European Actuarial Journal* 8(2),
353 407–436.
- 354 Wüthrich, Mario V and Christoph Buser. 2019. Data analytics for non-life insurance pricing. *Swiss Finance Institute*
355 *Research Paper* (16-68).
- 356 © 2019 by the authors. Submitted to *Risks* for possible open access publication under the terms and conditions
357 of the Creative Commons Attribution (CC BY) license (<http://creativecommons.org/licenses/by/4.0/>).