



# Article An Individual Claims History Simulation Machine

# Andrea Gabrielli \* and Mario V. Wüthrich

RiskLab, Department of Mathematics, ETH Zürich, Rämistrasse 101, 8092 Zürich, Switzerland; mario.wuethrich@math.ethz.ch

\* Correspondence: andrea.gabrielli@math.ethz.ch

Received: 5 March 2018; Accepted: 27 March 2018; Published: 30 March 2018



**Abstract:** The aim of this project is to develop a stochastic simulation machine that generates individual claims histories of non-life insurance claims. This simulation machine is based on neural networks to incorporate individual claims feature information. We provide a fully calibrated stochastic scenario generator that is based on real non-life insurance data. This stochastic simulation machine allows everyone to simulate their own synthetic insurance portfolio of individual claims histories and back-test thier preferred claims reserving method.

**Keywords:** claims reserving; individual claims; claims cash flows; micro-level stochastic reserving; loss reserving; claims simulation; neural network reserving; individual claims features; individual claims covariates; chain-ladder

# 1. Introduction

The aim of this project is to develop a stochastic simulation machine that generates individual claims histories of non-life insurance claims. These individual claims histories should depend on individual claims feature information such as the line of business concerned, the claims code involved or the age of the injured. This feature information should influence the reporting delay of the individual claim, the claim amount paid, its individual cash flow pattern as well as its settlement delay. The resulting (simulated) individual claims histories should be as 'realistic' as possible so that they may reflect a real insurance claims portfolio. These simulated claims then allow us to back-test classical aggregate claims reserving methods-such as the chain-ladder method-as well as to develop new claims reserving methods which are based on individual claims histories. The latter has become increasingly popular in actuarial science, see Antonio and Plat (2014), Hiabu et al. (2016), Jessen et al. (2011), Martínez-Miranda et al. (2015), Pigeon et al. (2013), Taylor et al. (2008), Verrall and Wüthrich (2016) and Wüthrich (2018a) for recent developments. A main shortcoming in this field of research is that there is no publicly available individual claims history data. Therefore, there is no possibility to back-test the proposed individual claims reserving methods. For this reason, we believe that this project is very beneficial to the actuarial community because it provides a common ground and publicly available (synthetic) data for research in the field of individual claims reserving.

This paper is divided into four sections. In this first section we describe the general idea of the simulation machine as well as the chosen data used for model calibration. In Section 2 we describe the design of our individual claims history simulation machine using neural networks. Section 3 focuses on the calibration of these neural networks. In Section 4 we carry out a use test by comparing the real data to the synthetically generated data in a chain-ladder claims reserving analysis. Appendix A presents descriptive statistics of the real data. Since the real insurance portfolio is confidential, we also design an algorithm to generate synthetic insurance portfolios of a similar structure as the real one, see Appendix B. Finally, in Appendix C we provide sensitivity plots of selected neural networks.

#### 1.1. Description of the Simulation Machine

The simulation machine is programmed in the language R. The corresponding .zip-folder can be downloaded from the website:

## https://people.math.ethz.ch/~wmario/simulation.html

This .zip-folder contains all parameters, a file readme.pdf which describes the use of our R-functions, as well as the two R-files Functions.V1 and Simulation.Machine.V1. The first R-file Functions.V1 contains the two R-functions Feature.Generation and Simulation.Machine. The former is used to generate synthetic insurance portfolios (this is described in more detail in Appendix B) and the latter to simulate the corresponding individual claims histories (this is described in the main body of this manuscript). The R-file Simulation.Machine.V1 demonstrates the use of these two R-functions, also providing a short chain-ladder claims reserving analysis.

#### 1.2. Procedure of Developing the Simulation Machine

In recent years, neural networks have become increasingly popular in all fields of machine learning. They have proved to be very powerful tools in classification and regression problems. Their drawbacks are that they are rather difficult to calibrate and, once calibrated, they act almost like black boxes between inputs and outputs. Of course, this is a major disadvantage in interpretation and getting deeper insight. However, the missing interpretation is not necessarily a disadvantage in our project because it implies—in back-testing other methods—that the true data generating mechanism cannot easily be guessed.

To construct our individual claims history simulation machine, we design a neural network architecture. This architecture is calibrated to real insurance data consisting of n = 9,977,298 individual claims that have occurred between 1994 and 2005. For each of these individual claims, we have full information of 12 years of claims development as well as the relevant feature information. Together with a portfolio generating algorithm (see Appendix B), one can then use the calibrated simulation machine to simulate as many individual claims development histories as desired.

## 1.3. The Chosen Data

The chosen data has been preprocessed correcting for wrong entries—for instance, an accident date that is bigger than the reporting date, etc. Moreover, we have dropped claims with missing feature components—for instance, if the age of the injured was missing. However, this was a negligible number of claims that we had to drop, and this does not distort the general calibration. The final (cleaned) data set consists of n = 9,977,298 individual claims histories. The following feature information is available for each individual claim:

- the claims number ClNr, which serves as a distinct claims identifier;
- the line of business LoB, which is categorical with labels in {1,...,4};
- the claims code cc, which is categorical with labels in {1,...,53} and denotes the labor sector of the injured;
- the accident year AY, which is in {1994, ..., 2005};
- the accident quarter AQ, which is in {1,...,4};
- the age of the injured age (in 5 years age buckets), which is in {15, 20, ..., 70};
- the injured part inj\_part, which is categorical with labels in {10,...,99} and denotes the part of the body injured;
- the reporting year RY, which is in  $\{1994, \ldots, 2016\}$ .

Not all values in  $\{10, \ldots, 99\}$  are needed for the labeling of the categorical classes of the feature component inj\_part. In fact, only 46 different values are attained, but for simplicity, we have decided to keep the original labeling received from the insurance company. 46 different labels may still seem to

be a lot and a preliminary classification could allow to reduce this number, here we refrain from doing so because each label has sufficient volume.

For all claims i = 1, ..., n, we are given the individual claims cash flow  $(C_i^{(j)})_{0 \le j \le 11}$ , where  $C_i^{(j)}$  is the payment for claim *i* in calendar year AY<sub>i</sub> + *j*—and where AY<sub>i</sub> denotes the accident year of claim *i*. Note that we only consider yearly payments, i.e., multiple payments and recovery payments within calendar year AY<sub>i</sub> + *j* are aggregated into a single, annual payment  $C_i^{(j)}$ . This single, annual payment can either be positive or negative, depending on having either more claim payments or more recovery payments in that year. The sum over all yearly payments  $\sum_j C_i^{(j)}$  of a given claim *i* has to be non-negative because recoveries cannot exceed payments (this is always the case in the considered data). Remark that our simulation machine will allow for recoveries.

Finally, for claims i = 1, ..., n, we are given the claim status process  $(I_i^{(j)})_{0 \le j \le 11}$  determining whether claim *i* is open or closed at the end of each accounting year. More precisely, if  $I_i^{(j)} = 1$ , claim *i* is open at the end of accounting year AY<sub>i</sub> + *j*, and if  $I_i^{(j)} = 0$ , claim *i* is closed at the end of that accounting year. Our simulation machine also allows for re-opening of claims, which is quite common in our real data. More description of the data is given in Appendix A.

## 2. Design of the Simulation Machine Using Neural Networks

In this section we describe the architecture of our individual claims history simulation machine. It consists of eight modeling steps: (1) reporting delay *T* simulation; (2) payment indicator *Z* simulation; (3) number of payments *K* simulation; (4) total claim size *Y* simulation; (5) number of recovery payments  $K^-$  simulation; (6) recovery size  $Y^-$  simulation; (7) cash flow  $(C_i^{(j)})_{0 \le j \le 11}$  simulation and (8) claim status  $(I_i^{(j)})_{0 \le j \le 11}$  simulation. Each of these eight modeling steps is based on one or several feed-forward neural networks. We introduce the precise setup of such a neural network in Section 2.1 for the simulation of the reporting delay *T*. Before, we present a global overview of the architecture of our simulation machine. Afterwards, in Sections 2.1–2.8, each single step is described in detail.

To start with, we define the initial feature space  $X_1$  consisting of the original six feature components as

$$\mathcal{X}_1 = \{ (LoB, cc, AY, AQ, age, inj_part) \}.$$
(1)

Observe that we drop the claims number C1Nr because it does not have explanatory power. Apart from these six feature values, the only other model-dependent input parameters of our simulation machine are the standard deviations for the total individual claim sizes and the total individual recoveries, see Sections 2.4 and 2.6 below. During the simulation procedure, not all of the subsequent steps (1)–(8) may be necessary—e.g., if we do not have any payments, then there is no need to simulate the claim size or the cash flow pattern. We briefly describe the eight modeling steps (1)–(8).

(1) In the first step, we use the initial feature space  $X_1$  to model the reporting delay *T* indicating the annualized difference between the reporting year and the accident year.

(2) For the second step, we extend the initial feature space  $X_1$  by including the additional information of the reporting delay *T*, i.e., we set

$$\mathcal{X}_2 = \{ (\text{LoB}, \text{cc}, \text{AY}, \text{AQ}, \text{age}, \text{inj}_{\text{part}}, T) \}.$$
(2)

We use  $X_2$  to model the payment indicator *Z* determining whether we have a payment or not.

(3) For the third step, we set  $X_3 = X_2$  and model the number of (yearly) payments *K*.

(4) In the fourth step, we extend the feature space  $X_3$  by including the additional information of the number of payments *K*, i.e., we set

$$\mathcal{X}_4 = \{ (\text{LoB}, \text{cc}, \text{AY}, \text{AQ}, \text{age}, \text{inj}_{\text{part}}, T, K) \},$$
(3)

which is used to model the total individual claim size *Y*.

(5) In the fifth step, we model the number of recovery payments  $K^-$ . We therefore work on the extended feature space

$$\mathcal{X}_5 = \{ (LoB, cc, AY, AQ, age, inj_part, T, K, Y) \}.$$
(4)

(6) In the sixth step, we model the total individual recovery  $Y^-$ . To this end, we set  $\mathcal{X}_6 = \mathcal{X}_5$ . We understand the total individual claim size Y to be net of recovery  $Y^-$ . Thus, the total payment from the insurance company to the insured is  $Y + Y^-$ , paid in  $K - K^-$  yearly payments. The total recovery from the insured to the insurance company is  $Y^-$ , paid in  $K^-$  yearly payments.

(7) In the seventh step, the task is to generate the cash flows  $(C_i^{(j)})_{0 \le j \le 11}$ . Therefore, we have to split the total gross claim amount  $Y + Y^-$  into  $K - K^-$  positive payments and the total recovery  $Y^-$  into  $K^-$  negative payments and distribute these K payments among the 12 development years. For this modeling step, we use different feature spaces  $\mathcal{X}_{7a}, \ldots, \mathcal{X}_{7g}$ , all being a subset of

$$\mathcal{X}_7 = \left\{ \left( \text{LoB}, \text{cc}, \text{AY}, \text{AQ}, \text{age}, \text{inj}_{\text{part}}, T, K, Y, K^-, Y^- \right) \right\},$$
(5)

see Section 2.7 below for more details.

(8) In the last step, we model the claim status process  $(I_i^{(j)})_{0 \le j \le 11}$ , where we use the feature space

$$\mathcal{X}_8 = \left\{ \left( \text{LoB, AQ}, T, \left( C^{(j)} \right)_{0 \le j \le 11} \right) \right\}$$

Each of these eight modeling steps (1)–(8) consists of one or even multiple feature-response problems, for which we design neural networks. In the end, the full individual claims history simulation machine consists of 35 neural networks. We are going to describe this neural network architecture in more detail next. We remark that some of these networks are rather similar. Therefore, we present the first neural network in full detail, and for the remaining neural networks we focus on the differences to the previous ones.

## 2.1. Reporting Delay Modeling

To model the reporting delay, we work with the initial feature space  $X_1$  given in (1). Let  $n_1 = n = 9,977,298$  be the number of individual claims in our data. We consider the (annualized) reporting delays  $T_i$ , for  $i = 1, ..., n_1$ , given by

$$T_i = \mathtt{R}\mathtt{Y}_i - \mathtt{A}\mathtt{Y}_i \in \mathcal{T} = \{0, \ldots, 11\},\$$

where  $AY_i$  is the accident year and  $RY_i$  the reporting year of claim *i*. For confidentiality reasons, we have only received data on a yearly time scale (with the additional information of the accident quarter AQ). A more accurate modeling would use a finer time scale.

The three feature components LoB, cc and inj\_part are categorical. For neural network modeling, we need to transform these categorical feature components to continuous ones. This could be done by dummy coding, but we prefer the following version because it leads to less parameters. We replace, for instance, the claims code cc by the sample mean of the reporting delay restricted to the corresponding feature label, i.e., for claims code cc = a, we set

$$a \mapsto a^* = a^*(a) = \frac{\sum_{i=1}^{n_1} T_i \mathbb{1}_{\{cc_i=a\}}}{\sum_{i=1}^{n_1} \mathbb{1}_{\{cc_i=a\}}} \in \mathbb{R},$$
(6)

where  $cc_1, ..., cc_{n_1}$  are the observed claims codes. By slight abuse of notation, we obtain a d = 6 dimensional feature space  $\mathcal{X}_1$  where we may assume that all feature components of  $\mathcal{X}_1$  are continuous. Such feature pre-processing as in (6) will be necessary throughout this section for the components LoB,

cc and  $inj_part$ : we just replace  $T_i$  in (6) by the respective response variable. Note that from now on this will be done without any further reference.

The above procedure equips us with the data

$$\mathcal{D}_1 = \{(x_1, T_1), \ldots, (x_{n_1}, T_{n_1})\},\$$

with  $x_1, \ldots, x_{n_1} \in \mathcal{X}_1$  being the observed features and  $T_1, \ldots, T_{n_1} \in \mathcal{T}$  the observed responses. For an insurance claim with feature  $x \in \mathcal{X}_1$ , the corresponding reporting delay T(x) is modeled by a categorical distribution

$$\mathbb{P}[T(\mathbf{x}) = t] = \pi_t(\mathbf{x}), \quad \text{for } t \in \mathcal{T}.$$

This requires that we model probability functions of the form

$$\pi_t: \mathcal{X}_1 \to [0,1], \qquad x \mapsto \pi_t(x),$$

satisfying normalization  $\sum_{t \in T} \pi_t(x) = 1$ , for all  $x \in \mathcal{X}_1$ . We design a neural network for the modeling of these probability functions and we estimate the corresponding network parameters from the observations  $\mathcal{D}_1$ .

We choose a classical feed-forward neural network with multiple layers. Each layer consists of several neurons, and weights connect all neurons of a given layer to all neurons of the next layer. Moreover, we use a non-linear activation function to pass the signals from one layer to the next. The first layer—consisting of the components  $x_1, \ldots, x_d$  of a feature  $\mathbf{x} = (x_1, \ldots, x_d) \in \mathcal{X}_1$ —is called input layer (blue circles in Figure 1). In our case, we have d = 6 neurons in this input layer. The last layer is called output layer (red circles in Figure 1) and it contains the categorical probabilities  $\pi_0(\mathbf{x}), \ldots, \pi_{11}(\mathbf{x})$ . In between these two layers, we choose two hidden layers having  $q_1$  and  $q_2$  hidden neurons, respectively (black circles in Figure 1 with  $q_1 = 11$  and  $q_2 = 15$ ).



**Figure 1.** Deep neural network with two hidden layers: the first column (blue circles) illustrates the d = 6 dimensional feature vector x (input layer), the second column gives the first hidden layer with  $q_1 = 11$  neurons, the third column gives the second hidden layer with  $q_2 = 15$  neurons and the fourth column gives the output layer (red circle) with 12 neurons.

More formally, we choose the  $q_1$  hidden neurons  $z_1^{(1)}, \ldots, z_{q_1}^{(1)}$  in the first hidden layer as follows

$$z_j^{(1)} = z_j^{(1)}(\mathbf{x}) = \phi\left(w_{j,0}^{(1)} + \sum_{l=1}^d w_{j,l}^{(1)} x_l\right), \quad \text{for all } j = 1, \dots, q_1,$$

for given weights  $(w_{j,l}^{(m)})_{j,l,m}$  and for the hyperbolic tangent activation function

$$\phi(x) = \tanh(x).$$

This is a centered version of the sigmoid activation function, with range (-1, 1). Moreover, we have  $\phi' = 1 - \phi^2$ , which is a useful property in the gradient descent method described in Section 3, below.

The activation is then propagated in an analogous fashion to the  $q_2$  hidden neurons  $z_1^{(2)}, \ldots, z_{q_2}^{(2)}$  in the second hidden layer, that is, we set

$$z_j^{(2)} = z_j^{(2)}(\mathbf{x}) = \phi\left(w_{j,0}^{(2)} + \sum_{l=1}^{q_1} w_{j,l}^{(2)} z_l^{(1)}(\mathbf{x})\right), \quad \text{for all } j = 1, \dots, q_2.$$

For the 12 neurons  $\pi_0(x), \ldots, \pi_{11}(x)$  in the output layer, we use the multinomial logistic regression assumption

$$\pi_t(\mathbf{x}) = \frac{\exp\left\{\mu_t(\mathbf{x})\right\}}{\sum_{s \in \mathcal{T}} \exp\left\{\mu_s(\mathbf{x})\right\}}, \quad \text{for all } t \in \mathcal{T},$$
(7)

with regression functions  $x \mapsto \mu_t(x)$  for all  $t \in \mathcal{T}$  given by

$$\mu_t(\mathbf{x}) = \beta_0^{(t)} + \sum_{j=1}^{q_2} \beta_j^{(t)} z_j^{(2)}(\mathbf{x}),$$

for given weights  $(\beta_i^{(t)})_{j,t}$ . We define the network parameter  $\alpha$  of all involved parameters by

$$\boldsymbol{\alpha} = \left( w_{1,0}^{(1)}, \dots, w_{q_2,q_1}^{(2)}, \beta_0^{(0)}, \dots, \beta_{q_2}^{(11)} \right)' \in \mathbb{R}^{q_1(d+1)+q_2(q_1+1)+12(q_2+1)}$$

The classification model for the tuples  $(x, T(x))_{x \in \mathcal{X}_1}$  is now fully defined and there remains the calibration of the network parameter  $\alpha$  and the choice of the hyperparameters  $q_1$  and  $q_2$ . Assume for the moment that  $q_1$  and  $q_2$  are given. In order to fit  $\alpha$  to our data  $\mathcal{D}_1$ , we aim to minimize a given loss function  $\alpha \mapsto \mathcal{L}(\alpha)$ . Therefore, we assume that  $(x_1, T_1), \ldots, (x_{n_1}, T_{n_1})$  are drawn independently from the joint distribution of (x, T(x)). The corresponding deviance statistics loss function of the categorical distribution of our data  $\mathcal{D}_1$  is then given by

$$\mathcal{L}(\boldsymbol{\alpha}) = \mathcal{L}_{\mathcal{D}_{1}}(\boldsymbol{\alpha}) = -2\log\left(\prod_{i=1}^{n_{1}}\sum_{t\in\mathcal{T}}\mathbb{1}_{\{T_{i}=t\}}\pi_{t}(\boldsymbol{x}_{i})\right) = -2\sum_{i=1}^{n_{1}}\sum_{t\in\mathcal{T}}\mathbb{1}_{\{T_{i}=t\}}\log\pi_{t}(\boldsymbol{x}_{i}).$$
(8)

The optimal network parameter  $\alpha$  is found by minimizing this deviance statistics loss function. We come back to this problem in Section 3.2.1, below. Since for different hyperparameters  $q_1$  and  $q_2$  we get different network structures, every pair ( $q_1$ ,  $q_2$ ) corresponds to a separate model. The choice of appropriate hyperparameters  $q_1$  and  $q_2$  is discussed in Section 3.3, below.

After the calibration of  $q_1$ ,  $q_2$  and  $\alpha$  to our data  $\mathcal{D}_1$ , we can simulate the reporting delay T(x) of a claim with given feature  $x \in \mathcal{X}_1$  by using the resulting categorical distribution given by (7). This simulated value will then allow us to go to the next modeling step (2), see (2).

We close this first part with the following remark: Our choice to work with two hidden layers may seem arbitrary since we could also have chosen more hidden layers or just one of them. From a theoretical point of view, one hidden layer would be sufficient to approximate a vast collection of regression functions to any desired degree of accuracy, provided that we have sufficiently many hidden neurons in that layer, see Cybenko (1989) and Hornik et al. (1989). However, these models with large-scale numbers of hidden neurons are known to be difficult to calibrate, and it is often

more efficient to use fewer neurons but more hidden layers to get an appropriate complexity in the regression function.

#### 2.2. Payment Indicator Modeling

In our real data, we observe that roughly 29% of all claims can be settled without any payment. For this reason, we model the claim sizes by compound distributions. First, we model a payment indicator Z that determines whether we have a payment or not. Then, conditionally on having a payment, we determine the exact number of payments K. Finally, we model the total individual claim size Y for claims with at least one payment.

In order to model the payment indicator, we work with the d = 7 dimensional feature space  $\mathcal{X}_2$  introduced in (2). Let  $n_2 = n_1$  and  $x_1, \ldots, x_{n_2} \in \mathcal{X}_2$  be the observed features, where this time the reporting delay *T* is also included. For all  $i = 1, \ldots, n_2$ , we define the number of payments  $K_i$  and the payments indicator  $Z_i$  by

$$K_i = \sum_{j=0}^{11} \mathbb{1}_{\left\{C_i^{(j)} \neq 0\right\}} \quad \text{and} \quad Z_i = \mathbb{1}_{\{K_i > 0\}}.$$
(9)

This provides us with the data

$$\mathcal{D}_2 = \{(x_1, Z_1), \dots, (x_{n_2}, Z_{n_2})\}$$

For a claim with feature  $x = (x_1, ..., x_d) \in \mathcal{X}_2$ , the corresponding payment indicator Z(x) is a Bernoulli random variable with

$$\mathbb{P}\left[Z(\boldsymbol{x})=1\right] = \pi(\boldsymbol{x}),$$

for a given (but unknown) probability function

$$\pi: \mathcal{X}_2 \rightarrow [0,1], \qquad x \mapsto \pi(x).$$

Note that this Bernoulli model is a special case of the categorical model of Section 2.1. Therefore, it can be calibrated completely analogously, as described above. However, we emphasize that instead of working with two probability functions  $\pi_0$  and  $\pi_1$  for the two categories  $\{0, 1\}$ , we set  $\pi(\cdot) = \pi_1(\cdot)$ , which implies  $1 - \pi(\cdot) = 1 - \pi_1(\cdot) = \pi_0(\cdot)$ . Moreover, the multinomial probabilities (7) simplify to the binomial case

$$\pi(\mathbf{x}) = \frac{\exp\{\mu_1(\mathbf{x})\}}{\exp\{\mu_0(\mathbf{x})\} + \exp\{\mu_1(\mathbf{x})\}} = \frac{1}{1 + \exp\{-(\mu_1(\mathbf{x}) - \mu_0(\mathbf{x}))\}} = \frac{1}{1 + \exp\{-\mu(\mathbf{x})\}},$$

with regression function

$$\mu: \mathcal{X}_2 \to \mathbb{R}, \qquad \mathbf{x} \mapsto \mu(\mathbf{x}) = \beta_0 + \sum_{j=1}^{q_2} \beta_j z_j^{(2)}(\mathbf{x})$$
 (10)

for a neural network with two hidden layers and network parameter  $\alpha$  given by

$$\boldsymbol{\alpha} = \left( w_{1,0}^{(1)}, \dots, w_{q_2,q_1}^{(2)}, \beta_0, \dots, \beta_{q_2} \right)' \in \mathbb{R}^{q_1(d+1)+q_2(q_1+1)+(q_2+1)}$$

Finally, the corresponding deviance statistics loss function to be minimized is given by

$$\mathcal{L}(\boldsymbol{\alpha}) = \mathcal{L}_{\mathcal{D}_2}(\boldsymbol{\alpha}) = -2\sum_{i=1}^{n_2} Z_i \log \pi(\boldsymbol{x}_i) + (1 - Z_i) \log(1 - \pi(\boldsymbol{x}_i)).$$
(11)

From this calibrated model, we simulate the payment indicator Z(x), which then allows us to go to the next modeling step. If this indicator is equal to one, we move to step (3), see Section 2.3; if this indicator is equal to zero, we directly go to step (8), see Section 2.8.

#### 2.3. Number of Payments Modeling

We use the d = 7 dimensional feature space  $\mathcal{X}_3 = \mathcal{X}_2$  to model the number of payments, conditioned on the event that the payment indicator Z is equal to one. We define  $n_3 \le n_2$  to be the number of claims with payment indicator equal to one and order the claims appropriately in i such that  $Z_i = 1$  for all  $i = 1, ..., n_3$ . Then, we define the number of payments  $K_i$  as in (9), for all  $i = 1, ..., n_3$ . This gives us the data

$$\mathcal{D}_3 = \{(x_1, K_1), \ldots, (x_{n_3}, K_{n_3})\}.$$

For a claim with feature  $x \in X_3$  and payment indicator Z = 1, we could now proceed as in Section 2.1 in order to model the number of payments K(x). However, the claims with  $K_i = 1$  are so dominant in the data that a good calibration of the categorical model (7) becomes difficult. For this reason, we choose a different approach: in a first step, we model the events {K(x) = 1} and {K(x) > 1}, conditioned on {Z = 1}, and, in a second step, we consider the conditional distribution of K(x), given K(x) > 1. In particular, in the first step we have a Bernoulli classification problem that is modeled completely analogously to Section 2.2, only replacing the data  $D_2$  by

$$\mathcal{D}_{3a} = \left\{ \left( x_1, \mathbb{1}_{\{K_1=1\}} \right), \dots, \left( x_{n_3}, \mathbb{1}_{\{K_{n_3}=1\}} \right) \right\}.$$

The case K(x) > 1 is then modeled analogously to the categorical case of Section 2.1, with 11 categories and data  $\mathcal{D}_{3b} \subset \mathcal{D}_3$  only considering the claims with more than one payment.

The simulation of the number of payments  $K(\mathbf{x})$  for a claim with feature  $\mathbf{x} \in \mathcal{X}_3$ , reporting delay T and payment indicator Z = 1 needs more care than the corresponding task in Section 2.1: here we have the restriction  $T + K(\mathbf{x}) \le 12$ . If T = 11, then we automatically need to have  $K(\mathbf{x}) = 1$ . For T < 11 and if the first neural network leads to  $\mathbb{1}_{\{K(\mathbf{x})=1\}} = 0$ , then the categorical conditional distribution for  $K(\mathbf{x})$ , given  $K(\mathbf{x}) > 1$ , can only take the values  $k \in \{2, ..., 12 - T\}$ . For this reason, instead of using the original conditional probabilities  $\pi_2(\mathbf{x}), ..., \pi_{12}(\mathbf{x})$  resulting from the second neural network, we use in that case the modified conditional probabilities  $\pi_k^*(\mathbf{x})$ , for  $k \in \{2, ..., 12 - T\}$ , given by

$$\pi_k^*(\mathbf{x}) = \frac{\pi_k(\mathbf{x})}{\sum_{l=2}^{12-T} \pi_l(\mathbf{x})}.$$
(12)

#### 2.4. Total Individual Claim Size Modeling

For the modeling of the total individual claim size, we add the number of payments *K* to the previous feature space and work with  $X_4$  given in (3). Let  $n_4 = n_3$  and consider the same ordering of the claims as in Section 2.3. Then, we define the total individual claim size  $Y_i$  of claim *i* as

$$Y_i = \sum_{j=0}^{11} C_i^{(j)} > 0$$

for all  $i = 1, ..., n_4$ . In particular, the total individual claim size  $Y_i$  is always to be understood net of recoveries. This leads us to the data

$$\mathcal{D}_4 = \{ (x_1, Y_1), \dots, (x_{n_4}, Y_{n_4}) \}.$$

For a claim with feature  $x \in \mathcal{X}_4$  and payment indicator Z = 1, we model the total individual claim size Y(x) with a log-normal distribution. We therefore choose a regression function

$$\mu: \mathcal{X}_4 \to \mathbb{R}$$

of type (10) for a neural network with two hidden layers. This regression function is used to model the mean parameter of the total individual claim sizes, i.e., we make the model assumption

$$Y(\mathbf{x}) \mid Z = 1 \sim \operatorname{LN}\left(\mu(\mathbf{x}), \sigma_{+}^{2}\right),$$
(13)

for given variance parameter  $\sigma_+^2 > 0$ . This choice implies

$$\mathbb{E} \left[ \log Y(x) \, | \, Z = 1 \right] = \mu(x) \quad \text{and} \quad \text{Var} \left( \log Y(x) \, | \, Z = 1 \right) = \sigma_+^2.$$

The density of  $\log Y(x) | Z = 1$  then motivates the choice of the square loss function (deviance statistics loss function)

$$\mathcal{L}(\boldsymbol{\alpha}) = \mathcal{L}_{\mathcal{D}_4}(\boldsymbol{\alpha}) = \sum_{i=1}^{n_4} (\log Y_i - \mu(\boldsymbol{x}_i))^2, \qquad (14)$$

with network parameter  $\alpha$ . The optimal model for the total individual claim size is then found by minimizing the loss function (14), which does not depend on  $\sigma_{+}^{2}$ .

This calibrated model together with the input parameter  $\sigma_+ > 0$  can be used to simulate the total individual claim size Y(x) from (13). Note that the expected claim amount is increasing in  $\sigma_+^2$ , as we have

$$\mathbb{E}\left[Y(\boldsymbol{x}) \mid Z=1\right] = \exp\left\{\mu(\boldsymbol{x}) + \frac{\sigma_+^2}{2}\right\}.$$

#### 2.5. Number of Recovery Payments Modeling

For the modeling of the number of recovery payments, we use the d = 9 dimensional feature space  $\mathcal{X}_5$  introduced in (4). Furthermore, we only consider claims *i* with  $K_i > 1$ , because recoveries may only happen if we have at least one positive payment. We define  $n_5 \le n_4$  to be the number of claims with more than one payment and order the claims appropriately in *i* such that  $K_i > 1$  for all  $i = 1, ..., n_5$ . Then, we define the number of recovery payments  $K_i^-$  of claim *i* as

$$K_i^- = \min\left(\sum_{j=0}^{11} \mathbb{1}_{\left\{C_i^{(j)} < 0\right\}}, 2\right), \tag{15}$$

for all  $i = 1, ..., n_5$ . In particular, for all observed claims *i* with more than two recovery payments, we set  $K_i^- = 2$ . This reduces combinatorial complexity in simulations (without much loss of accuracy) and provides us with the data

$$\mathcal{D}_5 = \{ (x_1, K_1^-), \dots, (x_{n_5}, K_{n_5}^-) \}.$$

For a claim with feature  $x \in \mathcal{X}_5$  and K > 1 payments, the corresponding number of recovery payments  $K^-(x)$ , conditioned on the event  $\{K > 1\}$ , is a categorical random variable taking values in  $\{0, 1, 2\}$ , i.e., we are in the same setup as in Section 2.1—with only three categorical classes. Thus, the calibration is done analogously.

This model then allows us to simulate the number of recovery payments  $K^-(x)$ . Note that also this simulation step needs additional care: if K = 2, then we can have at most one recovery payment. Thus, we have to apply a similar modification as given in (12) in this case.

#### 2.6. Total Individual Recovery Size Modeling

The modeling of the total individual recovery size is based on the feature space  $\mathcal{X}_6 = \mathcal{X}_5$ , given in (4), and we restrict to claims with  $K_i^- > 0$ . The number of these claims is denoted by  $n_6 \le n_5$ . Appropriate ordering provides us with the total individual recovery  $Y_i^-$  of claim *i* as

$$Y_i^- = -\sum_{j=0}^{11} C_i^{(j)} \, \mathbb{1}_{\left\{C_i^{(j)} < 0\right\}}$$

for all  $i = 1, ..., n_6$ . This gives us the data

$$\mathcal{D}_6 = \{ (x_1, Y_1^-), \dots, (x_{n_6}, Y_{n_6}^-) \}.$$

The remaining part is completely analogous to Section 2.4, we only need to replace the standard deviation parameter  $\sigma^+$  by a given  $\sigma^- > 0$ .

#### 2.7. Cash Flow Pattern Modeling

The modeling of the cash flow pattern is more involved, and we need to distinguish different cases. This distinction is done according to the total number of payments K = 1, ..., 12, the number of positive payments  $K^+ = K - K^- = 1, ..., 12$  as well as the number of recovery payments  $K^- = 0, 1, 2$ .

## 2.7.1. Cash Flow for Single Payments K = 1

The simplest case is the one of having exactly one payment  $K = K^+ = 1$ . In this case, we consider the payment delay after the reporting date. We define  $n_{7a} \le n_3$  to be the number of claims with exactly one payment and order the claims appropriately in *i* such that  $K_i = 1$  for all  $i = 1, ..., n_{7a}$ . Then, we define the payment delay  $S_i$  of claim *i* as

$$S_i = \sum_{j=0}^{11} j \mathbb{1}_{\left\{C_i^{(j)} > 0\right\}} - T_i \ge 0,$$

for all  $i = 1, ..., n_{7a}$ . In other words, we simply subtract the reporting year from the year in which the unique payment occurs. This provides us with the data

$$\mathcal{D}_{7a} = \{(x_1, S_1), \ldots, (x_{n_{7a}}, S_{n_{7a}})\},\$$

with  $x_1, \ldots, x_{n_{7a}} \in \mathcal{X}_{7a}$  being the observed features, where we use

$$\mathcal{X}_{7a} = \{ (LoB, cc, AY, AQ, age, inj_part, T, Y) \}$$
(16)

as d = 8 dimensional feature space. For a claim with feature  $x \in X_{7a}$  and K = 1 payment, the corresponding payment delay S(x) is a categorical random variable assuming values in  $\{0, ..., 11\}$ . Similarly as for the number of payments, the claims with  $S_i = 0$  are rather dominant. Therefore, we apply the same two-step modeling approach as in Section 2.3.

This calibrated model then allows us to simulate the payment delay  $S(\mathbf{x})$ . For given reporting delay T, we have the restriction  $T + S(\mathbf{x}) \le 11$ , which is treated in the same way as in (12). Finally, the cash flow is given by  $(C^{(j)}(\mathbf{x}))_{0 \le j \le 11}$  with

$$C^{(j)}(\mathbf{x}) = \begin{cases} Y, & \text{if } j = T + S(\mathbf{x}), \\ 0, & \text{else.} \end{cases}$$

2.7.2. Cash Flow for Two Payments K = 2

Now we consider claims with exactly two payments. Here we distinguish further between the two cases: (1) both payments are positive, and (2) one payment is positive and the other one negative.

(a) Two Positive Payments

We first consider the case where both payments are positive, i.e.,  $K = K^+ = 2$  and  $K^- = 0$ . In this case, we have to model the time points of the two payments as well as the split of the total individual claim size to the two payments. For both models, we use the d = 8 dimensional feature space  $\mathcal{X}_{7b} = \mathcal{X}_{7a}$ , see (16). We define  $n_{7b} \le n_3$  to be the number of claims with exactly two positive payments and no recovery and order them appropriately in *i* such that  $K_i = 2$  and  $K_i^- = 0$  for all  $i = 1, \ldots, n_{7b}$ . The time points  $R_i^{(1)}$  and  $R_i^{(2)}$  of the two payments are given by

$$R_i^{(1)} = \min\left\{ 0 \le j \le 11 \mid C_i^{(j)} \ne 0 \right\}$$
 and  $R_i^{(2)} = \max\left\{ 0 \le j \le 11 \mid C_i^{(j)} \ne 0 \right\}$ ,

for all  $i = 1, ..., n_{7b}$ . Then, we modify the two-dimensional vector  $(R_i^{(1)}, R_i^{(2)})$  to a one-dimensional categorical variable  $R_i$  by setting

$$R_{i} = \begin{cases} R_{i}^{(2)}, & \text{if } R_{i}^{(1)} = 0, \\ R_{i}^{(2)} - R_{i}^{(1)} + \sum_{k=12-R_{i}^{(1)}}^{11} k, & \text{else,} \end{cases}$$
(17)

for all  $i = 1, ..., n_{7b}$ . This leads us to the data

$$\mathcal{D}_{7b} = \{ (x_1, R_1), \dots, (x_{n_{7b}}, R_{n_{7b}}) \}.$$

Note that  $R_i$  is categorical with  $\binom{12}{2} = 66$  possible values. That is, we are in the same setup as in Section 2.1—with 66 different classes. Once again, the calibration is done in an analogous fashion as above.

Next, we model the split of the total individual claim size for claims with  $K = K^+ = 2$ . Let  $n_{7c} = n_{7b}$ ,  $\mathcal{X}_{7c} = \mathcal{X}_{7a}$ , see (16), and define the proportion  $P_i$  of the total individual claim size  $Y_i$  that is paid in the first payment by

$$P_i = \frac{C_i^{\left(R_i^{(1)}\right)}}{Y_i},$$

for all  $i \in 1, ..., n_{7c}$ . This gives us the data

$$\mathcal{D}_{7c} = \{(x_1, P_1), \dots, (x_{n_{7c}}, P_{n_{7c}})\}.$$

For a claim with feature  $x \in \mathcal{X}_{7c}$  and  $K = K^+ = 2$ , the corresponding proportion of its total individual claim size *Y* that is paid in the first payment is for simplicity modeled by a deterministic function P(x). Note that one could easily randomize P(x) using a Dirichlet distribution. However, at this modeling stage, the resulting differences would be of smaller magnitude. Hence, we directly fit the proportion function

$$P: \mathcal{X}_{7c} \to [0,1], \qquad \mathbf{x} \mapsto P(\mathbf{x})$$

Similarly to the calibration in Section 2.2, we assume a regression function  $\mu : \mathcal{X}_{7c} \to \mathbb{R}$  of type (10) for a neural network with two hidden layers. Then, for the output layer, we use

$$P(\mathbf{x}) = \frac{1}{1 + \exp\{-\mu(\mathbf{x})\}}$$
(18)

and as loss function the cross entropy function, see also (11),

$$\mathcal{L}(\boldsymbol{\alpha}) = \mathcal{L}_{\mathcal{D}_{7c}}(\boldsymbol{\alpha}) = -2\sum_{i=1}^{n_{7c}} P_i \log P(\boldsymbol{x}_i) + (1 - P_i) \log(1 - P(\boldsymbol{x}_i)),$$
(19)

where  $\alpha$  is the network parameter containing all the weights of the neural network.

From this model, we can then simulate the cash flow for a claim with  $K = K^+ = 2$ . First, we simulate  $R(\mathbf{x})$ . If  $R(\mathbf{x}) \in \{1, ..., 11\}$ , we have  $R^{(1)}(\mathbf{x}) = 0$  and  $R^{(2)}(\mathbf{x}) = R(\mathbf{x})$ . If  $R(\mathbf{x}) > 11$ , we have

$$R^{(1)}(\mathbf{x}) = \max\left\{1 \le k \le 10 \left| R(\mathbf{x}) > \sum_{u=12-k}^{11} u \right\} \text{ and } R^{(2)}(\mathbf{x}) = R(\mathbf{x}) + R^{(1)}(\mathbf{x}) - \sum_{k=12-R^{(1)}(\mathbf{x})}^{11} k.$$

The cash flow is given by  $(C^{(j)}(\mathbf{x}))_{0 \le j \le 11}$  with

$$C^{(j)}(\mathbf{x}) = \begin{cases} P(\mathbf{x}) Y, & \text{if } j = R^{(1)}(\mathbf{x}), \\ (1 - P(\mathbf{x})) Y, & \text{if } j = R^{(2)}(\mathbf{x}), \\ 0, & \text{else.} \end{cases}$$

(b) One Positive Payment, One Recovery Payment

Now we focus on the case where we have  $K^+ = 1$  positive and  $K^- = 1$  negative payment. Here we only have to model the time points of the two payments, since we know the total individual claim size as well as the total individual recovery and we assume that the positive payment precedes the recovery payment. The modeling of the time points of the two payments is done as above, except that this time we use the d = 9 dimensional feature space

$$\mathcal{X}_{7d} = \{ (LoB, cc, AY, AQ, age, inj_part, T, Y, Y^{-}) \},\$$

where we include the information of the total individual recovery  $Y^-$ . Moreover, we define  $n_{7d} \le n_3$  to be the number of claims with exactly one positive payment and one recovery payment and order the claims appropriately in *i* such that  $K_i = 2$  and  $K_i^- = 1$  for all  $i = 1, ..., n_{7d}$ . This provides us with the data

$$\mathcal{D}_{7d} = \{ (x_1, R_1), \dots, (x_{n_{7d}}, R_{n_{7d}}) \},\$$

with  $R_i$  defined as in (17). The rest is done as above. We obtain the cash flow  $(C^{(j)}(x))_{0 \le j \le 11}$  with

$$C^{(j)}(\mathbf{x}) = \begin{cases} Y + Y^{-}, & \text{if } j = R^{(1)}(\mathbf{x}), \\ -Y^{-}, & \text{if } j = R^{(2)}(\mathbf{x}), \\ 0, & \text{else.} \end{cases}$$

Remark that we again have combinatorial complexity of  $\binom{12}{2} = 66$  for the time points of the two payments. Since data is sparse, for this calibration we restrict to the 35 most frequent distribution patterns. More details on this restriction are provided in the next section.

2.7.3. Cash Flow for More than Two Payments K = 3, ..., 12

On the one hand, the models for the cash flows in the case of more than two payments depend on the exact number of payments K. On the other hand, they also depend on the respective numbers of positive payments  $K^+$  and negative payments  $K^-$ . If we have zero or one recovery payment ( $K^- = 0, 1$ ), then we need to model (a) the time points where the K payments occur and (b) the proportions of the total gross claim amount  $Y + Y^-$  paid in the  $K^+$  positive payments. If  $K^- = 0$ , then there are no recovery payments and, thus,  $Y^- = 0$ . If  $K^- = 1$ , the recovery payment is always set at the end. In the case of  $K^- = 2$  recovery payments, in addition to (a) and (b), we use another neural network to model (c) the proportions of the total individual recovery  $Y^-$  paid in the two recovery payments. The time point of the first recovery payment is for simplicity assumed to be uniformly distributed on the set of time points of the 2nd up to the (K - 1)-st payment. The second recovery payment is always set at the end. The time point of the first payment is excluded for recovery in our model since we first

require a positive payment before a recovery is possible. The three neural networks considered in this modeling part are outlined below in (a)–(c). Afterwards, we can model the cash flow for claims with K = 3, ..., 12 payments, see item (d) below.

#### (a) Distribution of the K Payments

If we have K = 12 payments, then the distribution of these payments to the 12 development years is trivial, as we have a payment in every development year. Since the model is pretty much the same in all other cases  $K \in \{3, ..., 11\}$ , we present here the case K = 6 as illustration.

For the modeling of the distribution of the payments to the development years, we slightly simplify our feature space by dropping the categorical feature components cc and inj\_part. Moreover, we simplify the feature LoB with its four categorical classes: since the lines of business one and four as well as the lines of business two and three behave very similarly w.r.t. the cash flow patterns, we merge these lines of business in order to get more volume (and less complexity). We denote this simplified lines of business by LoB<sup>\*</sup>. Thus, we work with the d = 8 dimensional feature space

$$\mathcal{X}_{7e} \;=\; \left\{ \left( \mathtt{LoB}^{*}, \mathtt{AY}, \mathtt{AQ}, \mathtt{age}, T, Y, K^{-}, Y^{-} 
ight) 
ight\}.$$

Let  $n_{7e} \le n_3$  be the number of claims with exactly six payments and order the claims appropriately in *i* such that  $K_i = 6$  for all  $i = 1, ..., n_{7e}$ . The time points  $R_i^{(1)}, ..., R_i^{(6)}$  of the six payments are given by

$$R_i^{(1)} = \min\left\{ 0 \le j \le 11 \ \middle| \ C_i^{(j)} \ne 0 \right\} \quad \text{and} \quad R_i^{(k)} = \min\left\{ R_i^{(k-1)} < j \le 11 \ \middle| \ C_i^{(j)} \ne 0 \right\},$$

for all k = 2, ..., 6 and  $i = 1, ..., n_{7e}$ . Then, we use the following binary representation

$$R_i = \sum_{k=1}^{6} 2^{R_i^{(k)} + 1},$$

for all  $i = 1, ..., n_{7e}$ , for the time points of the six payments. This leads us to the data

$$\mathcal{D}_{7e} = \{(x_1, R_1), \dots, (x_{n_{7e}}, R_{n_{7e}})\},\$$

where  $x_1, \ldots, x_{n_{7e}} \in \mathcal{X}_{7e}$  and  $R_1, \ldots, R_{n_{7e}} \in \mathcal{A}$  for some set  $\mathcal{A} \subset \mathbb{N}$ . Since there are  $\binom{12}{6} = 924$  possibilities to distribute the K = 6 payments to the 12 development years, we have  $|\mathcal{A}| = 924$  distribution patterns. To reduce complexity (in view of sparse data), we only allow for the most frequently observed distributions of the payments to the development years. For K = 6, we work with 21 different patterns, which cover 70% of all claims with K = 6. We denote the set containing these 21 patterns by  $\tilde{\mathcal{A}}$ . See Table 1 for an overview, for each  $K = 3, \ldots, 10$ , of the number of possible different patterns, the number of allowed different patterns and the percentage of all claims covered with this choice of allowed distribution patterns.

**Table 1.** Number of possible and allowed distribution patterns for K = 3, ..., 10 payments.

Number of Payments K	3	4	5	6	7	8	9	10
number of possible patterns $ \mathcal{A} $	220	495	792	924	792	495	220	66
number of allowed patterns $ \widetilde{\mathcal{A}} $	15	18	17	21	17	20	26	24
percentage of claims covered	91%	83%	73%	70%	62%	61%	64%	76%

Note that for K = 11, we allow for all the 12 possible distribution patterns. Going back to the case K = 6, we denote by  $\tilde{n}_{7e} \le n_{7e}$  the number of claims with exactly K = 6 payments and with a distribution of these six payments to the 12 development years contained in the set  $\tilde{A}$ . Then, we

modify the data  $\mathcal{D}_{7e}$  accordingly to  $\widetilde{\mathcal{D}}_{7e}$  by only considering the relevant observations in  $\widetilde{\mathcal{A}}$ . This provides us with a classification problem similar to the one in Section 2.1—with  $|\widetilde{\mathcal{A}}| = 21$  classes.

## (b) Proportions of the $K^+ = K - K^-$ Positive Payments

If the number of positive payments  $K^+$  is equal to one, then the amount paid in this unique positive payment is given by the total gross claim amount  $Y + Y^-$ . That is, we do not need to model the proportions of the positive payments. Since the model is basically the same in all other cases  $K^+ \in \{2, ..., 12\}$ , we present here the case  $K^+ = 6$  as illustration.

As in the previous part, we use the d = 8 dimensional feature space  $\mathcal{X}_{7f} = \mathcal{X}_{7e}$ . Let  $n_{7f} \le n_3$  be the number of claims with exactly six positive payments and order the claims appropriately in *i* such that  $K_i^+ = K_i - K_i^- = 6$  for all  $i = 1, ..., n_{7f}$ . We define

$$R_i^{+(1)} = \min\left\{ 0 \le j \le 11 \left| C_i^{(j)} > 0 \right\} \quad \text{and} \quad R_i^{+(k)} = \min\left\{ R_i^{+(k-1)} < j \le 11 \left| C_i^{(j)} > 0 \right\},\right\}$$

for all k = 2, ..., 6 and  $i = 1, ..., n_{7f}$ , to be the time points of the six positive payments. Then, we can define

$$P_i^{(k)} = \frac{C_i^{\left(R_i^{+(k)}\right)}}{Y_i + Y_i^{-}}$$

to be the proportion of the total gross claim amount  $Y_i + Y_i^-$  that is paid in the *k*-th positive, annual payment, for all k = 1, ..., 6 and  $i = 1, ..., n_{7f}$ . This equips us with the data

$$\mathcal{D}_{7f} = \left\{ \left( x_1, P_1^{(1)}, \dots, P_1^{(6)} \right), \dots, \left( x_{n_{7f}}, P_{n_{7f}}^{(1)}, \dots, P_{n_{7f}}^{(6)} \right) \right\}.$$

For a claim with feature  $x \in \mathcal{X}_{7f}$  and  $K^+ = K - K^- = 6$  positive payments, the corresponding proportions  $P^{(1)}(x), \ldots, P^{(6)}(x)$  of the total gross claim amount  $Y + Y^-$  that are paid in the six positive payments are for simplicity assumed to be deterministic. Note that we could randomize these proportions by simulating from a Dirichlet distribution, but—as in Section 2.7.2—we refrain from doing so. Hence, we consider the proportion functions

$$P^{(k)}: \mathcal{X}_{7f} \rightarrow [0,1], \qquad x \mapsto P^{(k)}(x),$$

for all k = 1, ..., 6, with normalization  $\sum_{k=1}^{6} P^{(k)}(x) = 1$ , for all  $x \in \mathcal{X}_{7f}$ . We use the same model assumptions as in (7) by setting for k = 1, ..., 6

$$P^{(k)}(\mathbf{x}) = \frac{\exp{\{\mu_k(\mathbf{x})\}}}{\sum_{l=1}^{6} \exp{\{\mu_l(\mathbf{x})\}}},$$
(20)

for appropriate regression functions  $\mu_k : \mathcal{X}_{7f} \to \mathbb{R}$  resulting as output layer from a neural network with two hidden layers. As in (19), we consider the cross entropy loss function

$$\mathcal{L}(\boldsymbol{\alpha}) = \mathcal{L}_{\mathcal{D}_{7f}}(\boldsymbol{\alpha}) = -2 \sum_{i=1}^{n_{7f}} \sum_{k=1}^{6} P_i^{(k)} \log P^{(k)}(\boldsymbol{x}_i),$$

where  $\alpha$  is the corresponding network parameter. This model is calibrated as described in Section 2.1. Remark that if  $K^+ = 2$ , the model (20) simplifies to the binomial case, see (18).

## (c) Proportions of the Recovery Payments if $K^- = 2$

In the case of  $K^- = 2$  recovery payments, we need to model the proportion of the total individual recovery  $Y^-$  that is paid in the first recovery payment. For this, we work with the d = 10 dimensional feature space

$$\mathcal{X}_{7g} = \left\{ \left( \text{LoB}, \text{cc}, \text{AY}, \text{AQ}, \text{age}, \text{inj}_{\text{part}}, T, K, Y, Y^{-} \right) \right\}.$$

We denote by  $n_{7g} \le n_6$  the number of claims with exactly two recovery payments and order the claims appropriately in *i* such that  $K_i^- = 2$  for all  $i = 1, ..., n_{7g}$ . Recall that we set  $K_i^- = 2$  for all claims *i* with two or more recovery payments, see (15). Moreover, we add all the amounts of the recovery payments done after the second recovery payment to the second one. Let

$$R_i^- = \min\left\{ 0 \le j \le 11 \, \middle| \, C_i^{(j)} < 0 \right\}$$

denote the time point of the first recovery payment, for all  $i \in \{1, ..., n_{7g}\}$ . Then, the proportion  $P_i^-$  of the total individual recovery  $Y_i^-$  that is paid in the first recovery payment is given by

$$P_i^- = \frac{-C_i^{(R_i^-)}}{Y_i^-},$$

for all  $i = 1, ..., n_{7g}$ . This provides us with the data

$$\mathcal{D}_{7g} = \left\{ \left( \mathbf{x}_1, P_1^- \right), \dots, \left( \mathbf{x}_{n_{7g}}, P_{n_{7g}}^- \right) \right\}.$$

The remaining modeling part is then done completely analogously to the second part of the two positive payments case (a) in Section 2.7.2.

#### (d) Cash Flow Modeling

Finally, using the three neural network models outlined above, we can simulate the cash flow for a claim with more than two payments and with feature  $x \in X_7$ , see (5). We illustrate the case K = 6. Note that we only allow for cash flow patterns in  $\tilde{A}$  that are compatible with the reporting delay T. We start by describing the case T = 0. In this case, there is no difficulty and we directly simulate the cash flow pattern  $R(x) \in \tilde{A}$ . This provides us six payments in the time points

$$R^{(6)}(\mathbf{x}) = \max\left\{ 0 \le j \le 11 \left| 2^{j+1} \le R(\mathbf{x}) \right\} \text{ and } \\ R^{(k)}(\mathbf{x}) = \max\left\{ 0 \le j < R^{(k+1)}(\mathbf{x}) \left| 2^{j+1} \le R(\mathbf{x}) - \sum_{l=k+1}^{6} 2^{R^{(l)}(\mathbf{x})+1} \right\}, \text{ for } k = 1, \dots, 5.$$

For reporting delay T = 1, the set  $\tilde{A}$  of potential cash flow patterns becomes smaller because some of them have to be dropped to remain compatible with T = 1. For this reason, we simulate with probability  $\frac{1}{2}$  a pattern from  $\tilde{A}$ , and with probability  $\frac{1}{2}$  the six time points  $R^{(1)}(x), \ldots, R^{(6)}(x)$  are drawn in a uniform manner from the remaining possible time points in  $\{T = 1, \ldots, 11\}$ . For T > 1, the potential subset of patterns in  $\tilde{A}$  becomes (almost) empty. For this reason, we simply simulate uniformly from the compatible configurations in  $\{T, \ldots, 11\}$ .

Having the six time points for the payments, we distinguish the three different cases  $K^- \in \{0, 1, 2\}$ : *Case*  $K^- = 0$ : we calculate the proportions  $P^{(1)}(\mathbf{x}), \ldots, P^{(6)}(\mathbf{x})$  according to point (b) above and we receive the cash flow  $(C^{(j)}(\mathbf{x}))_{0 \le j \le 11}$  with

$$C^{(j)}(\boldsymbol{x}) = \begin{cases} P^{(l)}(\boldsymbol{x}) Y, & \text{if } j = R^{(l)}(\boldsymbol{x}) \text{ for some } 1 \le l \le 6, \\ 0, & \text{else.} \end{cases}$$

*Case*  $K^- = 1$ : we have five positive payments with proportions  $P^{(1)}(x), \ldots, P^{(5)}(x)$  modeled according to point (b) above. This provides the cash flow  $(C^{(j)}(x))_{0 \le j \le 11}$  with

$$C^{(j)}(\mathbf{x}) = \begin{cases} P^{(l)}(\mathbf{x}) \ (Y+Y^{-}), & \text{if } j = R^{(l)}(\mathbf{x}) \text{ for some } 1 \le l \le 5, \\ -Y^{-}, & \text{if } j = R^{(6)}(\mathbf{x}), \\ 0, & \text{else.} \end{cases}$$

*Case*  $K^- = 2$ : we have four positive payments with proportions  $P^{(1)}(x), \ldots, P^{(4)}(x)$  according to point (b) above and two negative payments with proportions  $P^-(x)$  and  $1 - P^-(x)$  according to point (c) above. The time point of the first recovery  $R^-(x)$  is simulated uniformly from the set of time points  $\{R^{(2)}(x), \ldots, R^{(5)}(x)\}$ . Note that the time point  $R^{(1)}(x)$  is reserved for the first positive payment and the time point  $R^{(6)}(x)$  for the second recovery payment. We write  $\tilde{R}^{(1)}(x), \ldots, \tilde{R}^{(4)}(x)$  for the time points of the four positive payments. Summarizing, we get the cash flow  $(C^{(j)}(x))_{0 \le j \le 11}$  with

$$C^{(j)}(\mathbf{x}) = \begin{cases} P^{(l)}(\mathbf{x}) \ (Y+Y^{-}), & \text{if } j = \widetilde{R}^{(l)}(\mathbf{x}) \text{ for some } 1 \le l \le 4, \\ -P^{-}(\mathbf{x}) Y^{-}, & \text{if } j = R^{-}(\mathbf{x}), \\ -(1-P^{-}(\mathbf{x})) Y^{-}, & \text{if } j = R^{(6)}(\mathbf{x}), \\ 0, & \text{else.} \end{cases}$$

Of course, if K = 3 and  $K^- = 2$ , we do not need to simulate the proportions of the positive payments, as there is only one positive payment, which occurs in the beginning. Similarly, if K = 12, we do not need to simulate the time points of the payments, since there is a payment in every development year.

## 2.8. Claim Status Modeling

Finally, we design the model for the claim status process which indicates whether a claim is open or closed at the end of each accounting year. This process modeling will also allow for re-opening. Similarly to the payments, we do not model the status of a claim or its changes within an accounting year, but only focus on its status at the end of each accounting year. The modeling procedure of the claim status uses two neural networks, which are described below.

We remark that the closing date information was of lower quality in our data set compared to all other information. For instance, some of the dates have been modified retrospectively which, of course, destroys the time series aspect. For this reason, we have decided to model this process in a more crude form, however, still capturing predictive power.

## 2.8.1. Re-Opening Indicator

We start by modeling the occurrence of a re-opening, i.e., whether a claim gets re-opened after having been closed at an earlier date. We use the d = 15 dimensional feature space

$$\mathcal{X}_{8a} = \left\{ \left( \text{LoB, AQ}, T, \left( \widetilde{C}^{(j)} \right)_{0 \le j \le 11} \right) \right\},$$
(21)

where we do not consider the exact payment amounts, but the simplified version

$$\widetilde{C}^{(j)} = \begin{cases} -\frac{1}{2}, & \text{if } C^{(j)} = 0, \\ 0, & \text{if } C^{(j)} \neq 0 \text{ and } C^{(j)} \leq 1,000, \\ \frac{1}{2}, & \text{if } C^{(j)} > 1,000, \end{cases}$$
(22)

for all j = 0, ..., 11. Let  $n_{8a} \le n$  denote the number of claims *i* for which we have the full information  $(I_i^{(j)})_{T_i \le j \le 11}$ . For the ease of data processing, we set  $I_i^{(j)} = 1$  for all development years before claims reporting  $T_i$ . Then, we can define the re-opening indicator  $V_i$  as

$$V_i = \begin{cases} 1, & \text{if } \sum_{j=1}^{11} \mathbb{1}_{\left\{I_i^{(j)} - I_i^{(j-1)} = 1\right\}} \ge 1, \\ 0, & \text{else,} \end{cases}$$

for all  $i = 1, ..., n_{8a}$ . In particular, if  $V_i = 1$ , then claim *i* has at least one re-opening, and if  $V_i = 0$ , then claim *i* has not been re-opened. This leads us to the data

$$\mathcal{D}_{8a} = \{(x_1, V_1), \dots, (x_{n_{8a}}, V_{n_{8a}})\},\$$

where  $x_1, \ldots, x_{n_{8a}} \in \mathcal{X}_{8a}$ . For a given feature  $x \in \mathcal{X}_{8a}$ , the corresponding re-opening indicator V(x) is a Bernoulli random variable. Thus, model calibration is done analogously to Section 2.2 with, however, a neural network with only one hidden layer.

## 2.8.2. Closing Delay Indicator for Claims without a Re-Opening

For claims without a re-opening, we model the closing delay indicator determining whether the closing occurs in the same year as the last payment or if the closing occurs later. In case of no payments  $(Z_i = 0)$ , we replace the year of the last payment by the reporting year. We use the same d = 15 dimensional feature space as for the re-opening indicator and set  $\mathcal{X}_{8b} = \mathcal{X}_{8a}$ , see (21). Let  $n_{8b} \le n_{8a}$  be the number of claims without a re-opening and order them appropriately in *i* such that  $V_i = 0$  for all  $i = 1, \ldots, n_{8b}$ . Then, we define the closing delay indicator  $W_i$  as

$$W_{i} = \begin{cases} 1, & \text{if } Z_{i} = 1 \text{ and } \max\left\{0 \le j \le 11 \mid I_{i}^{(j)} = 1\right\} \ge \max\left\{0 \le j \le 11 \mid C_{i}^{(j)} \ne 0\right\}, \\ 1, & \text{if } Z_{i} = 0 \text{ and } \max\left\{0 \le j \le 11 \mid I_{i}^{(j)} = 1\right\} \ge T_{i}, \\ 0, & \text{else}, \end{cases}$$

for all  $i = 1, ..., n_{8b}$ . Hence, we have  $W_i = 1$  if the closing occurs in a later year compared to the year of the last payment (or in a later year compared to the claims reporting year in case there is no payment) and  $W_i = 0$  otherwise. This leads us to the data

$$\mathcal{D}_{8b} = \{(x_1, W_1), \dots, (x_{n_{8b}}, W_{n_{8b}})\}.$$

For a claim with feature  $x \in X_{8b}$ , the corresponding closing delay indicator W(x) is again a Bernoulli random variable. Therefore, model calibration is done analogously to Section 2.2. Similarly as for the re-opening indicator, we use a neural network with only one hidden layer.

#### 2.8.3. Simulation of the Claim Status

Based on the feature space  $\mathcal{X}_{8a}$ , we first simulate the re-opening indicator  $V(\mathbf{x})$  leading to the two cases (i) and (ii) described below. Note that before a claim is reported—for ease of data processing—we simply set its status to open (this has no further relevance).

(i) Case V(x) = 0 (without re-opening): for the given feature  $x \in \mathcal{X}_{8a}$ , we calculate the closing delay probability  $\pi(x) = \mathbb{P}[W(x) = 1]$  using the neural network of Section 2.8.2. The closing delay B(x) is then sampled from a categorical distribution on  $\{0, \ldots, 12\}$  with probabilities

$$\mathbb{P}[B(\mathbf{x}) = 0] = 1 - \pi(\mathbf{x}),$$
  

$$\mathbb{P}[B(\mathbf{x}) = 1] = \frac{9}{10}\pi(\mathbf{x}),$$
  

$$\mathbb{P}[B(\mathbf{x}) = k] = \frac{1}{10}\frac{1}{11}\pi(\mathbf{x}), \quad \text{for } k = 2, \dots, 12.$$

The resulting closing delay  $B(x) \in \{0, ..., 12\}$  is added to the year of the last payment (or to the reporting year if there is no payment). If this sum exceeds the value 11, the claim is still open at the end of the last modeled development year. This provides the claim status process  $(I^{(j)}(x))_{0 \le j \le 11}$  with

$$I^{(j)}(\mathbf{x}) = \begin{cases} 1, & \text{if } Z = 1 \text{ and } j < B(\mathbf{x}) + \max\left\{0 \le j \le 11 \mid C^{(j)} \ne 0\right\}, \\ 1, & \text{if } Z = 0 \text{ and } j < B(\mathbf{x}) + T, \\ 0, & \text{else.} \end{cases}$$

(ii) Case V(x) = 1 (with re-opening): if we have at least one payment for the considered claim, then the first settlement time  $B_1(x)$  is simulated from a uniform distribution on the set

$$\left\{T,\ldots,\max\left\{0\leq j\leq 11 \mid C^{(j)}\neq 0\right\}\right\}.$$

The second settlement time  $B_2(x)$  is simulated from a uniform distribution on the set

$$\left\{\max\left\{0\leq j\leq 11\,|\,C^{(j)}\neq 0\right\}+2,\ldots,13\right\}.$$

In particular, the first settlement arrives between the reporting year and the year of the last payment. Then, the claim gets re-opened in the year following the first settlement. The second settlement, if there is one, arrives between two years after the year of the last payment and the last modeled development year. In case the second settlement arrives after the last modeled development year, we simply cannot observe it and the claim is still open at the end of the last modeled development year. In case the first settlement happens in the last modeled development year, we do not even observe the re-opening.

If the claim does not have any payment, we set  $B_1(x) = T$  for the first settlement time. In particular, the claim gets closed for the first time in the same year as it is reported. The second settlement time  $B_2(x)$  is simulated from a uniform distribution on the set  $\{T + 2, ..., 13\}$ .

This leads to the claim status process  $(I^{(j)}(\mathbf{x}))_{0 \le j \le 11}$  with

$$I^{(j)}(\mathbf{x}) = \begin{cases} 1, & \text{if } j < B_1(\mathbf{x}) \text{ or } B_1(\mathbf{x}) < j < B_2(\mathbf{x}), \\ 0, & \text{if } j = B_1(\mathbf{x}) \text{ or } j \ge B_2(\mathbf{x}). \end{cases}$$

## 3. Model Calibration Using Momentum-Based Gradient Descent

In Section 2 we have introduced several neural networks that need to be calibrated to the data. This calibration involves the choice of the numbers of hidden neurons  $q_1$  and  $q_2$  as well as the choice of the corresponding network parameter  $\alpha$ . We first focus on the network parameter  $\alpha$  for given  $q_1$  and  $q_2$ .

#### 3.1. Gradient Descent Methods

State-of-the-art for finding the optimal network parameter  $\alpha$  w.r.t. a given differentiable loss function  $\alpha \mapsto \mathcal{L}(\alpha)$  is the gradient descent method (GDM). The GDM locally improves the loss in an iterative way. Consider the Taylor approximation of  $\mathcal{L}$  around  $\alpha$ , then

$$\mathcal{L}(\widetilde{\boldsymbol{\alpha}}) = \mathcal{L}(\boldsymbol{\alpha}) + (\nabla_{\boldsymbol{\alpha}} \mathcal{L}(\boldsymbol{\alpha}))' (\widetilde{\boldsymbol{\alpha}} - \boldsymbol{\alpha}) + o(\|\widetilde{\boldsymbol{\alpha}} - \boldsymbol{\alpha}\|),$$

as  $\|\tilde{\alpha} - \alpha\| \to 0$ . The locally optimal move points into the direction of the negative gradient  $-\nabla_{\alpha} \mathcal{L}(\alpha)$ . If we choose a learning rate  $\varrho > 0$  into that direction, we obtain a local loss decrease

$$\mathcal{L}\left(\boldsymbol{\alpha} - \varrho \nabla_{\boldsymbol{\alpha}} \mathcal{L}(\boldsymbol{\alpha})\right) \approx \mathcal{L}(\boldsymbol{\alpha}) - \varrho \|\nabla_{\boldsymbol{\alpha}} \mathcal{L}(\boldsymbol{\alpha})\|^{2}, \qquad (23)$$

for  $\varrho$  small. Iterative application of these locally optimal moves—with tempered learning rates—will converge ideally to the (local) minimum of the loss function. Note that (a) it is possible to end up in saddle points; (b) different starting points of this algorithm should be explored to see whether we converge to different (local) minima resp. saddle points and (c) the speed of convergence should be

18 of 32

fine-tuned. An improved version of the GDM is the so-called momentum-based GDM introduced in Rumelhart et al. (1986). Consider a velocity vector **v** with the same dimensions as  $\alpha$  and initialize **v** = **0**, corresponding to zero velocity in the beginning. Then, in every iteration step of the GDM, we are building up velocity to achieve a faster convergence. In formulas, this provides

$$\mathbf{v} \leftarrow \mu \mathbf{v} - \varrho \, \nabla_{\boldsymbol{\alpha}} \mathcal{L}(\boldsymbol{\alpha}) \\ \boldsymbol{\alpha} \leftarrow \boldsymbol{\alpha} + \mathbf{v},$$

where  $\mu \in [0, 1]$  is the momentum coefficient controlling how fast velocity is built up. By choosing  $\mu = 0$ , we get the original GDM without a velocity vector, see (23). Fine-tuning  $0 < \mu \le 1$  may lead to faster convergence. We refer to the relevant literature for more on this topic.

#### 3.2. Gradients of the Loss Functions Involved

In Section 2 we have met three different model types of neural networks:

- categorical case with more than two categorical classes;
- Bernoulli case with exactly two categorical classes;
- log-normal case.

In order to apply the momentum-based GDMs, we need to calculate the gradients of the corresponding loss functions of these three model types. As illustrations, we choose the reporting delay T for the categorical case, the payment indicator Z for the Bernoulli case and the total individual claim size Y for the log-normal case.

## 3.2.1. Categorical Case (with More than Two Categorical Classes)

The loss function  $\alpha \mapsto \mathcal{L}(\alpha)$  for the modeling of the reporting delay *T* is given in (8). The gradient  $\nabla_{\alpha} \mathcal{L}(\alpha)$  can be calculated as

$$\nabla_{\boldsymbol{\alpha}} \mathcal{L}(\boldsymbol{\alpha}) = -2 \sum_{i=1}^{n_1} \sum_{t \in \mathcal{T}} \mathbb{1}_{\{T_i = t\}} \nabla_{\boldsymbol{\alpha}} \log \pi_t(\boldsymbol{x}_i) = -2 \sum_{i=1}^{n_1} \sum_{t \in \mathcal{T}} \mathbb{1}_{\{T_i = t\}} \frac{1}{\pi_t(\boldsymbol{x}_i)} \nabla_{\boldsymbol{\alpha}} \pi_t(\boldsymbol{x}_i).$$

We have for the last gradients

$$\nabla_{\boldsymbol{\alpha}} \pi_t(\boldsymbol{x}_i) = \nabla_{\boldsymbol{\alpha}} \frac{\exp\left\{\mu_t(\boldsymbol{x}_i)\right\}}{\sum_{s \in \mathcal{T}} \exp\left\{\mu_s(\boldsymbol{x}_i)\right\}} = \pi_t(\boldsymbol{x}_i) \left(\nabla_{\boldsymbol{\alpha}} \mu_t(\boldsymbol{x}_i) - \sum_{s \in \mathcal{T}} \pi_s(\boldsymbol{x}_i) \nabla_{\boldsymbol{\alpha}} \mu_s(\boldsymbol{x}_i)\right),$$

for all  $t \in \mathcal{T}$  and  $i = 1, ..., n_1$ . Collecting all terms, we conclude

$$\begin{aligned} \nabla_{\boldsymbol{\alpha}} \mathcal{L}(\boldsymbol{\alpha}) &= -2 \sum_{i=1}^{n_1} \sum_{t \in \mathcal{T}} \mathbb{1}_{\{T_i = t\}} \left( \nabla_{\boldsymbol{\alpha}} \mu_t(\boldsymbol{x}_i) - \sum_{s \in \mathcal{T}} \pi_s(\boldsymbol{x}_i) \nabla_{\boldsymbol{\alpha}} \mu_s(\boldsymbol{x}_i) \right) \\ &= -2 \sum_{i=1}^{n_1} \sum_{t \in \mathcal{T}} \left( \mathbb{1}_{\{T_i = t\}} - \pi_t(\boldsymbol{x}_i) \right) \nabla_{\boldsymbol{\alpha}} \mu_t(\boldsymbol{x}_i). \end{aligned}$$

There remains to calculate the gradients  $\nabla_{\alpha} \mu_t(\mathbf{x}_i)$ , for all  $t \in \mathcal{T}$  and  $i = 1, ..., n_1$ . This is done using the back-propagation algorithm, which in today's form goes back to Werbos (1982).

#### 3.2.2. Bernoulli Case (Two Categorical Classes)

We calculate the gradient  $\nabla_{\alpha} \mathcal{L}(\alpha)$  for the modeling of the payment indicator *Z* with corresponding loss function  $\mathcal{L}(\alpha)$  given in (11). We get as in the categorical case above

$$\nabla_{\boldsymbol{\alpha}} \mathcal{L}(\boldsymbol{\alpha}) = -2 \sum_{i=1}^{n_2} \left( \frac{Z_i}{\pi(\boldsymbol{x}_i)} - \frac{1-Z_i}{1-\pi(\boldsymbol{x}_i)} \right) \nabla_{\boldsymbol{\alpha}} \pi(\boldsymbol{x}_i),$$

Risks 2018, 6, 29

with gradient

$$\nabla_{\boldsymbol{\alpha}} \pi(\boldsymbol{x}_i) = \frac{\exp\left\{-\mu(\boldsymbol{x}_i)\right\}}{\left(1 + \exp\left\{-\mu(\boldsymbol{x}_i)\right\}\right)^2} \nabla_{\boldsymbol{\alpha}} \mu(\boldsymbol{x}_i) = \pi(\boldsymbol{x}_i) \left(1 - \pi(\boldsymbol{x}_i)\right) \nabla_{\boldsymbol{\alpha}} \mu(\boldsymbol{x}_i),$$

for all  $i = 1, ..., n_2$ . Collecting all terms, we obtain

$$abla_{\boldsymbol{\alpha}} \mathcal{L}(\boldsymbol{\alpha}) = -2 \sum_{i=1}^{n_2} (Z_i - \pi(\boldsymbol{x}_i)) \nabla_{\boldsymbol{\alpha}} \mu(\boldsymbol{x}_i).$$

We again apply back-propagation to calculate the gradient  $\nabla_{\alpha}\mu(x_i)$ , for all  $i = 1, ..., n_2$ .

#### 3.2.3. Log-Normal Case

Finally, the loss function  $\mathcal{L}(\alpha)$  for the modeling of the total individual claim size *Y* is given in (14). Hence, for the gradient  $\nabla_{\alpha} \mathcal{L}(\alpha)$ , we have

$$\nabla_{\boldsymbol{\alpha}} \mathcal{L}(\boldsymbol{\alpha}) = \sum_{i=1}^{n_4} \nabla_{\boldsymbol{\alpha}} \left( \log Y_i - \mu(\boldsymbol{x}_i) \right)^2 = -2 \sum_{i=1}^{n_4} \left( \log Y_i - \mu(\boldsymbol{x}_i) \right) \nabla_{\boldsymbol{\alpha}} \mu(\boldsymbol{x}_i),$$

where the last gradient  $\nabla_{\alpha} \mu(\mathbf{x}_i)$ , for all  $i = 1, ..., n_4$ , is again calculated using back-propagation.

## 3.3. Choice of the Numbers of Hidden Neurons

For each modeling step of our simulation machine, we still need to determine the optimal neural network in terms of the numbers  $q_1$  and  $q_2$  of hidden neurons. These hyperparameters are determined by splitting the original data set into a training set and a validation set, where for each calibration we choose at random 90% of the data for the training set. The training set is then used to fit the models for the different choices of hyperparameters  $q_1$  and  $q_2$  by minimizing the corresponding (training) in-sample losses of the functions  $\alpha \mapsto \mathcal{L}(\alpha)$ . This is done as described in the previous sections—for given  $q_1$  and  $q_2$ . The hyperparameter choices  $q_1$  and  $q_2$ —and model choices, respectively—are then done by choosing the model with the smallest (validation) out-of-sample loss on the validation set.

## 4. Chain-Ladder Analysis

In this section we use the calibrated stochastic simulation machine to perform a small claims reserving analysis. We generate data from the simulation machine and compare it to the real data. For both data sets, we analyze the resulting claims reporting patterns and the corresponding claims cash flow patterns. For claims reportings, we separate the individual claims i = 1, ..., n by accident year AY  $\in$  {1994, ..., 2005} and reporting delays  $T \in \{0, ..., 11\}$ . For claims cash flows, we separate the individual claims i = 1, ..., n again by accident year AY  $\in$  {1994, ..., 2005} and aggregate the corresponding payments over the development delays j = 0, ..., 11. The reported claims and the claims payments that are available by the end of accounting year 2005 then provide the so-called upper claims reserving triangles. These triangles of reported claims of real and simulated data are shown in Tables 2 and 3, the triangles of cumulative claims payments of real and simulated data are given in Tables 4 and 5. At a first glance, these triangles show that the simulated data looks very similar to the real data, with a slightly bigger similarity for claims reportings than for claims cash flows.

These data sets can be used to perform a chain-ladder (CL) claims reserving analysis. We therefore use Mack's chain-ladder model, for details we refer to Mack (1993). We calculate the chain-ladder reserves for both the real and the simulated data, and we also calculate Mack's square-rooted conditional mean square error of prediction  $\sqrt{\text{msep}}$ .

We start the analysis on the claims reportings. Using the chain-ladder method, we predict the number of incurred but not yet reported (IBNYR) claims. These are the predicted numbers of late reported claims in the lower triangles in Tables 2 and 3. The resulting predictions are provided in the

2nd and 5th columns of Table 6. We observe a high similarity between the results on the real and the simulated data. In particular, for all the individual accident years, the chain-ladder predicted numbers of IBNYR claims of the real data and the simulated data are very close to each other. Aggregating over all accident years, the chain-ladder predicted number of the total IBNYR claims is only 0.2% higher for the simulated data compared to the real data. This similarity largely carries over to the prediction uncertainty analysis illustrated by the columns  $\sqrt{\text{msep}}$  in Table 6. Indeed, comparing the real and the simulated data, we see that  $\sqrt{\text{msep}}$  is of similar magnitude for most accident years. Only for the accident years 2003 and 2004 it seems notably higher for the real data. From this, we conclude that, at least from a chain-ladder reserving point of view, our stochastic simulation machine provides very reasonable claims reporting patterns.

Finally, Table 7 shows the results of the chain-ladder analysis for claims payments. Columns 2 and 5 of that table provide the chain-ladder reserves. These are the payment predictions for the cash flows paid after accounting year 2005 and complete the lower triangles in Tables 4 and 5. Also here we see high similarities between the real data and the simulated data analysis: the corresponding total chain-ladder reserves as well as the corresponding reserves for most of the individual accident years are rather close to each other. In particular, the total chain-ladder reserves are only 1.2% higher for the simulated data. We only observe slightly shorter cash flow patterns in the simulated data, which partially carries over to the prediction uncertainties illustrated by the columns  $\sqrt{msep}$  in Table 7.

Accident				Rep	orting	Delay	T					
Year AY	0	1	2	3	4	5	6	7	8	9	10	11
1994	861,899	59,056	1540	460	230	154	84	56	50	28	32	12
1995	850,297	64,733	1568	562	216	124	94	62	44	34	32	
1996	781,875	61,465	1742	414	252	153	76	62	38	22		
1997	756,147	59,269	1466	496	210	147	54	48	40			
1998	753,552	60,249	1660	530	248	136	98	44				
1999	754,992	59,690	1625	468	208	100	68					
2000	766,684	61,120	1274	320	136	88						
2001	758,443	61,449	1024	286	90							
2002	745,125	55,246	876	200								
2003	757,843	53,272	956									
2004	733,785	51,742										
2005	730,978											

Table 2. Triangle of reported claims of the real data.

Table 3. Triangle of reported claims of the simulated data.

Accident				Rep	orting	g Delay	уT					
Year AY	0	1	2	3	4	5	6	7	8	9	10	11
1994	860,337	60,143	1837	553	256	155	101	70	57	35	36	21
1995	851,877	62,924	1776	499	263	164	80	58	52	29	36	
1996	783,006	60,511	1557	477	182	122	87	49	53	39		
1997	756,015	59,556	1434	399	169	129	77	42	42			
1998	752,454	61,913	1422	380	164	100	65	37				
1999	753,635	61,552	1279	362	153	116	49					
2000	768,180	59,636	1215	338	150	71						
2001	759,501	60,131	1124	304	132							
2002	744,478	55 <i>,</i> 577	1012	270								
2003	757,635	53,352	937									
2004	732,884	52,586										
2005	731,357											

Accident	Development Delay <i>j</i>											
Year AY	0	1	2	3	4	5	6	7	8	9	10	11
1994	78,433	120,396	130,167	134,749	137,143	138,798	139,994	141,052	142,106	142,913	143,652	144,247
1995	79,372	124,532	135,488	140,338	143,145	144,859	146,408	147,624	148,735	149,686	150,578	
1996	71,398	113,335	123,336	128,052	130,779	132,692	134,175	135,359	136,381	137,311		
1997	68,600	107,716	117,556	122,331	125,304	127,477	129,037	130,171	131,216			
1998	68,055	109,906	120,706	126,443	129,727	131,998	133,745	135,176				
1999	71,989	114,344	127,311	134,174	138,075	140,614	142,530					
2000	72,225	118,418	131,615	138,216	142,024	144,450						
2001	74,891	126,244	141,008	147,923	151,917							
2002	78,167	129,105	143,731	150,560								
2003	82,668	134,010	148,161									
2004	80,630	130,390										
2005	82,015											

Table 4. Triangle of cumulative claims payments (in 10,000 CHF) of the real data.

 Table 5. Triangle of cumulative claims payments (in 10,000 CHF) of the simulated data.

Accident						Developn	ent Delay	j				
Year AY	0	1	2	3	4	5	6	7	8	9	10	11
1994	80,491	117,807	129,673	135,331	138,591	140,745	142,268	143,394	144,344	145,090	145,694	146,076
1995	79,170	116,943	129,313	135,330	138,785	141,079	142,743	144,045	145,115	146,002	146,760	
1996	71,675	107,228	119,578	125,407	128,771	131,009	132,566	133,718	134,657	135,472		
1997	68,857	104,291	116,406	122,177	125,495	127,662	129,194	130,388	131,270			
1998	67,418	103,415	116,194	122,262	125,673	127,877	129,452	130,588				
1999	69,308	107,587	121,381	127,930	131,534	133,924	135,502					
2000	73,359	113,266	127,878	134,804	138,806	141,431						
2001	73,338	115,626	131,197	138,481	142,591							
2002	74,887	117,602	133,850	141,539								
2003	81,921	127,461	144,444									
2004	81,394	128,864										
2005	86,837											

**Table 6.** Chain-ladder predicted numbers of incurred but not yet reported (IBNYR) claims and Mack's  $\sqrt{\text{msep}}$  for the real and the simulated data.

Accident	CL Predicted	$\sqrt{msep}$	in %	CL Predicted	$\sqrt{msep}$	in %
Year AY	IBNYR Claims			IBNYR Claims		
	Real	Data		Simula	ted Data	
1994	0			0		
1995	12	0	0.0%	21	0	0.0%
1996	40	0	0.4%	52	0	0.4%
1997	65	5	8.4%	82	7	8.6%
1998	105	7	6.2%	129	9	6.6%
1999	156	10	6.1%	178	14	7.7%
2000	235	19	8.0%	255	21	8.3%
2001	357	32	9.0%	370	35	9.4%
2002	536	65	12.2%	535	53	9.9%
2003	944	135	14.3%	925	95	10.2%
2004	2201	330	15.0%	2170	265	12.2%
2005	57,734	3542	6.1%	57,789	3410	5.9%
total	62,385	3565	5.7%	62,506	3425	5.5%

Accident	CL Reserves	$\sqrt{msep}$	in %	CL Reserves	√msep	in %
Year AY						
	Re	al Data		Simul	lated Data	
1994	0			0		
1995	624	117	18.8%	385	109	28.2%
1996	1337	146	10.9%	991	161	16.2%
1997	2112	168	8.0%	1723	179	10.4%
1998	3224	177	5.5%	2636	187	7.1%
1999	4686	259	5.5%	3943	209	5.3%
2000	6476	394	6.1%	5826	230	3.9%
2001	9275	599	6.5%	8846	290	3.4%
2002	13,049	889	6.8%	12,489	421	3.4%
2003	19,973	1421	7.1%	20,817	837	4.0%
2004	32,532	2394	7.4%	36,647	2050	5.6%
2005	82,706	5039	6.1%	84,175	4968	5.9%
total	175,994	6275	3.6%	178,076	5732	3.2%

**Table 7.** Chain-ladder reserves for claims payments (in 10,000 CHF) and Mack's  $\sqrt{\text{msep}}$  for the real and the simulated data.

## 5. Conclusions

We have developed a stochastic simulation machine that generates individual claims histories of non-life insurance claims. This simulation machine is based on neural networks which have been calibrated to real non-life insurance data. The inputs of the simulation machine are a portfolio of non-life insurance claims—for which we want to simulate the corresponding individual claims histories—and the two variance parameters  $\sigma^2_+$  (for the total individual claim size, see (13)) and  $\sigma^2_-$  (for the total individual recovery, see Section 2.6). Together with a portfolio generating algorithm, see Appendix B, one can use this simulation machine to simulate as many individual claims histories as desired. In a chain-ladder analysis we have seen that the simulation machine leads to reasonable results, at least from a chain-ladder reserving point of view. Therefore, our simulation machine may serve as a stochastic scenario generator for individual claims histories, which provides a common ground for research in this area, we also refer to the study in Wüthrich (2018b).

Acknowledgments: Our greatest thanks go to Suva, Peter Blum and Olivier Steiger for providing data, their insights and for their immense support.

Author Contributions: Both authors contributed equally to this work.

Conflicts of Interest: The authors declare no conflict of interest.

#### Appendix A. Descriptive Statistics of the Chosen Data Set

In this appendi we provide descriptive statistics of the data used to calibrate the individual claims history simulation machine. For confidentiality reasons, we can only show aggregate statistics of the claims portfolio, see Figures A1–A4 below.



Figure A1. Portfolio distributions w.r.t. the features (a) LoB; (b) cc; (c) AY; (d) AQ; (e) age and (f) inj\_part.



**Figure A2.** (a) Logarithmic number of claims; (b) average claim size and (c) average number of payments w.r.t. the reporting delay *T*; the red lines show the averages.



**Figure A3.** (a) Logarithmic number of claims; (b) average claim size and (c) number of claims with recoveries w.r.t. the number of payments *K*; the red lines show the averages.



Figure A4. Average claim size w.r.t. the features (a) LoB; (b) cc; (c) AY; (d) AQ; (e) age and (f) inj\_part; the red lines show the averages.

# Appendix B. Procedure of Generating a Synthetic Portfolio

In order to use the stochastic simulation machine derived above, we require a portfolio of features  $x_1, \ldots, x_n \in \mathcal{X}_1$ , see (1). Therefore, we need an additional scenario generator that simulates reasonable synthetic portfolios. In this appendix we describe the design of our portfolio scenario generator which provides portfolios similar in structure to the original portfolio.

Our algorithm of synthetic portfolio generation uses the following input parameters:

• V = totally expected number of claims;

- $(p_l)_{1 \le l \le 4}$  = categorical distribution for the allocation of the claims to the four lines of business;
- $(r_l)_{1 \le l \le 4}$  = growth parameters for the numbers of claims in the 12 accident years for each of the four lines of business.

In a first step, we use these parameters to simulate the total number of claims and allocate them to the lines of business LoB and the accident years AY. We start by simulating  $(V_l)_{1 \le l \le 4}$  according to

$$(V_l)_{1 \le l \le 4} \sim \text{Multinomial}(V, (p_l)_{1 \le l \le 4})$$

To determine the distribution of the claims among the 12 accident years within each line of business l = 1, ..., 4, we simulate  $(X_j^{(l)})_{1 \le l \le 4, 2 \le j \le 12}$  from a normal distribution according to

$$X_j^{(l)} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(r_l, r_l^2).$$

Then, we define the weights  $W_1^{(l)} = 1$  and

$$W_{j}^{(l)} = W_{j-1}^{(l)} \exp \left\{ X_{j}^{(l)} \right\},$$

for all  $l = 1, \dots, 4$  and  $j = 2, \dots, 12$ . Finally, we set

$$V_{l,j} = V_l \frac{W_j^{(l)}}{\sum_{i'=1}^{12} W_{i'}^{(l)}},$$

for all l = 1, ..., 4 and j = 1, ..., 12, to be the expected number of claims in line of business l with accident year j. Conditionally given  $\mathbf{V} = (V_{1,1}, V_{1,2}, ..., V_{4,12})$ , we simulate the number of claims  $N_{l,j}$  in line of business l with accident year j from a Poisson distribution according to

$$N_{l,j} \mid \mathbf{V} \stackrel{\text{ind.}}{\sim} \operatorname{Poi}(V_{l,j}),$$

for all l = 1, ..., 4 and j = 1, ..., 12. Note that we have  $\mathbb{E}[\sum_{j=1}^{12} N_{l,j}] = V p_l$ , which justifies the above modeling choices.

After having simulated the number of claims  $N_{l,j}$  for each line of business l and accident year j, we need to establish these claims with the remaining feature components cc, AQ, age and inj\_part. This is achieved by choosing a multivariate distribution having a Gaussian copula and appropriate marginal densities. These densities and the covariance parameters of the Gaussian copula have been estimated from the real data. For the explicit parametrization, we refer to the R-function Feature.Generation in our simulation package.

## Appendix C. Sensitivities of Selected Neural Networks

In this final appendix we consider 11 selected neural networks of our simulation machine and present the impact on the response variable of the respective most influential features. For each neural network considered, we use the corresponding calibration data set, fix a feature component—e.g., the accident quarter AQ—and vary its value over its entire domain—e.g.,  $\{1, \ldots, 4\}$  for the accident quarter AQ—to analyze the sensitivities in this feature component.

In Figure A5 we analyze the reporting delay T as a function of the features AQ, age and inj\_part. Not surprisingly, the accident quarter has the biggest influence, because a claim occurring in December is likely to be reported only in the next accounting year.



Figure A5. Reporting delay T w.r.t. the features (a) AQ; (b) age and (c) inj\_part.

Figure A6 tells us that claims in lines of business one and four almost always have a payment. In contrast, we expect only roughly half of the claims in lines of business two and three to have a payment. Furthermore, the claims code cc causes some variation in the probability of having a payment, and claims with either a small or a large reporting delay *T* have a higher probability of having a payment than claims with a medium reporting delay.



Figure A6. Payment indicator Z w.r.t. the features (a) LoB; (b) cc and (c) reporting delay T.

Recall that in determining the number of payments K, we use two neural networks, where in the first one we model whether we have K = 1 or K > 1 payments. According to Figure A7, claims that occur later in a year tend to have a higher probability of having more than one payment. The same holds true with increasing age of the injured. In passing from reporting delay T = 0 to T = 1, the probability of having only one payment increases. But then we observe a sinus curve shape in that probability as a function of T.

The second neural network used to determine the number of payments K models the distribution of K, conditioned on K > 1. As we see in Figure A8, claims in line of business two tend to have more payments than claims in other lines of business, and both inj\_part and reporting delay T heavily influence the number of payments.



**Figure A7.** Indicator whether we have K = 1 or K > 1 payments w.r.t. the features (**a**) AQ; (**b**) age and (**c**) reporting delay *T*.



**Figure A8.** Conditional distribution of the number of payments *K*, given K > 1, w.r.t. the features (a) LoB; (b) inj\_part and (c) reporting delay *T*.

In Figure A9 we present sensitivities for the expected total individual claim size *Y* on the log scale. The main drivers here are the line of business LoB and the number of payments *K*.



**Figure A9.** Total individual claim size *Y* (on log scale) w.r.t. the features (**a**) LoB; (**b**) reporting delay *T* and (**c**) number of payments *K*.

Figure A10 tells us that claims in lines of business one and four almost never have a recovery. Moreover, the probability of having at least one recovery payment first increases with the number of payments *K* but then slightly decreases again. Finally, up to 50% of the claims with a small total individual claim size Y (of less than 10 CHF) have a recovery. This also comprises claims whose recovery is almost equal to the total gross claim amount, leading to a small net claim size. In general, the higher the total individual claim size, the less likely are recovery payments.



**Figure A10.** Number of recovery payments  $K^-$  w.r.t. the features (**a**) LoB; (**b**) number of payments K and (**c**) total individual claim size Y.

According to Figure A11, the total individual recovery  $Y^-$  is substantially higher for claims in lines of business two and three, compared to claims in lines of business one and four. Furthermore, if we have a recovery, then the higher the number of payments *K* and the total individual claim size *Y*, the higher also the recovery, where the increase w.r.t. the number of payments is decisively more pronounced.



**Figure A11.** Total individual recovery  $Y^-$  w.r.t. the features (**a**) LoB; (**b**) number of payments *K* and (**c**) total individual claim size *Y*.

In determining the payment delay *S* for claims with exactly one payment, we use two neural networks. In the first one, we model whether S = 0 or S > 0, and in the second one, we consider the conditional distribution of *S*, given S > 0. Here we only present sensitivities for the first neural network. We observe, see Figure A12, that the probability of a payment delay equal to zero decreases with increasing accident quarter AQ and increasing total individual claim size *Y*. In particular, claims that occur in the last quarter of a year have a considerably higher probability of having a payment delay. This might be explained by claims for which the short time lag between the accident date and

the end of the year only suffices for claims reporting but not for claims payments, leading to a payment delay. Finally, claims with a reporting delay T > 0 almost never have an additional payment delay.



**Figure A12.** Indicator whether we have payment delay S = 0 or S > 0 in the case of K = 1 payment w.r.t. the features (a) AQ; (b) reporting delay *T* and (c) total individual claim size *Y*.

As a representative of the neural networks that calculate the proportions with which the total gross claim amount  $Y + Y^-$  is distributed among the  $K^+$  positive payments, we choose the one for  $K^+ = 6$ . According to Figure A13, we see some monotonicity, but apart from that these proportions do not vary considerably. For claims which occur early during a year or have a high reporting delay *T* or a comparably small total individual claim size *Y*, the biggest proportion of the total gross claim amount is paid in the first (positive) payment.



**Figure A13.** Proportions  $P^{(1)}, \ldots, P^{(6)}$  of the total gross claim amount  $Y + Y^-$  paid in the  $K^+ = 6$  positive payments w.r.t. the features (**a**) AQ; (**b**) reporting delay *T* and (**c**) total individual claim size *Y*.

According to Figure A14, in the case of  $K^- = 2$  recovery payments, the proportion  $P^-$  of the total individual recovery  $Y^-$  that is paid in the first recovery payment varies substantially for the different values of the features cc and inj\_part. We also observe that the higher the total individual recovery, the higher the proportion paid in the first recovery.



**Figure A14.** Proportion  $P^-$  of the total individual recovery  $Y^-$  paid in the first recovery payment in the case of  $K^- = 2$  recovery payments w.r.t. the features (**a**) cc; (**b**) inj\_part and (**c**) total individual recovery  $Y^-$ .

In Figure A15 we see that claims in lines of business one and four have a higher re-opening probability than claims in lines of business two and three. Moreover, the higher the reporting delay *T* of a claim, the lower the rate of reopening. Finally, the probability of re-opening heavily depends on the cash flow. In order to not overload the plot, we only show sensitivities w.r.t. the payments  $C^{(0)}, C^{(1)}, C^{(2)}, C^{(3)}, C^{(8)}, C^{(10)}$ . Recall that for this neural network, the yearly payments are coded with the values  $-\frac{1}{2}$ , 0 and  $\frac{1}{2}$ , see (22). Summarizing, one can say that if we have a payment after the first development year, then the probability of re-opening is quite high.



**Figure A15.** Re-opening indicator *V* w.r.t. the features (**a**) LoB; (**b**) reporting delay *T* and (**c**) yearly payments  $C^{(0)}$ ,  $C^{(1)}$ ,  $C^{(2)}$ ,  $C^{(3)}$ ,  $C^{(8)}$ ,  $C^{(10)}$ .

# References

- Antonio, Katrien, and Richard Plat. 2014. Micro-Level Stochastic Loss Reserving for General Insurance. *Scandinavian Actuarial Journal* 7: 649–69.
- Cybenko, George. 1989. Approximation by Superpositions of a Sigmoidal Function. *Mathematics of Control, Signals, and Systems (MCSS)* 2: 303–14.
- Hiabu Munir, Carolin Margraff, Maria D. Martínez-Miranda, and Jens P. Nielsen. 2016. The Link between Classical Reserving and Granular Reserving through Double Chain-Ladder and its Extensions. *British Actuarial Journal* 21: 97–116.
- Hornik, Kurt, Maxwell Stinchcombe, and Halbert White. 1989. Multilayer Feedforward Networks are Universal Approximators. *Neural Networks* 2: 359–66.

- Jessen, Anders H., Thomas Mikosch, and Gennady Samorodnitsky. 2011. Prediction of Outstanding Payments in a Poisson Cluster Model. *Scandinavian Actuarial Journal* 3: 214–37.
- Mack, Thomas. 1993. Distribution-Free Calculation of the Standard Error of Chain Ladder Reserve Estimates. *ASTIN Bulletin* 23: 213–25.
- Martínez-Miranda, Maria D., Jens P. Nielsen, Richard J. Verrall, and Mario V. Wüthrich. 2015. The Link between Classical Reserving and Granular Reserving through Double Chain-Ladder and its Extensions. *Scandinavian Actuarial Journal* 5: 383–405.
- Pigeon, Mathieu, Katrien Antonio, and Michel Denuit. 2013. Individual Loss Reserving with the Multivariate Skew Normal Framework. *ASTIN Bulletin* 43: 399–428.
- Rumelhart, David E., Geoffrey E. Hinton, and Ronald J. Williams. 1986. Learning Representations by Back-Propagating Errors. *Nature* 323: 533–36.
- Taylor, Greg, Gráinne McGuire, and James Sullivan. 2008. Individual Claim Loss Reserving Conditioned by Case Estimates. *Annals of Actuarial Science* 3: 215–56.

Verrall, Richard J., and Mario V. Wüthrich. 2016. Understanding Reporting Delay in General Insurance. Risks 4: 25.

- Werbos, Paul J. 1982. Applications of Advances in Nonlinear Sensitivity Analysis. In System Modeling and Optimization. Paper Presented at the 10th IFIP Conference, New York City, NY, USA, 31 August–4 September 1981.
   Edited by Rudolf F. Drenick and Frank Kozin. Berlin and Heidelberg: Springer, pp. 762–70.
- Wüthrich, Mario V. 2018. Machine Learning in Individual Claims Reserving. *To appear in Scandinavian Actuarial Journal* 25: 1–16.
- Wüthrich, Mario V. 2018. Neural Networks Applied to Chain-Ladder Reserving. SSRN Manuscript. doi:10.2139/ssrn.2966126.



© 2018 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (http://creativecommons.org/licenses/by/4.0/).