**ORIGINAL RESEARCH PAPER**

# Combined modelling of micro-level outstanding claim counts and individual claim frequencies in non-life insurance

Axel Bücher[1,2] · Alexander Rosenstock[1,3]

© The Author(s) 2024

## Abstract
Usually, the actuarial problems of predicting the number of claims incurred but not reported (IBNR) and of modelling claim frequencies are treated successively by insurance companies. New micro-level methods designed for large datasets are proposed that address the two problems simultaneously. The methods are based on an elaborated occurrence process model that includes both a claim intensity model and a claim development model. The influence of claim feature variables is modelled by suitable neural networks. Extensive simulation experiments and a case study on a large real data set from a motor legal insurance portfolio show accurate predictions at both the aggregate and individual policy level, as well as appropriate fitted models for claim frequencies. Moreover, a novel alternative approach combining data from classic triangle-based methods with a micro-level intensity model is introduced and compared to the full micro-level approach.

✉ Axel Bücher
   axel.buecher@rub.de

   Alexander Rosenstock
   alexander.rosenstock@hhu.de

1   Mathematisches Institut, Heinrich-Heine-Universität Düsseldorf, Universitätsstraße 1, 40225 Düsseldorf, Germany

2   Fakultät für Mathematik, Ruhr-Universität Bochum, Universitätsstraße 150, 44780 Bochum, Germany

3   ARAG SE, ARAG-Platz 1, 40472 Düsseldorf, Germany

Springer

## 1 Introduction

Actuarial departments in insurance companies are responsible for many different tasks related to risk modelling of insurance operations. Two of these tasks are closely inter-linked: reserving actuaries need to compute accurate predictions of incurred liabilities for insurance portfolios. Pricing actuaries use these results to develop risk models, which allow them to estimate expected losses for a range of policy parameters and thus to determine the prices at which these policies are sold. While developing models on per-policy data (subsequently referred to as the micro-level) is common practice for risk modelling, micro-level methods for reserving are still not widely adopted in the industry. Adoption is slow in part due to the requirements necessary for micro-level reserving (such as fine-grained claims development information, the necessary data warehousing and compute infrastructure as well as advanced modelling skills such as machine learning in combination with actuarial expertise) and the strong governing bodies that will require tried and tested reserving methods for many business processes in the foreseeable future.

In the research literature on reserving in non-life insurance, there are two main approaches to micro-level reserving. On the one hand, there are methods that operate in discrete time (such as development years), which have similarities with macro-level reserving methods like Bornhuetter–Ferguson. Micro-level methods of this kind typically deal with claim-level information which is only available for reported claims, and therefore work on the reserve of claims that have been reported but not settled (RBNS). Examples of this kind are [6, 9, 11, 12, 24], among others. On the other hand, there are continuous time claim development models which are based on the assumption that claim development is regarded as a point process (e.g., a position-dependent marked Poisson process); respective model parameters are then estimated from observed data. This strain of research goes back to [17, 18]. More recent papers studying a continuous time claim development model are [4, 5, 19, 21], among others.

The current paper aims to combine the task of predicting IBNR claim counts on a policy level with the development of a matching micro-level claim intensity model, thus addressing the two tasks of reserving and risk modelling simultaneously in the hope of improving accuracy in both disciplines. Two main approaches are proposed to achieve this goal. Our first approach is based on an explicit micro-level claim occurrence process model inspired by [17]. Suitable methods are proposed for estimating all model parameters, which results in a fully fitted continuous-time process model that includes a claim intensity model. We thereby extend a related approach in [7], where a submodel for reporting delays was studied. Secondly, we propose a new method that uses classical triangle-level reserving methods as input to the estimation of a micro-level claim intensity model. This approach allows for micro-level allocation of IBNR claim counts that is consistent with the triangle. Due to the nature of the underlying triangle-based reserving methods, the approach only allows discrete time steps.

For each of the two approaches, we discuss the design and estimation of suitable (parametric) sub-models involving classical multilayer perceptron (MLP) neural networks, see [14]. For the estimation, it is taken into account that the observed data is subject to random truncation (indeed, a claim is only observed at a given calendar time if the sum of its associated (random) reporting delay and its (random) accident time

does not exceed the given calendar time). Predictors for the number of IBNR claims on a per-policy level are derived from the estimated models. They are compared with each other as well as with classical chain ladder methods in an extensive simulation study as well as on a real dataset. It is found that the new predictors provide accurate predictions as well as appropriate fitted models for claim intensities even in simulation scenarios involving non-homogeneous portfolios and in the real-life example. Moreover, in contrast to classical factor-based reserving methods, the predictor based on the first approach may yield a non-zero number of claims even for policies without already reported claims. This is natural as it allows for interpreting the prediction as the expected number of unreported claims for that particular policy.

For learning the neural networks, we utilize the TensorFlow framework [1] with custom loss functions that take random truncation into account. The implementation is written using the R language and its interface binding packages `keras` and `tensorflow` [2, 10] to utilize TensorFlow. Core functionality used for estimation and prediction is freely available as an R package `reservr` [20].

The remaining parts of the paper are organized as follows. In Sect. 2, we introduce the micro-level claim process used throughout as well as the observation setting. Modelling and fitting of the individual components of the claim process are discussed in Sect. 3. Based on a completely estimated micro-level model, we derive a corresponding micro-level predictor for the IBNR claim count in Sect. 4. Section 5 introduces an alternative micro-level predictor based on similar ideas, but using a triangle-based global reserving method and discrete time steps. After defining performance metrics in Sect. 6, we present results on a large-scale simulation study in Sect. 7. An application to a real dataset from a motor legal insurance portfolio is presented in Sect. 8. Finally, Sect. 9 concludes.

## 2 Preliminaries on insurance portfolio data

The general model for the claim arrivals in a given insurance portfolio is the same as in [7], and builds on the notion of *(position-dependent) marked Poisson processes* [17]. More precisely, we consider an insurance portfolio containing $\mathcal{I} \in \mathbb{N}$ independent risks. Each risk is described by a coverage period $C = [t_{\text{start}}, t_{\text{end}}]$, and by risk features $\bar{x} \in \bar{\mathfrak{X}}$, where $\bar{\mathfrak{X}}$ is a feature space containing both discrete and continuous features; for example, information on the insured product and chosen options such as deductibles. We write $x = (C, \bar{x}) \in \mathfrak{X} = \{\text{intervals on } [0, \infty)\} \times \bar{\mathfrak{X}}$, and assume that $x$ is constant over the course of the contract. In practice, risk features do change over time, but not very often, whence such a contract could be modelled as two separate risks.

Each risk can potentially incur claims during its coverage period, which will formally be modelled by a claim arrival process. Note that (marked) claim arrival processes provide a natural mathematical model for random event analysis, which have been proposed for actuarial risk modelling in [17]; see also the textbook [16]. Each claim in the claim arrival process occurs at some (calendar) accident time $t_{\text{acc}} \in [t_{\text{start}}, t_{\text{end}}]$, and will be associated with several claim features $y \in \mathfrak{Y}$ like the type of the claim or the value of the cars involved in some accident; here, $\mathfrak{Y}$ denotes some suitable feature space. Moreover, the claim will not be immediately known to
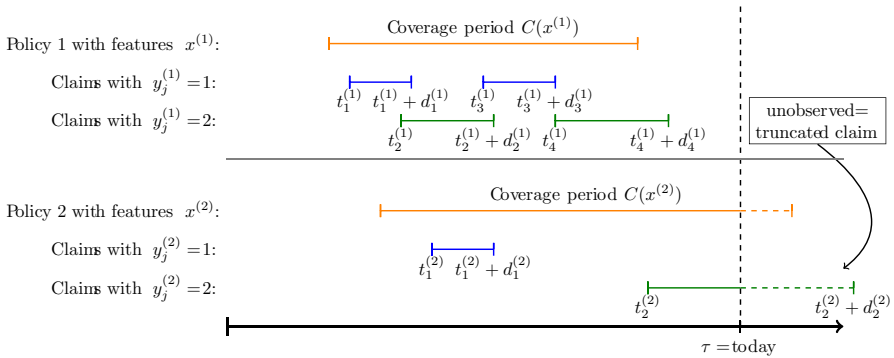
**Fig. 1** Illustration of a portfolio dataset consisting of two risks with 6 claims in total for the case where the claim covariate $y_j^{(i)}$ is a binary variable with values in $\{1, 2\}$. The second claim of the second policy, $t_2^{(2)}$, is unobserved by time $\tau$ and hence an IBNR claim (see Sect. 4 for details). It is the main goal of this paper to predict the number of IBNR claim counts for each policy, based on all observations available at time $\tau$

the insurer, but it will rather be reported at (calendar) reporting time $t_{\text{report}} \in [t_{\text{acc}}, \infty)$. Formally, both the claim features $y$ and the reporting delay $d_{\text{report}} := t_{\text{report}} - t_{\text{acc}}$ will be assumed to be a mark on the claim arrival process.

Let $\delta_z$ denote the dirac-measure at $z$. Following the notation in [15], we arrive at the following definition of a claim arrival process. The process is illustrated in Fig. 1 for the case that the claim feature is binary.

**Definition 2.1** (*Claim arrivals and portfolio*) Consider the $i$th risk in a portfolio, $i \in \{1, \ldots, \mathcal{I}\}$, with risk features $x^{(i)} \in \mathfrak{X}$ among which we find the coverage period $C(x^{(i)})$. The claim arrival process associated with that risk is a position-dependent marked Poisson process with $N^{(i)} \sim \text{Poi}\left(\int_{C(x^{(i)})} \lambda(x^{(i)}, t)\, dt\right)$ points

$$\xi^{(i)} = \sum_{j=1}^{N^{(i)}} \delta_{(T_{\text{acc},j}^{(i)}, Y_j^{(i)}, D_{\text{report},j}^{(i)})}$$

on $[0, \infty) \times \mathfrak{Y} \times [0, \infty)$ with:

(i) Intensity $\lambda(x^{(i)}, t)\mathbf{1}(t \in C(x^{(i)}))$, i.e., for all intervals $[t_0, t_1] \subseteq [0, \infty)$, we have

$$\sum_{j=1}^{N^{(i)}} \mathbf{1}(T_{\text{acc},j}^{(i)} \in [t_0, t_1]) = \int_{t_0}^{t_1} \xi^{(i)}(dt, \mathfrak{Y}, [0, \infty))$$

$$\sim \text{Poi}\left(\int_{t_0}^{t_1} \mathbf{1}(t \in C(x^{(i)}))\lambda(x^{(i)}, t)\, dt\right).$$

(ii) Conditional claim feature distribution $P_{Y|x^{(i)},t} = P_{Y|X=x^{(i)}, T_{\text{acc}}=t}$. Here, $Y$ denotes a generic $\mathfrak{Y}$-valued claim feature variable containing all claim features except for the reporting delay, while $X$ and $T_{\text{acc}}$ are generic risk feature and accident time variables, respectively.

(iii) Conditional reporting delay distribution $P_{D|x^{(i)},t,y} = P_{D|X=x^{(i)},T_{acc}=t,Y=y}$. Here, $D = D_{report}$ denotes a generic reporting delay variable in $[0, \infty)$.

A portfolio consists of $\mathcal{I}$ risks $\xi^{(1)}, \ldots, \xi^{(\mathcal{I})}$ that are mutually independent.

Note that the claim intensity $\lambda(x^{(i)}, t)$ controls the expected number of claims per exposure (the *claim frequency*), i.e., the expected amount of claims of the $i$th risk in the period $A \subset C(x^{(i)})$ is given by $\int_A \lambda(x^{(i)}, t) \, dt$.

In practice, the three building blocks of Definition 2.1, i.e., $\lambda(x, t)$, $P_{Y|x,t}$ and $P_{D|x,t,y}$, are unknown and must be estimated based on observational data that is available to the insurer at some given calendar time $\tau > 0$ of observing the portfolio. More precisely, the insurer observes data from Observation Scheme 2.2, which is again illustrated in Fig. 1.

**Observation Scheme 2.2** At given calendar time $\tau$, the available dataset $\mathfrak{D} = \mathfrak{D}_\tau$ consists of all risk features $x^{(i)}$, $i \in \{1, \ldots, \mathcal{I}\}$, and all reported claim data up to calendar time $\tau$, i.e.

$$\left\{ (x^{(i)}, t^{(i)}_{acc,j}, y^{(i)}_j, d^{(i)}_{report,j}) \mid t^{(i)}_{acc,j} + d^{(i)}_{report,j} \leq \tau \right\}. \tag{1}$$

Equivalently, we observe, for each $i \in \{1, \ldots, \mathcal{I}\}$, the risk feature $x^{(i)}$ and the restriction $\xi^{(i)}_r(\cdot) = \xi^{(i)}(\cdot \cap R_\tau)$, where $R_\tau = \{(t, y, d) : t + d \leq \tau\}$ and where the lower index $r$ stands for '*reported*'.

Clearly, estimating the building blocks of Definition 2.1 can only be feasible if additional model assumptions are made. Those assumptions, as well as respective fitting procedures, will be described in the next section.

## 3 Modelling and fitting claim arrival processes

Modelling and estimating the claim intensity $\lambda(x, t)$, the claim feature distribution $P_{Y|x,t}$ and the reporting delay distribution $P_{D|x,t,y}$ from Definition 2.1 will be done iteratively, starting with the latter. We discuss each aspect in a separate subsection.

### 3.1 Modelling and fitting the reporting delay distribution

Modelling and fitting the reporting delay distribution has recently been considered in the accompanying paper [7]. Throughout this section, we take a slightly more general approach as in that paper, as we do not restrict ourselves to a specific distribution family.

**Model 3.1** (*Micro-level model for reporting delays*) Let $\mathcal{P}^D = \{P^D_\theta : \theta \in \Theta^D\}$ denote a parametric distribution family that is dominated by a sigma-finite measure $\mu^D$ and that has a finite-dimensional parameter space $\Theta^D$; the $\mu^D$-densities of $P^D_\theta$ will be written as $f^D_\theta$. Further, let $\mathcal{G}^D$ denote a set of MLPs $g^D : \mathfrak{X} \times [0, \infty) \times \mathfrak{Y} \to \Theta^D$. We

assume that the conditional reporting delay distribution satisfies, for some $g^D \in \mathcal{G}^D$,

$$P_{D|x,t,y} = P^D_{g^D(x,t,y)} \qquad \forall x, t, y.$$

Details on fitting Model 3.1 to right-truncated observations as in Observation Scheme 2.2 can be found in Section 3 in [7]. A publicly available implementation is provided in [20]; it can be used for a variety of parametric families. Based on a discussion of stylized facts of reporting delays, [7] propose to work with a specific parametric mixture family for $\mathcal{P}^D$, the Blended Dirac-Erlang-Generalized Pareto Distribution family.

## 3.2 Modelling and fitting the claim feature distribution

Throughout this section, we assume that $P_{D|x,t,y}$ is available, for instance since it has been estimated as described in the previous section. For modelling and estimating the claim feature distribution $P_{Y|x,t}$, we assume that $\mathfrak{Y}$ can be written as a $Q$-fold cartesian product $\mathfrak{Y} = \mathfrak{Y}_1 \times \cdots \times \mathfrak{Y}_Q$ with $\mathfrak{Y}_q \subset \mathbb{R}$ for each $q \in \{1, \dots, Q\}$. Note that this assumption is of no practical concern, since claims data typically consists of a combination of real-valued and categorical features (if the $q$th feature is categorical, its categories may be identified with $1, \dots, n_q$). As a consequence, we may write a generic claim feature variable $Y$ as $Y = (Y_1, \dots, Y_Q)$, and by the chain rule for conditional distributions we can decompose the conditional claim feature distribution as

$$P_{Y|X,T} = P_{Y_Q|X,T,Y_1,\dots,Y_{Q-1}} \cdots P_{Y_2|X,T,Y_1} P_{Y_1|X,T}, \tag{2}$$

where $T = T_{\text{acc}}$. This decomposition will be crucial for defining a flexible model for which iterative estimation (from left to right) is feasible.

**Model 3.2** (*Micro-level model for claim features*) Assume that $P_{Y|X,T}$ allows for a decomposition as in (2) with $Q \in \mathbb{N}$. For each $q \in \{1, \dots, Q\}$, let $\mathcal{P}^{(q)} = \{P^{(q)}_\theta : \theta \in \Theta^{(q)}\}$ denote a parametric distribution family that is dominated by a sigma-finite measure $\mu^{(q)}$ and that has a finite-dimensional parameter space $\Theta^{(q)}$; the $\mu^{(q)}$-densities of $P^{(q)}_\theta$ will be written as $f^{(q)}_\theta$. Further, let $\mathcal{G}^{(q)}$ denote a set of MLPs $g^{(q)} : \mathfrak{X} \times [0, \infty) \times \mathfrak{Y}_1 \cdots \times \mathfrak{Y}_{q-1} \to \Theta^{(q)}$. We assume that there exists $g^{(q)} \in \mathcal{G}^{(q)}$ such that

$$P_{Y_q|x,t,y_1,\dots,y_{q-1}} = P^{(q)}_{g^{(q)}(x,t,y_1,\dots,y_{q-1})} \qquad \forall x, t, y_1, \dots, y_{q-1}.$$

A natural choice for $\mathcal{P}^{(q)}$ in case of a categorical component is the multinomial distribution, which gives full flexibility. Distributions for continuous components must be decided case-by-case, bearing in mind that integration with respect to the chosen distribution will need to be performed. Therefore, choosing an overly flexible distribution could lead to numerical problems in application.

We will now describe how to iteratively estimate the unknown components of Model 3.2, taking into account the fact that observations are subject to random truncation. Suppose we have already fitted $P_{D|x,t,y}$, $P_{Y_Q|x,t,y_1,\ldots,y_{Q-1}}$, $\ldots$, $P_{Y_{q+1}|x,t,y_1,\ldots,y_q}$, and we are to estimate $P_{Y_q|x,t,y_1,\ldots,y_{q-1}}$ next.

For that purpose, we propose to maximize the following weighted conditional likelihood function over all functions $g \in \mathcal{G}^{(q)}$:

$$\tilde{\mathcal{L}}(g|\mathfrak{D}_\tau) = \sum_{(x,t,y,d)\in\mathfrak{D}_\tau} \tilde{\ell}_{(x,t,y_1,\ldots,y_{q-1})}(g|y_q), \tag{3}$$

where

$$
\begin{aligned}
&\tilde{\ell}_{(x,t,y_1,\ldots,y_{q-1})}(g|y_q)\\
&= \frac{\log f^{(q)}_{g(x,t,y_1,\ldots,y_{q-1})}(y_q)}{P(T + D \le \tau | X = x, T = t, Y_1 = y_1, \ldots, Y_q = y_q)}.
\end{aligned}
\tag{4}
$$

Thereby, the log-likelihood contribution of each observation is essentially weighted with the reciprocal of the probability of observing that particular observation, i.e., observations that are more likely to be truncated (i.e., less likely to be observed) get a higher weight in the log-likelihood. Further, note that the denominator in the definition of $\tilde{\ell}$ does not depend on $g$, but only on objects that have already been fitted. Indeed, writing $y^{(q)} = (y_1, \ldots, y_q)$, we have

$$
\begin{aligned}
&P(T \le \tau | X = x, T = t, Y^{(q)} = y^{(q)})\\
&= \int_{\mathfrak{Y}_{q+1}} \cdots \int_{\mathfrak{Y}_Q} P_{D|x,t,y}([0, \tau - t])\, \mathrm{d}P_{Y_Q|x,t,y^{(Q-1)}}(y_Q) \ldots \mathrm{d}P_{Y_{q+1}|x,t,y^{(q)}}(y_{q+1}),
\end{aligned}
\tag{5}
$$

which can readily be computed for each observation $(x, t, y, d) \in \mathfrak{D}_\tau$, subject to computational constraints.

The estimator for $g$ may be motivated as follows: first of all, standard heuristics underlying the M-estimation principle suggest that a maximizer of $\tilde{\mathcal{L}}$ may be considered as an estimator for the maximizer of the (conditional) expected value $g \mapsto L(g) := \mathbb{E}_{Y_q|T+D\le\tau,x,t,y^{(q-1)}}[\tilde{\mathcal{L}}(g|\mathfrak{D}_\tau)]$, where $\mathbb{E}_{Y_q|T+D\le\tau,x,t,y^{(q-1)}}$ refers to integration with respect to the true conditional distribution of $Y_q$ given $T + D \le \tau$, $X = x$, $T = t$, $Y^{(q-1)} = y^{(q-1)}$, for each observation. More precisely, we have

$$
\begin{aligned}
L(g) = &\sum_{(x,t,y,d)\in\mathfrak{D}_\tau} \mathbb{E}[\tilde{\ell}_{(x,t,y_1,\ldots,y_{q-1})}(g|Y_q) \mid T\\
&+ D \le \tau, X = x, T = t, Y^{(q-1)} = y^{(q-1)}].
\end{aligned}
\tag{6}
$$

The following lemma characterizes the maximizers of $g \mapsto L(g)$; its proof is given in Appendix A.

**Lemma 3.3** *Assume that there is a true function $g_0 = g_0^{(q)} \in \mathcal{G}^{(q)}$ such that*

$$P_{Y_q|x,t,y_1,\ldots,y_{q-1}} = P_{g_0(x,t,y_1,\ldots,y_{q-1})}^{(q)} \qquad \forall x, t, y_1, \ldots, y_{q-1}.$$

*Then, for each fixed value of $x, t, y^{(q-1)}$, the summands of the objective function from* (6)

$$g \mapsto \mathbb{E}[\tilde{\ell}_{(x,t,y_1,\ldots,y_{q-1})}(g|Y_q) \mid T + D \leq \tau, X = x, T = t, Y^{(q-1)} = y^{(q-1)}]$$

*attain their maximal value at $g = g_0$.*

During preliminary simulation experiments we found that more reliable estimates with a smaller variance may be obtained by smoothing the denominator in (4). This requires additional assumptions on top of Definition 2.1, the *local homogeneity* assumptions.

**Assumption 3.4** (Local homogeneity of claims developement) Let $p > 0$ be a given period length measured in days; e.g., $p = 365$ days. For all intervals $I_p(k) = [kp, kp + p)$ with midpoints $t_k = kp + \frac{p}{2}, k \in \mathbb{N}_0$, we have:

(i) $t \mapsto \lambda(x, t) = \lambda(x, t_k) > 0$ is constant on $I_p(k)$ for any $x$.
(ii) $t \mapsto P_{Y|x,t} = P_{Y|x,t_k}$ is constant on $I_p(k)$ for any $x$.
(iii) $t \mapsto P_{D|x,t,y} = P_{D|x,y,t_k}$ is constant on $I_p(k)$ for any $x, y$.

Heuristically, the smaller $p$ is, the closer the "unknown truth" is to a model that fulfills the homogeneity assumptions. For our final predictors, $p$ may be regarded as a hyperparameter to be chosen by the statistician to balance model bias and variance: a smaller choice for $p$ increases estimation variance while allowing for a more flexible model and hence less bias. We will often work with $p = 365$ for simplicity, which was found to provide reasonable predictions in applications.

Implicitly assuming Assumption 3.4 for some given period length $p > 0$, we propose to replace the denominator $P(T + D \leq \tau | X = x, T = t, Y_1 = y_1, \ldots, Y_q = y_q)$ in (4), see also the alternative expression in (5), by

$$\text{denom}_{q-1}(x, t_{k_t}, y_1, \ldots, y_{q-1})$$
$$:= \frac{1}{\text{Leb}(I_p(k_t) \cap C)} \int_{I_p(k_t) \cap C} \int_{\mathfrak{Y}_{q+1}} \cdots \int_{\mathfrak{Y}_Q} P_{D|x,t_{k_t},y}([0, \tau - s])$$
$$\times dP_{Y_Q|x,t_{k_t},y^{(Q-1)}}(y_Q) \ldots dP_{Y_{q+1}|x,t_{k_t},y^{(q)}}(y_{q+1}) \, ds, \tag{7}$$

where $k_t = \lfloor \frac{t}{p} \rfloor$ denotes the number of the period of length $p$ containing $t$, which in turn, using the notation from Assumption 3.4, is $I_p(k_t) = [k_t p, k_t p + p)$ with midpoint $t_{k_t} = k_t p + \frac{p}{2}$. Moreover, $C = C(x)$ is the coverage period associated with $x$ and Leb refers to the Lebesgue measure. Note that both the denominator and the integral are non-zero for observed values $(x, t, y, d) \in \mathfrak{D}_\tau$. Overall, we aim at maximizing

$$\mathcal{L}^{(p)}(g|\mathfrak{D}_\tau) = \sum_{(x,t,y,d) \in \mathfrak{D}_\tau} \ell_{(x,t,y_1,\ldots,y_{q-1})}^{(p)}(g|y_q)$$

instead of (3), where, recalling $\text{denom}_{q-1}(x, t_{k_t}, y_1, \ldots, y_{q-1})$ from (7),

$$\ell^{(p)}_{(x,t,y_1,\ldots,y_{q-1})}(g|y_q) = \frac{\log f^{(q)}_{g(x,t,y_1,\ldots,y_{q-1})}(y_q)}{\text{denom}_{q-1}(x, t_{k_t}, y_1, \ldots, y_{q-1})}.$$

### 3.3 Modelling and fitting the claim intensity

Once a distribution for $P_{Y|X,T}$ has been fitted, the only unknown object in the model from Definition 2.1 is the claim intensity $\lambda = \lambda(x, t)$.

By the restriction theorem (Theorem 5.2 in [15]), the reported claims process $\xi_r^{(i)} = \xi^{(i)}(\cdot \cap R_\tau)$ with $R_\tau = \{(t, y, d) : t + d \leq \tau\}$ has intensity measure

$$\mu_r^{(i)}(S) = \mu^{(i)}(S \cap R_\tau) = \mathbb{E}[\xi^{(i)}(S \cap R_\tau)]$$
$$= \int_{C(x^{(i)})} \int_{\mathfrak{Y}} \int_{[0,\tau-t]} \mathbf{1}_S(t, y, d)\lambda(x^{(i)}, t)\, dP_{D|x^{(i)},t,y}(d)\, dP_{Y|x^{(i)},t}(y)\, dt,$$

where $S \subset [0, \infty) \times \mathfrak{Y} \times [0, \infty)$.

For some period length $p > 0$, let $I_p(k) = [kp, (k+1)p)$ denote the $k$th period of length $p$ with midpoint $t_k = kp + \frac{p}{2}$ and let

$$I_p(x, k) := C(x) \cap I_p(k) = C(x) \cap [kp, (k+1)p) \tag{8}$$

denote the coverage time of policy $x$ within the $k$th period. Assuming local homogeneity as in Assumption 3.4 for period length $p > 0$, and letting $S(k) = I_p(k) \times \mathfrak{Y} \times [0, \infty)$ denote the set of all claims $(y, t, d)$ that occur in the $k$th period, we obtain that, for each $k \in \mathbb{N}_0$,

$$\xi_r^{(i)}(S(k)) \sim \text{Poi}\left(\lambda(x^{(i)}, t_k) \int_{\mathfrak{Y}} \int_{I_p(x^{(i)},k)} P_{D|x^{(i)},t_k,y}([0, \tau - s])\, ds\, dP_{Y|x^{(i)},t_k}(y)\right).$$

Given a set of policies, this allows fitting a Poisson model (e.g., a Poisson GLM with log link) to the number of reported claims per period in the usual way by specifying a fixed offset $o$ for each observation and estimating the common intensity factor $\lambda(x, t)$. Compared to a classical claim intensity model without truncation, where $\xi^{(i)}(S(k)) \sim \text{Poi}\left(\lambda(x^{(i)}, t_k)\text{Leb}(I_p(x^{(i)}, k))\right)$ with $\text{Leb}(I_p(x^{(i)}, k))$ usually called the *exposure*, the offset term $\log(\text{Leb}(I_p(x^{(i)}, k)))$ must be adjusted by the reporting probability; see (10) below. See [13, 23] for a more detailed introduction to the classical intensity modelling approach and offsets.

**Model 3.5** (*Micro-level model for claim intensity*) Let $\mathcal{G}$ denote a set of MLPs $g : \mathfrak{X} \times [0, \infty) \to \mathbb{R}_+$. We assume that there exists $g \in \mathcal{G}$ such that $\lambda(x, t) = g(x, t_{k_t})$ for all $x \in \mathfrak{X}$ and all $t > 0$, i.e., the claim intensity $\lambda$ is given by a piecewise constant extension of $g(x, t_k)$ to the intervals $I_p(k)$ for $k = 0, 1, \ldots$, which is consistent with using Assumption 3.4 for period length $p$.

The claim intensity $\lambda(x, t)$ can hence be estimated by maximizing the Poisson loglikelihood

$$\mathcal{L}(g|\mathfrak{D}_\tau) = \sum_{i=1}^{\mathcal{I}} \sum_{k=0}^{\infty} \xi_r^{(i)}(S(k)) \log\left(g(x^{(i)}, t_k)\mathrm{ex}(x^{(i)}, k)\right) - g(x^i, t_k)\mathrm{ex}(x^{(i)}, k),$$

$$(9)$$

where

$$\mathrm{ex}(x^{(i)}, k) := \int_{\mathfrak{Y}} \int_{I_p(x^{(i)}, k)} P_{D|x^{(i)}, t_k, y}([0, \tau - s)) \, ds \, dP_{Y|x^{(i)}, t_k}(y). \qquad (10)$$

Note that $\mathbb{E}[\xi_r^{(i)}(S(k))] = \lambda(x^{(i)}, t_k)\mathrm{ex}(x^{(i)}, k)$, whence $\mathrm{ex}(x^{(i)}, k)$ may be interpreted as the expected number of claims that have occurred in the accident period $I_p(k)$ and are reported by calendar time $\tau$ for a policy with constant claim intensity $\lambda(x^{(i)}, t_k) = 100\%$. If $\hat{g} \in \mathcal{G}$ maximizes $\mathcal{L}(g|\mathfrak{D}_\tau)$, we write

$$\hat{\lambda}^{\mathrm{NNet}}(x, t) = \hat{g}(x, t_{k_t}).$$

Note that maximization of $\mathcal{L}(g|\mathfrak{D}_\tau)$ is straight-forward once the exposures $\mathrm{ex}(x^{(i)}, k)$ have been computed. The latter requires numerical integration over $\mathfrak{Y}$, after replacing $P_{D|x,t,y}$ and $P_{Y|x,t}$ by estimated versions thereof. Care must be taken in the choice of $\mathfrak{Y}$ during modelling, so this integral remains feasible: choosing continuous covariates necessitates computation of possibly challenging (and maybe indefinite) integrals with respect to $P_{Y_q|x,t,y_1,\ldots,y_{q-1}}$, choosing too many discrete covariates results in combinatorial explosion of the number of summands to be computed when performing integration with respect to the counting measure.

## 4 Individual claims count prediction based on estimated claim arrival processes

The models and estimators from the previous section can be used in various ways to define predictors for IBNR claim numbers; see [7] for an example that only involves the reporting delay model. Throughout this section, we describe a predictor that is based on the full (estimated) claim arrival model. Alternative intermediate predictors will be defined in the simulation study.

More precisely, for each given period $I_p(k) = [kp, (k+1)p)$ of length $p > 0$ and each claim feature set $\mathfrak{Y}' \subset \mathfrak{Y}$ and each reporting interval $(\tau_0, \tau_1] \subset [0, \infty]$, we derive a predictor for the number of claims policy $i$ has incurred within period $I_p(k)$ with claim features in $\mathfrak{Y}'$ and with a reporting time in $(\tau_0, \tau_1]$. For that purpose, let $S'(k) := I_p(k) \times \mathfrak{Y}' \times [0, \infty) = \{(t, y, d) : t \in I_p(k)\}$ denote the set of claims that occurred in the $k$th period with claim features in $\mathfrak{Y}'$ and let

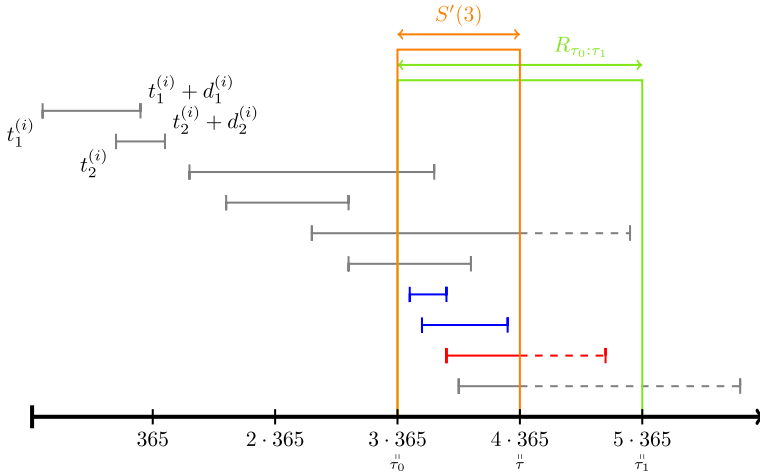$$R_{\tau_0:\tau_1} := \{(t, y, d) : \tau_0 < t + d \le \tau_1\} \qquad (11)$$

**Fig. 2** Illustration of claim counts for the $i$th policy occurring in the 3rd year (i.e., claims in $S'(3)$) and reported in year 3 or 4 (i.e., claims in $R_{\tau_0:\tau_1}$ with $\tau_0 = 3 \cdot 365$, $\tau_1 = 5 \cdot 365$). With today $\tau = 4 \cdot 365$, we have two reported claims (blue color) and 1 unreported IBNR claim (red color)

denote the set of claims reported between times $\tau_0$ and $\tau_1$, see Fig. 2 for an illustration. For completeness, let $R_{\tau_0:\tau_1} = \varnothing$ if $\tau_0 \geq \tau_1$.

Note that the target number of claims for the $i$th policy can then be written as

$$N_{\tau_0:\tau_1}^{(i)}(S'(k)) := \xi^{(i)}(S'(k) \cap R_{\tau_0:\tau_1}),$$

and that we observe, under Observation Scheme 2.2, the respective number of reported claims $\xi_r^{(i)}(S'(k) \cap R_{\tau_0:\tau_1}) = \xi^{(i)}\big(S'(k) \cap R_{\tau_0:\min(\tau_1, \tau)}\big)$, which is zero if $\tau_0 > \tau$.

Now, if Assumption 3.4 is met for the given period length $p > 0$, we obtain that, by the restriction theorem (Theorem 5.2 in [15]),

$$\mathbb{E}[\xi^{(i)}(S'(k) \cap R_{\tau_0:\tau_1}) \mid \xi_r^{(i)}(S'(k) \cap R_{\tau_0:\tau_1})]$$
$$= \xi_r^{(i)}(S'(k) \cap R_{\tau_0:\tau_1}) + \mathbb{E}[\xi_{nr}^{(i)}(S'(k) \cap R_{\tau_0:\tau_1})]$$
$$= \xi_r^{(i)}(S'(k) \cap R_{\tau_0:\tau_1}) + \mathbb{E}[\xi^{(i)}(S'(k) \cap R_{\max(\tau_0,\tau):\tau_1})]$$
$$= \xi_r^{(i)}(S'(k) \cap R_{\tau_0:\tau_1}) + \lambda(x^{(i)}, t_k)$$
$$\times \int_{\mathfrak{Y}'} \int_{I_p(x^{(i)}, k)} P_{D|x^{(i)}, t_k, y}(I_{\tau_0:\tau_1}(\tau, s)) \, ds \, P_{Y|x^{(i)}, t_k}(dy).$$

Here, $\xi_{nr}^{(i)} = \xi^{(i)} - \xi_r^{(i)}$ denotes the unknown number of unreported claims, $I_p(x, k)$ is the coverage time of the policy associated with $x$ within the $k$th period, see (8), and $I_{\tau_0:\tau_1}(\tau, s) := (\max(\tau, \tau_0) - s, \tau_1 - s]$, with the convention that the interval is the empty set if $\max(\tau, \tau_0) > \tau_1$. As is well-known, if $\lambda(x, t)$, $P_{Y|x,t}$ and $P_{D|x,t,y}$ were known, this would be the best $L^2$-predictor for $\xi^{(i)}(S'(k) \cap R_{\tau_0:\tau_1})$ (the total number of claims in $S'(k)$ that are reported between $\tau_0$ and $\tau_1$) in terms of the observed counterpart of

reported claims $\xi_r^{(i)}(S'(k) \cap R_{\tau_0:\tau_1})$. Replacing the unknown objects on the right-hand side by suitable estimators as in the previous sections, we arrive at the predictor

$$
\begin{aligned}
&\hat{N}_{\tau_0:\tau_1}^{(i)}(S'(k)) \\
&\quad := \hat{\xi}^{(i)}(S'(k) \cap R_{\tau_0:\tau_1}) := \xi_r^{(i)}(S'(k) \cap R_{\tau_0:\tau_1}) \\
&\qquad + \hat{\lambda}(x^{(i)}, t_k) \int_{\mathfrak{Y}'} \int_{I_p(x^{(i)},k)} \hat{P}_{D|x^{(i)},t_k,y}(I_{\tau_0:\tau_1}(\tau, s)) \, \mathrm{d}s \, \mathrm{d}\hat{P}_{Y|x^{(i)},t_k}(y).
\end{aligned} \tag{12}
$$

In contrast to classical factor-based reserving methods, this predictor may yield a non-zero expected number of claims even for policies without already reported claims. This allows for the individual-level count predictions to have a meaningful interpretation as the expected number of unreported claims for that particular policy.

**Remark 4.1** The predictor in (12) can be adapted to a general set of claims $S = I \times \mathfrak{Y}' \times [0, \infty)$ by summing over the intervals covering $I$ as follows:

$$
\begin{aligned}
\hat{N}_{\tau_0:\tau_1}^{(i)}(S) = \xi_r(S \cap R_{\tau_0:\tau_1}) + \sum_{k=0}^{\tau/p-1} \hat{\lambda}(x^{(i)}, t_k) \\
\times \int_{\mathfrak{Y}'} \int_{I_p(x^{(i)},k) \cap I} \hat{P}_{D|x^{(i)},t_k,y}(I_{\tau_0:\tau_1}(\tau, s)) \, \mathrm{d}s \, \mathrm{d}\hat{P}_{Y|x^{(i)},t_k}(y).
\end{aligned}
$$

## 5 Individual claim count prediction based on chain ladder networks for the claim intensity

We propose an alternative estimator for the claim intensity $\lambda$, which is similar to the estimator from Sect. 3.3. However, instead of being based on preliminary estimators of the claim feature and reporting delay distributions, the new estimator is based on classical chain ladder factors. In a second step, the estimator is used to define a new predictor for IBNR claims, similarly as in Sect. 4.

Given a partition $\mathfrak{Y} = \mathfrak{Y}_1 \cup \mathfrak{Y}_2 \cup \cdots \cup \mathfrak{Y}_M$ into groups of claim features and a development period length $p$ where $\tau = \bar{\tau} p$ for some $\bar{\tau} \in \mathbb{N}_{\geq 2}$, we make the following classical chain ladder assumption: for any group index $m \in \{1, \ldots, M\}$ and any development period $j \in \{1, \ldots, \bar{\tau} - 1\}$, there exists a factor $f_j^{\mathrm{CL},\mathfrak{Y}_m}$ called *chain ladder factor* such that, for any policy $i$ and any accident period $k \in \{0, \ldots, \bar{\tau} - 1\}$,

$$
\mathbb{E}[\xi^{(i)}(S_m(k) \cap R_{0:(k+j+1)p})] = f_j^{\mathrm{CL},\mathfrak{Y}_m} \mathbb{E}[\xi^{(i)}(S_m(k) \cap R_{0:(k+j)p})],
$$

where $S_m(k) := I_p(k) \times \mathfrak{Y}_m \times [0, \infty)$ denotes the set of claims with claim features from $\mathfrak{Y}_m$ occurring in the $k$th period and $R_{0:kp}$ from (11) denotes the set of claims reported until time $kp$. Iterating the equation for fixed $k$ and with $j = \bar{\tau} - 1, \ldots, \bar{\tau} - k$, we obtain that

$$
\mathbb{E}[\xi^{(i)}(S_m(k) \cap R_{0:(k+\bar{\tau})p})] = \mathrm{FtU}_k^{\mathrm{CL},\mathfrak{Y}_m} \mathbb{E}[\xi^{(i)}(S_m(k) \cap R_{0:\bar{\tau}p})], \tag{13}
$$

where

$$\mathrm{FtU}_k^{\mathrm{CL},\mathfrak{Y}_m} := \prod_{j=\bar{\tau}-k}^{\bar{\tau}-1} f_j^{\mathrm{CL},\mathfrak{Y}_m}$$

is the *chain ladder factor-to-ultimate*. Equation (13) has the following interpretation: the expected number of claims from $\mathfrak{Y}_m$ with accident period $k$ that are reported within the next $k$ periods from today is equal to the $\mathfrak{Y}_m$-specific factor-to-ultimate multiplied with the expected number of claims from $\mathfrak{Y}_m$ with accident period $k$ that have been reported until today (which is observable). Under the additional assumption that every claim is developed within at most $\bar{\tau}$ periods, we have that $\xi^{(i)}(S_m(k) \cap R_{0:(k+\bar{\tau})p}) = \xi^{(i)}(S_m(k))$. Hence, if Assumption 3.4 is met for $p > 0$ specified above, the left-hand side of (13) can be written as

$$\mathbb{E}[\xi^{(i)}(S_m(k) \cap R_{0:(k+\bar{\tau})p})] = \mathbb{E}[\xi^{(i)}(S_m(k))]$$
$$= P_{Y|x^{(i)},t_k}(\mathfrak{Y}_m)\mathrm{Leb}(I_p(x^{(i)},k))\lambda(x^{(i)},t_k).$$

On the other hand, for the expression on the right-hand side of (13), we observe that $\xi^{(i)}(S_m(k) \cap R_{0:\bar{\tau}p}) = \xi_r^{(i)}(S_m(k))$ is the reported number of claims from $\mathfrak{Y}_m$ with accident period $k$. Hence, combining the previous equations with (13), we obtain that

$$\mathbb{E}[\xi_r^{(i)}(S_m(k))] = \frac{1}{\mathrm{FtU}_k^{\mathrm{CL},\mathfrak{Y}_m}} P_{Y|x^{(i)},t_k}(\mathfrak{Y}_m)\mathrm{Leb}(I_p(x^{(i)},k))\lambda(x^{(i)},t_k).$$

In view of the basic Poisson assumption on $\xi^{(i)}$ from Definition 2.1 (and hence on the reported and unreported counterparts $\xi_r^{(i)}$ and $\xi_{nr}^{(i)} = \xi^{(i)} - \xi_r^{(i)}$ from Observation Scheme 2.2), we obtain that

$$\xi_r^{(i)}(S_m(k)) \sim \mathrm{Poi}\left(P_{Y|x^{(i)},t_k}(\mathfrak{Y}_m)\mathrm{Leb}(I_p(x^{(i)},k))\lambda(x^{(i)},t_k)\frac{1}{\mathrm{FtU}_k^{\mathrm{CL},\mathfrak{Y}_m}}\right),$$

$$\xi_{nr}^{(i)}(S_m(k)) \sim \mathrm{Poi}\left(P_{Y|x^{(i)},t_k}(\mathfrak{Y}_m)\mathrm{Leb}(I_p(x^{(i)},k))\lambda(x^{(i)},t_k)\left(1 - \frac{1}{\mathrm{FtU}_k^{\mathrm{CL},\mathfrak{Y}_m}}\right)\right).$$

This can be used to estimate the unknown claim intensities on $\mathfrak{Y}_m$, i.e.,

$$\lambda^{\mathfrak{Y}_m}(x,t) := P_{Y|x,t}(\mathfrak{Y}_m)\lambda(x,t).$$

Indeed, let $\mathcal{G}$ denote a set of MLPs $g : \mathfrak{X} \times [0,\infty) \to [0,\infty)$ as in Model 3.5. We assume that the claim intensity on $\mathfrak{Y}_m$ satisfies, for some $g^{\mathfrak{Y}_m} \in \mathcal{G}$,

$$\lambda^{\mathfrak{Y}_m}(x,t) = g^{\mathfrak{Y}_m}(x,t_{k_t}) \qquad \forall x,t,$$

where $k_t = \lfloor \frac{t}{p} \rfloor$ indicates that the $k_t$-th period of length $p$ contains $t$ and where $t_k = kp + \frac{p}{2}$ denotes the mid-point of the $k$th period.

As in (9), we arrive at the per-triangle loss

$$
\mathcal{L}^{\mathrm{CL},\mathfrak{Y}_m}(g^{\mathfrak{Y}_m}|\mathfrak{D}_\tau) \;=\; \sum_{i=1}^{\mathcal{I}} \sum_{k=0}^{\bar{\tau}-1} \xi_r^{(i)}(S_m(k)) \log\left(g^{\mathfrak{Y}_m}(x^{(i)},t_k)\mathrm{ex}^{\mathrm{CL},\mathfrak{Y}_m}(x^{(i)},k)\right)
$$
$$
- g^{\mathfrak{Y}_m}(x^{(i)},t_k)\mathrm{ex}^{\mathrm{CL},\mathfrak{Y}_m}(x^{(i)},k),
$$

where

$$
\mathrm{ex}^{\mathrm{CL},\mathfrak{Y}_m}(x^{(i)},k) := \mathrm{Leb}(I_p(x^{(i)},k))\frac{1}{\mathrm{FtU}_k^{\mathrm{CL},\mathfrak{Y}_m}}.
$$

As in (10), $\mathrm{ex}^{\mathrm{CL},\mathfrak{Y}_m}(x^{(i)},k)$ may now be interpreted as the expected number of claims that have occurred in the accident period $I_p(k)$ and are reported by calendar time $\tau$ for a policy with constant claim intensity $\lambda^{\mathfrak{Y}_m}(x^{(i)},t_k) = 100\%$. In practice, the chain ladder factors within the loss must be estimated, for which we apply the well-known estimators

$$
\hat{f}_j^{\mathrm{CL},\mathfrak{Y}_m}
$$
$$
:= \frac{\#\{(x,t,y,d)\in\mathfrak{D}_\tau \mid y\in\mathfrak{Y}_m,\ \lfloor t/p\rfloor \le (\bar{\tau}-j-1),\ \lfloor (t+d)/p\rfloor - \lfloor t/p\rfloor \le j\}}{\#\{(x,t,y,d)\in\mathfrak{D}_\tau \mid y\in\mathfrak{Y}_m,\ \lfloor t/p\rfloor \le (\bar{\tau}-j-1),\ \lfloor (t+d)/p\rfloor - \lfloor t/p\rfloor \le j-1\}},
$$
$$
\widehat{\mathrm{FtU}}_k^{\mathrm{CL},\mathfrak{Y}_m} := \prod_{j=\bar{\tau}-k}^{\bar{\tau}-1} \hat{f}_j^{\mathrm{CL},\mathfrak{Y}_m}.
$$

Note that in contrast to the micro-level approach from the previous sections, there is no explicit model for the distribution of claim features on $\mathfrak{Y}$. If $g^{\mathfrak{Y}_m} = \hat{\lambda}^{\mathrm{CL},\mathfrak{Y}_m}(x,t)$ are maxima of the per-triangle losses $\mathcal{L}^{\mathrm{CL},\mathfrak{Y}_m}$, the triangle-level claim intensity estimates can be aggregated to a common intensity estimator

$$
\hat{\lambda}^{\mathrm{CL}}(x,t) := \sum_{m=1}^{M} \hat{\lambda}^{\mathrm{CL},\mathfrak{Y}_m}(x,t).
$$

Finally, exploiting $\mathbb{E}[\xi^{(i)}(S)\mid\xi_r^{(i)}(S)] = \xi_r^{(i)}(S) + \mathbb{E}[\xi_{nr}^{(i)}(S)]$ similar as in Sect. 4, we may define an IBNR-predictor as follows: recalling $S_m(k) = I_p(k)\times\mathfrak{Y}_m\times[0,\infty)$, the claims from $\mathfrak{Y}_m$ occurring in period $k$, let

$$
\hat{\xi}^{(i),\mathrm{CL}}(S_m(k)) := \xi_r^{(i)}(S_m(k)) + \hat{\lambda}^{\mathrm{CL},\mathfrak{Y}_m}(x^{(i)},t_k)\mathrm{Leb}(I_p(x^{(i)},k))\left(1 - \frac{1}{\widehat{\mathrm{FtU}}_k^{\mathrm{CL},\mathfrak{Y}_m}}\right),
$$
$$
\hat{\xi}^{(i),\mathrm{CL}}(S(k)) := \sum_{m=1}^{M} \hat{\xi}^{(i),\mathrm{CL}}(S_m(k)).
$$

Recalling the notation $R_{\tau_0:\tau_1}:=\{(t, y, d) : \tau_0 < t + d \leq \tau_1\}$, similar derivations show that this predictor can also be extended to a predictor for reporting times $\tau_0 = \bar{\tau}_0 p$ to $\tau_1 = \bar{\tau}_1 p$ with $\bar{\tau}_0 < \bar{\tau}_1 \in \mathbb{N}_0 \cup \{+\infty\}$:

$$\hat{N}_{\tau_0:\tau_1}^{(i),\text{CL}}(S_m(k)):=\xi_r^{(i)}(S_m(k) \cap R_{\tau_0:\tau_1})$$

$$+\hat{\lambda}^{\text{CL},\mathfrak{Y}_m}(x^{(i)}, t_k)\text{Leb}(I_p(x^{(i)}, k)) \left( \frac{1}{\widehat{\text{FtU}}_{k+\bar{\tau}-\bar{\tau}_1}^{\text{CL},\mathfrak{Y}_m}} - \frac{1}{\widehat{\text{FtU}}_{k+\bar{\tau}-\bar{\tau}_0}^{\text{CL},\mathfrak{Y}_m}} \right) \quad (14)$$

Here, we define the empty product as 1, i.e., $\widehat{\text{FtU}}_j^{\text{CL},\mathfrak{Y}_m}:=1$ for all $j \leq 0$. Finally, it is worthwhile to mention that in contrast to the predictor described in Sect. 4, it is not possible to use the chain ladder networks to obtain predictions for arbitrary $\tau_0, \tau_1$ that are not whole multiples of $p$.

## 6 Evaluating individual claim count predictors

The quality of competing predictors may be assessed by suitable error measures. In this section, we define two such measures: an individual mean squared prediction error, and an aggregated global mean squared prediction error.

We start by considering the individual error measure. For $S = [0, \tau) \times \mathfrak{Y}' \times [0, \infty) \subset [0, \infty) \times \mathfrak{Y} \times [0, \infty)$ and $0 \leq \tau_0 < \tau_1 \leq \infty$, let $\hat{N}_{\tau_0:\tau_1}^{(i)}(S)$ denote individual claim count predictions for $N_{\tau_0:\tau_1}^{(i)}(S) = \xi^{(i)}(S \cap R_{\tau_0:\tau_1})$, the number of claims in $S$ incurred by policy $x^{(i)}$ that are reported between $\tau_0$ and $\tau_1$; recall $R_{\tau_0:\tau_1} = \{(t, y, d) : \tau_0 < t + d \leq \tau_1\}$. Let $q > 0$ denote an evaluation period length (for instance, $q = 365$ corresponding to a year; note that there should be no confusion with the running index $q$ used in Sect. 3.2), which is assumed to be a divisor of the total observation length $\tau$ from Observation Scheme 2.2, and let $\mathfrak{Y}' \subset \mathfrak{Y}$ denote an evaluation set of claim features. We then define

$$\text{RMSE}_{\tau_0:\tau_1}^{\text{expo}}(\mathfrak{Y}', q)$$

$$:=\left( \frac{1}{\sum_{j=0}^{\tau/q-1} \#\mathcal{P}(q, j)} \sum_{\ell=0}^{\tau/q-1} \sum_{i \in \mathcal{P}(q,\ell)} \left\{ \hat{N}_{\tau_0:\tau_1}^{(i)}(S_q'(\ell)) - N_{\tau_0:\tau_1}^{(i)}(S_q'(\ell)) \right\}^2 \right)^{\frac{1}{2}}, \quad (15)$$

where, for $\ell \in \mathbb{N}_0$, recalling the notation $I_q(x, \ell) = C(x) \cap [\ell q, (\ell + 1)q)$ for the covering time of policy $x$ within the $\ell$th period of length $q$,

$$S_q'(\ell):=[\ell q, (\ell + 1)q) \times \mathfrak{Y}' \times [0, \infty), \quad \mathcal{P}(q, \ell):=\{i \in \{1, \ldots, \mathcal{I}\} : I_q(x^{(i)}, \ell) \neq \emptyset\}.$$

Note that in practice the measure can only be calculated for $\tau_1 \leq \tau$ (with $\tau$ the most recent date for which data is available) and on selected tests sets (for instance, in a back-testing approach). In controlled simulation experiments, see Sect. 7, we may and will use $\tau_1 = \infty$, thereby aiming at predicting the total number of unreported claims for each policy. Moreover, for using the error measures with the predictors from Sect. 4

and 5, the evaluation period $q$ must be a multiple of the homogeneity period length $p$ (unless one is willing to use the extension discussed in Remark 4.1).

The quality of individual claim count predictors may alternatively be assessed by first aggregating the individual predictions and then using standard global error measures; the predictors may then even be compared with classical methods for aggregated data like the standard chain ladder approach. Aggregated predictions are obtained from individual predictions straightforwardly: for $A = \mathfrak{X}' \times S$ with policies from $\mathfrak{X}' \subset \mathfrak{X}$ and claims from $S$ as above, let

$$\hat{N}_{\tau_0:\tau_1}(A) := \sum_{\substack{i \in \{1,\dots,\mathcal{I}\}: \\ x^{(i)} \in \mathfrak{X}'}} \hat{N}^{(i)}_{\tau_0:\tau_1}(S),$$

which is to be considered a predictor for the aggregated claim number

$$N_{\tau_0:\tau_1}(A) := \sum_{\substack{i \in \{1,\dots,\mathcal{I}\}: \\ x^{(i)} \in \mathfrak{X}'}} \xi^{(i)}(S \cap R_{\tau_0:\tau_1}).$$

For $q$ and $\mathfrak{Y}'$ as in (15), we then define

$$\mathrm{RMSE}_{\tau_0:\tau_1}(\mathfrak{Y}', q) := \left( \frac{q}{\tau} \sum_{\ell=0}^{\tau/q-1} \left\{ \hat{N}_{\tau_0:\tau_1}(A_{q,\ell,\mathfrak{Y}'}) - N_{\tau_0:\tau_1}(A_{q,\ell,\mathfrak{Y}'}) \right\}^2 \right)^{\frac{1}{2}}, \qquad (16)$$

where $A_{q,\ell,\mathfrak{Y}'} := \mathfrak{X} \times [\ell q, (\ell+1)q) \times \mathfrak{Y}' \times [0, \infty)$ comprises all claims in the portfolio from $\mathfrak{Y}'$ that have occurred in the $\ell$th period of length $q$. Note that this measure has also been used in [7], Formula (20). Its application is limited to the constraints mentioned above for $\tau_1$ and $q$.

# 7 Simulation study

In this section, we will study the performance of the new estimators and predictors within nine different simulation scenarios taken from [7]. We start by restating a brief, partially verbatim summary of the simulation models taken from the last-named paper:

The underlying portfolios build upon the car insurance data set described in Appendix A in [23]. The latter data set provides claim counts for 500,000 insurance policies, where each policy is associated with the risk features

<p align="center"><code>(age, ac, power, gas, brand, area, dens, ct)</code>,</p>

which correspond to age of driver, age of car, power of car, fuel type of car, brand of car, and area code, respectively; see also (A.1) in [23] for further details. Next to that, the data set also provides the variable `truefreq`, which corresponds to the claim intensity $\lambda(x)$ in our model.

Each portfolio is considered over ten periods of 365 days, that is, the portfolio coverage period is the interval $[0, 3650]$. The different scenarios are as follows:

***The baseline scenario.*** The baseline scenario/portfolio is characterized by a homogeneous *exposure* as well as position-independent *claim intensity*, *occurrence process*, and *reporting process*. It may be considered the vanilla portfolio that practitioners often aim at by careful selection of considered risks and suitable transformations, e.g., adjustment for inflation. More precisely:

- *Exposure.* New risks arrive according to a homogeneous Poisson process with intensity $50,000/365 \approx 137$ and contracts run for exactly one year. Moreover, the portfolio starts with exactly 50,000 policies with $t_{\text{start}} = 0$ and with remaining contract duration that is uniform on $[0, 365]$. As a consequence, the total exposure is constant in expectation and we have $\mathcal{I} \sim 50,000 + \text{Poi}(500,000)$. Finally, for each risk in the portfolio we randomly draw (with replacement) risk features from the aforementioned data set from [23].
- *Claim Intensity.* The claim intensity $\lambda(t, x) = \lambda(x)$ is independent of $t$ and $t_{\text{start}}$ and given by the variable `truefreq` that belongs to the risk selected in the previous paragraph.
- *Occurrence Process.* The occurrence process is position-independent, i.e., $P_{Y|X=x,T=t} = P_{Y|X=x}$ for all $t$. We choose to work with two claim variables, $y = (\text{cc}, \text{severity})$, with claims code $\text{cc} \in \{\text{injury, material}\}$, and an initial continuous proxy for the severity of the claim $\text{severity} \in \mathbb{R}_+$ (this variable should not be confused with the final claim amount, which we do not assess in this paper at all). The claim feature distribution of `cc` is chosen to be a function of the policy features `ac`, `power`, and `dens` in such a way that material damages are more likely to occur in densely populated areas and with low-powered and newer cars (see Appendix D in [8] for details on the precise relationship). The initial claim severity distribution of `severity` is log-normal with $\sigma$ constant and with $\mu$ depending on `cc`, `brand`, `ac` and `power` in such a way that injury claims, especially with older high-powered cars have a higher initial severity estimate. Moreover, material damages for certain premium brands are also more severe. Again, details are provided in Appendix D in [8].
- *Reporting Process.* The reporting process is position-independent, i.e, $P_{D|X=x,T=t,Y=y} = P_{D|X=x,Y=y}$. We choose to work with the Blended Dirac-Erlang-Generalized Pareto distribution from [7], with parameters specified in such a way that claims with higher initial severity, material claims with new cars, and claims with younger drivers in populated areas will be reported sooner, while low initial severity injuries will be reported later; see Appendix D in [8] for details.

***Eight non-homogeneous scenarios.***

Eight non-homogeneous scenarios are obtained by altering a single element of the baseline scenario:

1. *Exposure*: The distribution of `ac` changes continuously (drift) or abruptly (shock).
2. *Intensity*: $\lambda(x, t)$ decreases continuously (drift) or abruptly (shock).
3. *Occurrence*: The distribution of `cc` changes continuously (drift) or abruptly (shock).
4. *Reporting delay*: The distribution of $D$ is altered by moving probability mass to shorter reporting delays, continuously (drift) or abruptly (shock).
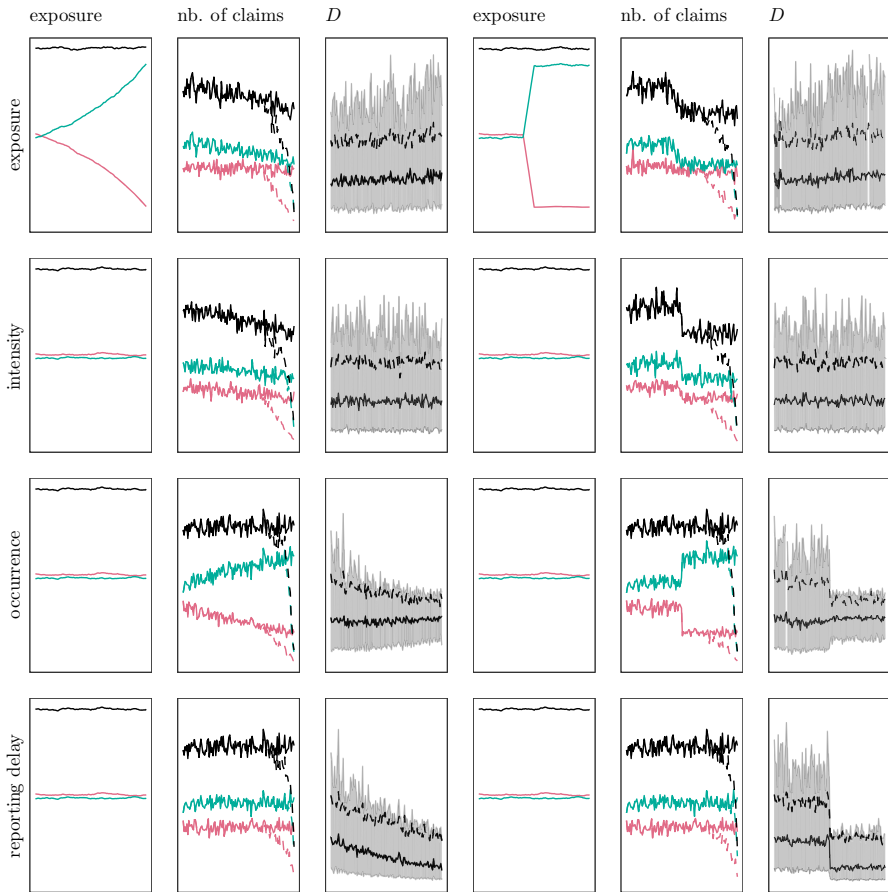
**Fig. 3** Overview of all scenarios, taken from Figure 4 in [7]. Rows show Scenarios 1-4, the left three columns showing the drift variation and the right three columns showing the shock variation. Within the scenarios, the panels show, from left to right, the exposure at risk (aggregated and split by `ac ≤ 5` shown in red and `ac > 5` shown in blue), the number of claims (aggregated and split by `cc` with injury shown in red and material shown in blue; dashed line: reported, solid line: occurred), the reporting delay distribution (dashed line: mean, solid line: median, ribbon: first and third quartiles)

Figure 3 illustrates the effect of the different scenarios on exposure, claim counts and reporting delays. The precise functional relationships are documented in Appendix D in [8].

The simulation study was conducted with 50 data seeds for each of the 9 scenarios, i.e., with 450 simulated portfolio datasets in total.

## 7.1 Training procedure

In this section, we describe details on the training and model selection procedure for the various neural networks used in the predictors. All networks were trained for 5, 000

epochs using the adam optimizer with fixed parameters $\alpha = 0.05$, $\beta_0 = \beta_1 = 0$ and an adaptive learning rate, halving the learning rate on plateaus (patience $= 2$) down to a minimum of $\underline{\alpha} = 10^{-4}$. In addition to this, the available (truncated) data was randomly split into 75% training data and 25% validation data. The validation data was not used for model calibration and instead kept aside to assess the generalization error.

### 7.1.1 Estimating the claim arrival process

Fitting of the claim arrival process was done in four steps, each requiring a slightly different neural network architecture. First, $\hat{P}_{D|x,t,y}$ was estimated as described in [7], compare Sect. 3.1, where we use the correct Blended Dirac-Erlang-Generalized Pareto distribution family for $\mathcal{P}^D$ (with unknown parameters). Reporting delay networks were trained for 100 starting seeds for parameter initialization. The top 10 performing reporting delay networks were chosen by computing $\mathrm{RMSE}_{\tau-365:\tau}(\mathfrak{Y}, 365; \mathfrak{D}_{\tau-365})$ for the predictor based on $\hat{P}_{D|x,t,y_1,y_2}$ (denoted NNet in Sect. 7.2), i.e., by the back-testing error for one year in the past.

Next, as described in Sect. 3.2, for each of the 10 networks from the previous step, coordinate distributions for $\hat{P}_{Y|x,t}$ were estimated from the decomposition $\mathfrak{Y} = \mathfrak{Y}_1 \times \mathfrak{Y}_2$ with $\mathfrak{Y}_1 = \{\text{injury, material}\}$ describing the claims code and $\mathfrak{Y}_2 = \mathbb{R}_+$ describing the initial severity estimate. First, $\hat{P}_{Y_2|x,y,y_1}$ was estimated using an MLP for the two parameters of a log-normal distribution, i.e., $\hat{P}_{Y_2|x,t,y_1} = \log \mathcal{N}(g^{(2)}(x, t, y_1))$. The network architecture $\mathcal{G}^{(2)}$ for $g^{(2)}(x, t, y_1)$ consisted of a $(10, 5)$ MLP with a softplus activation function adapted to an output in $\mathbb{R} \times (0, \infty)$, matching the two parameters $(\mu, \sigma)$ defining a log-normal distribution. For each of the 10 reporting delay networks, 10 starting seeds were used for training the initial severity feature network $g^{(2)}(x, t, y_1)$, resulting in a total of 100 estimates for the pair $(\hat{P}_{D|x,t,y_1,x_2}, \hat{P}_{Y_2|x,t,y_1})$. Similar as in the previous step, the ten best estimates were chosen based on back-testing the error one year in the past, using the predictor based on $\hat{P}_{D|x,t,y_1,x_2}$ and $\hat{P}_{Y_2|x,t,y_1}$ (denoted NNet$_{\texttt{severity}}$ in Sect. 7.2). It should be noted that this predictor performed worse than the underlying NNet predictor based solely on $\hat{P}_{D|x,t,y_1,y_2}$. Nonetheless, using $\hat{P}_{Y_2|x,t,y_1}$ is necessary for the subsequent steps.

For each of the ten estimates for $(\hat{P}_{D|x,t,y_1,x_2}, \hat{P}_{Y_2|x,t,y_1})$ from the previous step, we next estimated $\hat{P}_{Y_1|x,t}$. The associated network architecture $\mathcal{G}^{(1)}$ consisted of a $(10, 5)$ MLP with a softplus activation function outputting probability masses for a discrete distribution on $\{\text{injury, material}\}$. It was trained for 10 different starting seeds, resulting in a total of 100 estimates for $(\hat{P}_{D|x,t,y_1,y_2}, \hat{P}_{Y_2|x,t,y_1}, \hat{P}_{Y_1|x,t})$ and hence for $(\hat{P}_{D|x,t,y}, \hat{P}_{Y|x,t})$ using the definition $\hat{P}_{Y|x,t}(\mathrm{d}y_1, \mathrm{d}y_2) = \hat{P}_{Y_1|x,t}(\mathrm{d}y_1) \hat{P}_{Y_2|x,t,y_1}(\mathrm{d}y_2)$. Again, the 10 best estimates were chosen by evaluating the associated predictor (denotes NNet$_{\texttt{cc}}$ in Sect. 7.2) using the backtesting error $\mathrm{RMSE}_{\tau-365:\tau}(\mathfrak{Y}, 365; \mathfrak{D}_{\tau-365})$.

Finally, for each of the ten estimates for $(\hat{P}_{D|x,t,y}, \hat{P}_{Y|x,t})$, ten intensity estimates $\hat{\lambda}^{\mathrm{NNet}}(x, t)$ were obtained as described in Sect. 3.3, with ten different starting seeds. The underlying network architecture consisted of a $(10, 5)$ MLP with a softplus activation function and a parameter-free skip connection for the offset term as described

in [22], leading to a single Poisson parameter in $\mathbb{R}_+$. After training, the bias regularization method described in [22] was applied. From the resulting 100 estimates for $(\hat{P}_{D|x,t,y}, \hat{P}_{Y|x,t}, \hat{\lambda}^{\mathrm{NNet}}(x,t))$, the final estimate was chosen according to the backtesting error for its associated predictor, denoted $\mathrm{NNet}_{\mathrm{FreqNet}}$ in Sect. 7.2.

Overall, the number of trained networks for each data set is 400, resulting in a total of $450 \times 400 = 180,000$ trained networks for the simulation study.

### 7.1.2 Fitting chain ladder networks

Training a chain ladder network requires fitting $M$ neural networks, $g^{\mathfrak{Y}_m}$ for $m = 1, \ldots, M$. As described in next section, we use both $M = 1$ (predictor CLFreqNet in Sect. 7.2) and $M = 2$ (CLFreqNet$_{\mathrm{cc}}$ in Sect. 7.2), resulting in three networks to be trained for each data set. The MLP architecture was fixed as $(10, 5)$ with a softplus activation function and a parameter-free skip connection for the offset term. After training, the bias regularization method described in [22] was applied using the training dataset. For each data set, ten starting seeds were used, resulting in 30 networks for each data set, from which a best predictor was chosen based on the backtesting error $\mathrm{RMSE}_{\tau-365:\tau}(\mathfrak{Y}, 365; \mathfrak{D}_{\tau-365})$. For the entire simulation study, $450 \times 30 = 13,500$ networks were trained.

### 7.2 Predictors

We provide a detailed overview of the predictors, tailored to the specific portfolios described at the beginning of Sect. 7. For the macro-level error measure from (16), we will compare a total of eight different predictors, three of which provide reasonable micro-level predictions as measured by (15). All predictors target the number of claims in $A = \mathfrak{X}' \times I_p(k) \times \mathfrak{Y}' \times [0, \infty)$ with some $\mathfrak{X}' \subset \mathfrak{X}$, $\mathfrak{Y}' \subset \mathfrak{Y}$ and $I_p(k)$ the $k$th period of length $p$.

For index m encoding one of the methods specified below, let

$$\hat{N}_{\mathrm{m}}^{\mathrm{cw}}(A; \mathfrak{D}_\tau) := \sum_{(x,t,y,d) \in A \cap \mathfrak{D}_\tau} \hat{c}_{\mathrm{m}}(x, t, y, d), \qquad \hat{N}_{\mathrm{m}}^{\mathrm{pw}}(A; \mathfrak{D}_\tau) := \sum_{\substack{i \in \{1, \ldots, \mathcal{I}\}: \\ x^{(i)} \in \mathfrak{X}'}} \hat{c}_{\mathrm{m}}(i, A),$$

where the upper index cw and pw stand for *claim-wise* and *policy-wise*, respectively, and where $\hat{c}_{\mathrm{m}}(x, t, y, d)$ and $\hat{c}_{\mathrm{m}}(i, A)$ are suitable numbers, additionally depending on $\mathfrak{D}_\tau$, $\tau_0$ and $\tau_1$, as specified below. For $k \in \mathbb{N}_0$, $t \geq 0$ and $x \in \mathfrak{X}$, recall the notations $t_k = kp + \frac{p}{2}$, $k_t = \lfloor \frac{t}{p} \rfloor$ and $I_p(x, k) = C(x) \cap [kp, (k+1)p)$ with $C(x)$ the coverage period of policy $x$, see (8).

- **Predictor NNet**. We consider the original method from [7] that only relies on modeling and estimating reporting delays, see formula (19) in that paper. More precisely, we define $\hat{N}_{\mathrm{NNet}} := \hat{N}_{\mathrm{NNet}}^{\mathrm{cw}}$ with constants

$$\hat{c}_{\mathrm{NNet}}(x, t, y, d) := \frac{\int_{I_p(x, k_t)} \hat{P}_{D|X=x, T=t_{k_t}, Y=y}((\tau_0 - s, \tau_1 - s]) \, ds}{\int_{I_p(x, k_t)} \hat{P}_{D|X=x, T=t_{k_t}, Y=y}([0, \tau - s]) \, ds}.$$

- **Predictor NNet$_{\text{severity}}$.** We additionally incorporate the estimated first stage claim feature model from Sect. 3.2, based upon the decomposition $\mathfrak{Y} = \{\text{injury, material}\} \times \mathbb{R}_+ =: \mathfrak{Y}_1 \times \mathfrak{Y}_2$ into claims code and initial claim severity. More precisely, we define $\hat{N}_{\text{NNet}_{\text{severity}}} := \hat{N}_{\text{NNet}_{\text{severity}}}^{\text{cw}}$ with

$$\hat{c}_{\text{NNet}_{\text{severity}}}(x, t, y, d) := \frac{\hat{p}_{\text{rep;severity}}(\tau_0, \tau_1, x, k_t, y_1)}{\hat{p}_{\text{rep;severity}}(0, \tau, x, k_t, y_1)}.$$

with $\hat{p}_{\text{rep;severity}}(\tau_0, \tau_1, x, k_t, y_1)$ defined as

$$\int_{\mathfrak{Y}'_{y_1}} \int_{I_p(x,k_t)} \hat{P}_{D|X=x,T=t_{k_t},Y_1=y_1,Y_2=w}((\tau_0 - s, \tau_1 - s])\, \mathrm{d}s$$
$$\mathrm{d}\hat{P}_{Y_2|X=x,T=t_{k_t},Y_1=y_1}(w),$$

where $\mathfrak{Y}'_{y_1} = \{w : (y_1, w) \in \mathfrak{Y}'\}$.

- **Predictor NNet$_{\text{cc}}$.** This predictor is built on the full estimated claim feature model, see Sect. 3.2. More precisely, we define $\hat{N}_{\text{NNet}_{\text{cc}}} := \hat{N}_{\text{NNet}_{\text{cc}}}^{\text{cw}}$

$$\hat{c}_{\text{NNet}_{\text{cc}}}(x, t, y, d) := \frac{\hat{p}_{\text{rep;cc}}(\tau_0, \tau_1, x, k_t)}{\hat{p}_{\text{rep;cc}}(0, \tau, x, k_t)}$$

with $\hat{p}_{\text{rep;cc}}(\tau_0, \tau_1, x, k_t)$ defined as

$$\int_{\mathfrak{Y}'} \int_{I_p(x,k_t)} \hat{P}_{D|X=x,T=t_{k_t},Y=y}((\tau_0 - s, \tau_1 - s])\, \mathrm{d}s\, \mathrm{d}\hat{P}_{Y|X=x,T=t_{k_t}}(y).$$

For $\mathfrak{Y}' = \mathfrak{Y}$, this can be written as

$$\sum_{u \in \{\text{injury,material}\}} \int_0^\infty \int_{I_p(x,k_t)} \hat{P}_{D|X=x,T=t_{k_t},Y_1=u,Y_2=w}((\tau_0 - s, \tau_1 - s])\, \mathrm{d}s$$
$$\mathrm{d}\hat{P}_{Y_2|X=x,T=t_{k_t},Y_1=u}(w)\, \hat{P}_{Y_1|X=x,T=t_{k_t}}(\{u\}).$$

- **Predictor FreqNet.** This predictor is the one from Sect. 4 that builds upon the full estimated model for the claim arrival process. More precisely, $\hat{N}_{\text{FreqNet}} = \hat{N}_{\text{FreqNet}}^{\text{pw}}$ with

$$\hat{c}_{\text{FreqNet}}(i, A) := \xi_r^{(i)}(I_p(k) \times \mathfrak{Y}' \times [0, \infty)) + \hat{\lambda}^{\text{NNet}}(x^{(i)}, t_k)$$
$$\int_{\mathfrak{Y}'} \int_{I_p(x^{(i)},k)} \hat{P}_{D|X=x^{(i)},T=t_k,Y=y}((\max(\tau, \tau_0) - s, \tau_1 - s])\, \mathrm{d}s\, \mathrm{d}\hat{P}_{Y|X=x^{(i)},T=t_k}(y),$$

which corresponds to (12).

- **Predictor CL**. This predictor is the basic chain ladder predictor. More precisely, $\hat{N}_{\mathrm{CL}} = \hat{N}_{\mathrm{CL}}^{\mathrm{cw}}$ with

$$\hat{c}_{\mathrm{CL}}(x, t, y, d) = \mathbf{1}(\bar{\tau}_0 \leq \lfloor (T + D)/p \rfloor \leq \bar{\tau}_1) + \mathrm{FtU}_{k_t + \bar{\tau} - \bar{\tau}_0}^{\mathfrak{Y}} - \mathrm{FtU}_{k_t + \bar{\tau} - \bar{\tau}_1}^{\mathfrak{Y}},$$

  where $\tau = \bar{\tau} p$, $\tau_0 = \bar{\tau}_0 p$ and $\tau_1 = \bar{\tau}_1 p$ must be whole multiples of the development period and $\mathrm{FtU}_k^{\mathfrak{Y}} := 1$ for $k < 0$.

- **Predictor CLFreqNet**. This is the basic chain ladder network based predictor, and is only defined for $\mathfrak{Y}' = \mathfrak{Y}$. More precisely, $\hat{N}_{\mathrm{CLFreqNet}} = \hat{N}_{\mathrm{CLFreqNet}}^{\mathrm{pw}}$ with

$$\hat{c}_{\mathrm{CLFreqNet}}(i, A) := \xi_r^{(i)}(I_p(k) \times \mathfrak{Y} \times [0, \infty) \cap R_{\tau_0:\tau_1})$$

$$+ \mathrm{Leb}(I_p(x^{(i)}, k)) \hat{\lambda}^{\mathrm{CL}, \mathfrak{Y}}(x^{(i)}, t_k) \left( \frac{1}{\mathrm{FtU}_{k + \bar{\tau} - \bar{\tau}_1}^{\mathfrak{Y}}} - \frac{1}{\mathrm{FtU}_{k + \bar{\tau} - \bar{\tau}_0}^{\mathrm{CL}, \mathfrak{Y}}} \right).$$

  This formula comes from (14) with the trivial partition using $M = 1$ component.

- **Predictor $\mathrm{CL}_{\mathrm{cc}}$**. This predictor is the chain ladder predictor based on splitting by claims code $\mathfrak{Y} = \mathfrak{Y}_1^{\mathrm{cc}} \cup \mathfrak{Y}_2^{\mathrm{cc}} := \{\mathrm{injury}\} \times \mathbb{R}_+ \cup \{\mathrm{material}\} \times \mathbb{R}_+$. More precisely, $\hat{N}_{\mathrm{CL}_{\mathrm{cc}}} = \hat{N}_{\mathrm{CL}_{\mathrm{cc}}}^{\mathrm{cw}}$ with

$$\hat{c}_{\mathrm{CL}_{\mathrm{cc}}}(x, t, y, d) = \mathbf{1}(\bar{\tau}_0 \leq \lfloor (T + D)/p \rfloor \leq \bar{\tau}_1)$$

$$+ (\mathrm{FtU}_{k_t + \bar{\tau} - \bar{\tau}_0}^{\mathfrak{Y}_1^{\mathrm{cc}}} - \mathrm{FtU}_{k_t + \bar{\tau} - \bar{\tau}_1}^{\mathfrak{Y}_1^{\mathrm{cc}}}) \mathbf{1}(y_1 = \mathrm{injury})$$

$$+ (\mathrm{FtU}_{k_t + \bar{\tau} - \bar{\tau}_0}^{\mathfrak{Y}_2^{\mathrm{cc}}} - \mathrm{FtU}_{k_t + \bar{\tau} - \bar{\tau}_1}^{\mathfrak{Y}_2^{\mathrm{cc}}}) \mathbf{1}(y_1 = \mathrm{material}).$$

- **Predictor $\mathrm{CLFreqNet}_{\mathrm{cc}}$**. This is the chain ladder network based predictor for the partition $\mathfrak{Y} = \mathfrak{Y}_1^{\mathrm{cc}} \cup \mathfrak{Y}_2^{\mathrm{cc}}$ defined in the description of $\mathrm{CL}_{\mathrm{cc}}$. It is only defined for $\mathfrak{Y}' \in \{\mathfrak{Y}_1^{\mathrm{cc}}, \mathfrak{Y}_2^{\mathrm{cc}}, \mathfrak{Y}\}$. More precisely, $\hat{N}_{\mathrm{CLFreqNet}_{\mathrm{cc}}} = \hat{N}_{\mathrm{CLFreqNet}_{\mathrm{cc}}}^{\mathrm{pw}}$ with

$$\hat{c}_{\mathrm{CLFreqNet}_{\mathrm{cc}}}(i, A) := \xi_r^{(i)}(I_p(k) \times \mathfrak{Y}' \times [0, \infty) \cap R_{\tau_0:\tau_1})$$

$$+ \mathrm{Leb}(I_p(x^{(i)}, k)) \sum_{\substack{\mathbf{Y} \in \{\mathfrak{Y}_1^{\mathrm{cc}}, \mathfrak{Y}_2^{\mathrm{cc}}\} \\ \mathbf{Y} \subset \mathfrak{Y}'}} \hat{\lambda}^{\mathrm{CL}, \mathbf{Y}}(x^{(i)}, t_k) \left( \frac{1}{\mathrm{FtU}_{k + \bar{\tau} - \bar{\tau}_1}^{\mathbf{Y}}} - \frac{1}{\mathrm{FtU}_{k + \bar{\tau} - \bar{\tau}_0}^{\mathrm{CL}, \mathbf{Y}}} \right).$$

  The formula again stems from applying (14), this time with $M = 2$ and the partition by claim code.

- **Predictor cheating**. This is the predictor using the true parameters of the simulated model for prediction; it is not available in practice and only serves as a benchmark for evaluating the other predictors. More precisely, $\hat{N}_{\mathrm{cheating}} = \hat{N}_{\mathrm{cheating}}^{\mathrm{pw}}$ with
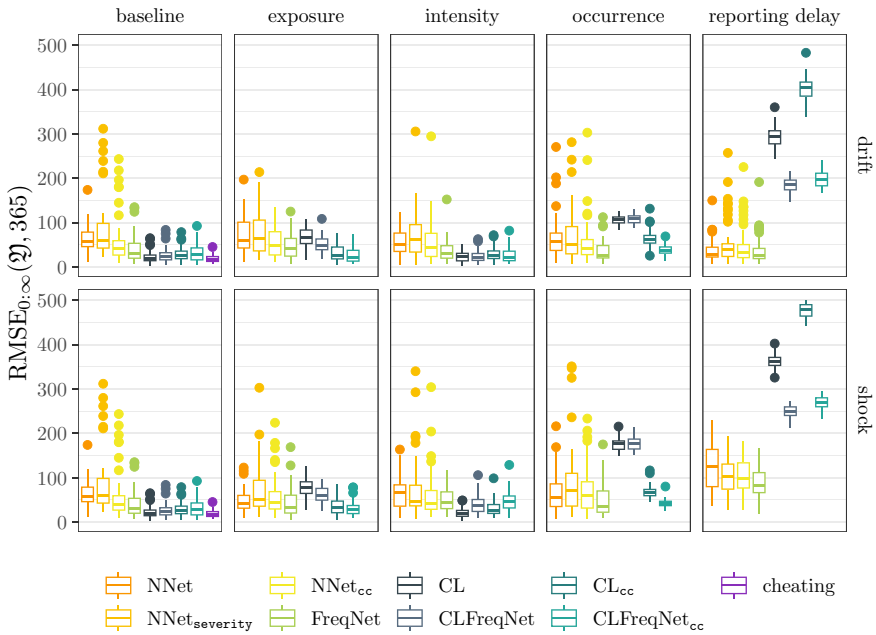
**Fig. 4** Boxplots of the overall error measure $\text{RMSE}_{0:\infty}(\mathfrak{Y}, 365)$, each based on $n = 50$ simulated paths. Legend is shown in column major order

$$\hat{c}_{\text{cheating}}(i, A) = \xi_r^{(i)}(I_p(k) \times \mathfrak{Y}' \times [0, \infty) \cap R_{\tau_0:\tau_1})$$

$$+\lambda(x^{(i)}, t_k) \int_{\mathfrak{Y}'} \int_{I_p(x^{(i)},k)} P_{D|X=x^{(i)},T=t_k,Y=y}(I_{\tau_0:\tau_1}(\tau, s)) \, \mathrm{d}s \, \mathrm{d}P_{Y|X=x^{(i)},T=t_k}(y),$$

corresponding to (12) with estimated distributions replaced by true distributions.

## 7.3 Results

Throughout the simulation study, we use an evaluation period of $q = 365$ days. Figure 4 shows partly the same results as Figure 5 in [7], extended by the methods described in this paper and using the same color keys and predictor names, if applicable. The underlying error measure is the one from (16). Regarding the baseline scenario, we can see that modelling more and more parts of the claim arrival process, i.e., going from NNet to $\text{NNet}_{\text{cc}}$ and then finally to FreqNet, reduces the overall error with $\text{NNet}_{\text{cc}}$ seemingly exhibiting slightly larger variance. Only applying a partial model for the distribution of $Y$ as in $\text{NNet}_{\text{severity}}$ increases the prediction error and its variance. We can also see that the chain ladder predictor provides close-to-optimal predictions on par with those obtained from the true model in this setting where the underlying chain ladder assumptions are exactly met.

For the chain ladder based methods, the error of the neural network predictors in the baseline scenario increases when compared to the pure factor based prediction—at the advantage of providing individual reserve predictions for each policy in the portfolio.

This behavior can be explained by the training method used: all neural network fitting procedures with a Poisson loss use the GLM skip connection described by [22], but only on the 75% of the available data chosen for training. The 25% of the data used for hold-out validation therefore did not take part in the bias regularization, whereas the factor based methods had no hold-out data. If bias regularization was done on 100% of the data, the difference in errors would be smaller but not zero, because the bias regularization only ensures the total number of claims to remain constant, but not their allocation to accident years.

The results for the exposure scenario show that predictions can be improved when using portfolio information (in particular, exposure data); in fact, we see this improvement across all three approaches, i.e., NNet $\rightsquigarrow$ FreqNet, CL $\rightsquigarrow$ CLFreqNet and $CL_{cc} \rightsquigarrow CLFreqNet_{cc}$. Heuristically, this can be explained by the fact that the drift in exposure is directly reflected by a drift in expected claim counts from the claim intensity models (see Fig. 3), thereby influencing IBNR claim counts. The observed improvement is most prominent for the unpartitioned chain ladder approach, because the other basic methods can at least partially detect the changes via changes in the distribution of cc, which is also influenced by the exposure shift.

Changes in the claim intensity make it harder to train the underlying intensity, $\lambda(x, t)$. Due to this disadvantage, one might expect to see a deterioration in prediction error for the intensity based approaches. Surprisingly, this is only found to be the case for the intensity shock scenario, and only so for the chain ladder based CLFreqNet and $CLFreqNNet_{cc}$. The intensity drift exhibits no noticeable deterioration in error and FreqNet shows a smaller improvement in error when confronted with an intensity shock compared to NNet, but an improvement nonetheless.

Regarding drifts and shocks in the occurrence process (i.e. in the distribution of cc), we do not observe a substantial effect on the prediction errors of the NNet based approach (i.e., they are similar as in the baseline scenario). Unpartitioned chain ladder does not deal well with these changes to the claims process and the intensity based extension doesn't manage to reduce the problem. Substantial improvements are found when moving from NNet to FreqNet and from $CL_{cc}$ to $CLFreqNet_{cc}$.

When the reporting delay distribution changes, chain ladder based methods start to perform very badly, even with partitioning. Since the introduced change effectively reduced the time-to-report, plain chain ladder approaches are confronted with higher claim counts upfront, which amplifies the error due to the multiplicative structure. This effect is dampened by intensity based extensions, because here the expected number of IBNR claims is based on the expected (long-term) intensity and not on short-term observations. Again, FreqNet shows a similar improvement compared to NNet as seen in other scenarios.

In summary, we can see that FreqNet performs well across all scenarios, improving on the method developed in [7]. The robustness to changes in exposure and reporting delays of chain ladder based estimates can be improved at little cost to overall accuracy by employing the CLFreqNet method.

We will now move our attention to the individual level results as measured by $RMSE_{0:\infty}^{expo}$ defined in (15), which are summarized in Fig. 5. Within that figure, we do not display results for straightforward individual factor based predictions (i.e.,
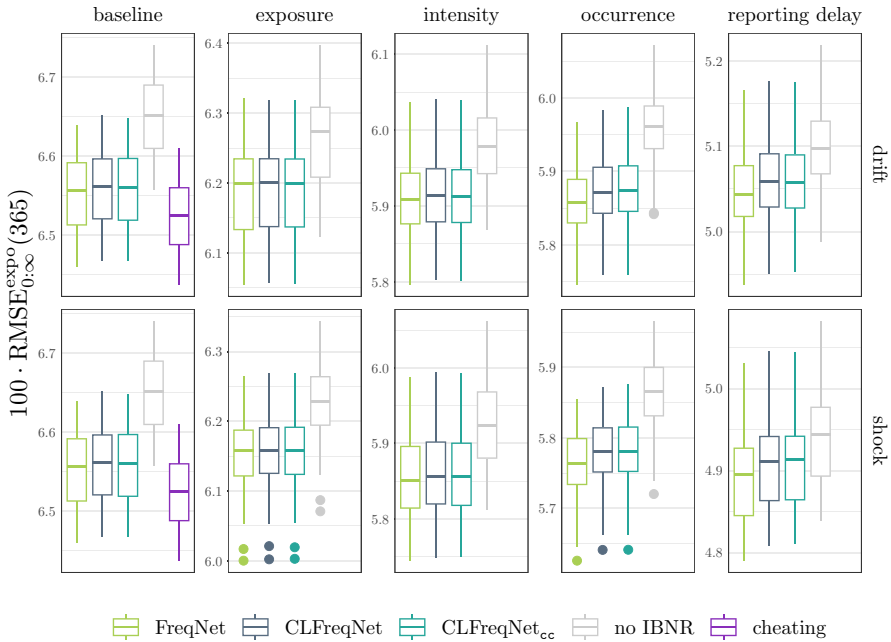
**Fig. 5** Boxplots of the individual-level error measure $\text{RMSE}_{0:\infty}^{\text{expo}}(365)$, each based on $n = 50$ simulated paths. The trivial "predictor" $\hat{N}^{\text{no IBNR}}(A) := N_r(A)$, predicting no IBNR claims, is also shown. Note that the predictors $\hat{N}^{\text{CL}}$ and $\hat{N}^{\text{CL}_{\text{CC}}}$ are not suitable for individual level claim count predictions as their $\text{RMSE}_{0:\infty}^{\text{expo}}$ is worse than $\hat{N}^{\text{no IBNR}}$. In the baseline scenario, $\text{RMSE}_{0:\infty}^{\text{expo}}(365)$ has a median of $8.45 \cdot 10^{-2}$ for $\hat{N}^{\text{CL}}$ and $10.8 \cdot 10^{-2}$ for $\hat{N}^{\text{CL}_{\text{CC}}}$, compared to $6.65 \cdot 10^{-2}$ for $\hat{N}^{\text{no IBNR}}$

multiplying the number of reported claims on an individual level by a factor-to-ultimate) because of their generally poor performance: for instance, in the baseline scenario, the mean $\text{RMSE}_{0:\infty}^{\text{expo}}(365)$ for CL, $\text{CL}_{\text{cc}}$ and noIBNR is 0.0843, 0.108 and 0.0665, respectively, where noIBNR refers to simply predicting no IBNR claims at all. However, the results in Fig. 5 show that the chain ladder based neural network predictors CLFreqNet and CLFreqNet$_{\text{cc}}$ provide viable solutions for allocating the IBNR claims from a chain ladder triangle to individual policies, albeit without yielding a full distributional model. In general, CLFreqNet and CLFreqNet$_{\text{cc}}$ exhibit very similar errors whereas FreqNet shows slightly smaller errors than the other two methods. Comparing the mean $\text{RMSE}_{0:\infty}^{\text{expo}}(365)$ for the methods noIBNR(0.0665), CLFreqNet(0.0656), CLFreqNet$_{\text{cc}}$(0.0656), FreqNet(0.0655) and cheating(0.0653), we see that the chain ladder based methods score 72% of the performance of cheating when compared to noIBNR and FreqNet even achieves 78% of the improvement from noIBNR to cheating. It is interesting to note that the results are quite similar for all five scenarios, with the reporting delay scenario exhibiting the smallest difference between the noIBNR error and the error of the other three methods.
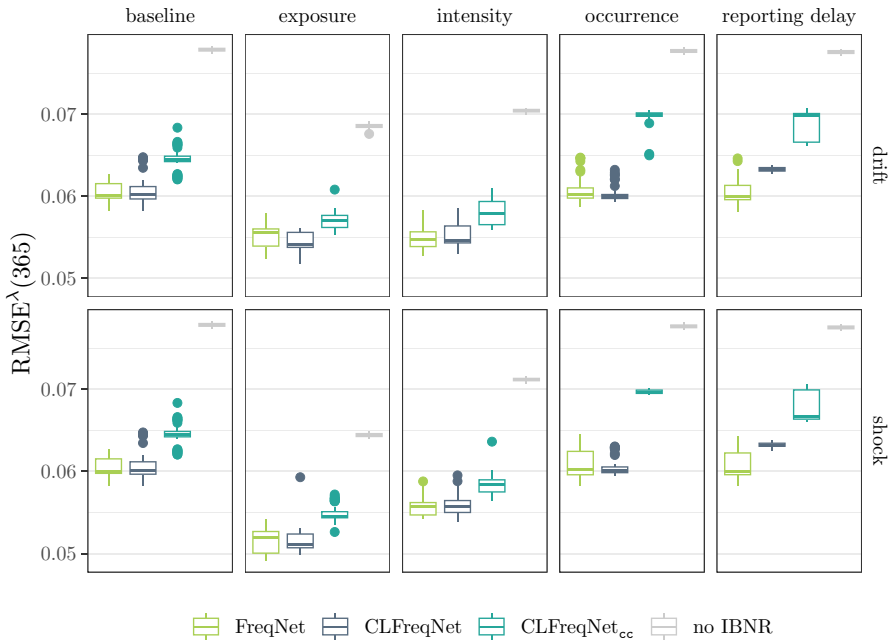
**Fig. 6** Boxplots of the intensity error measure $\text{RMSE}^\lambda(365)$, each based on $n = 50$ simulated paths. For no IBNR, the estimate $\hat{\lambda}(x, t) \equiv \text{const} = \hat{N}_{0:\infty}(\mathfrak{X} \times [0, \tau) \times \mathfrak{Y} \times [0, \infty)) / \sum_{i=1}^{\mathcal{I}} \text{Leb}(C(x^{(i)}) \cap [0, \tau])$ was used

Of primal importance in insurance pricing is an accurately estimated risk model (which includes the intensity $\lambda$), for instance for reducing or avoiding cross-subsidisation within a portfolio. In Fig. 6 we study the quality of the estimated intensity models obtained for the methods FreqNet, CLFreqNet and CLFreqNet$_{\text{cc}}$. As a measure for the quality of the estimates, we use, for some evaluation period length $q$ which is divisor of $\tau$ (as before, we fix $q = 365$ throughout),

$$\text{RMSE}^\lambda(q)$$

$$:= \left( \frac{1}{\sum_{j=0}^{\tau/q-1} \#\mathcal{P}(q, j)} \sum_{\ell=0}^{\tau/q-1} \sum_{i \in \mathcal{P}(q,\ell)} \left( \hat{\lambda}(x^{(i)}, (\ell + \tfrac{1}{2})q) - \lambda(x^{(i)}, (\ell + \tfrac{1}{2})q) \right)^2 \right)^{\frac{1}{2}},$$

(17)

where we have used the notation from Sect. 6. Note that, when using $\hat{\lambda}(x, t) \equiv \text{const} = \hat{N}_{0:\infty}(\mathfrak{X} \times [0, \tau) \times \mathfrak{Y} \times [0, \infty)) / \sum_{i=1}^{\mathcal{I}} \text{Leb}(C(x^{(i)}) \cap [0, \tau])$ as an estimate for $\lambda$ using the predictors CL and CL$_{\text{cc}}$ for $\hat{N}_{0:\infty}$ yields very similar $\text{RMSE}^\lambda(365)$ as noIBNR, so they were left out of the plots in Fig. 6 for readability. As an example, the baseline scenario has a mean $\text{RMSE}^\lambda(365)$ of 0.0737 for noIBNR and of 0.0734 for both CL and CL$_{\text{cc}}$ whereas the intensity networks yield values from 0.0578 to 0.0632.

A priori, one would expect that $\text{RMSE}^\lambda$ correlates with $\text{RMSE}_{0:\infty}^{\text{expo}}$, since both are error measures at the individual policy level. Surprisingly, the results in Fig. 6 show

that this correlation breaks down for the chain ladder-based intensity models: while $\text{RMSE}_{0:\infty}^{\text{expo}}$ is very similar for CLFreqNet and CLFreqNet$_{\text{CC}}$, $\text{RMSE}^{\lambda}$ is smaller for CLFreqNet than it is for CLFreqNet$_{\text{CC}}$ in all scenarios. Heuristically, a larger variance of CLFreqNet$_{\text{CC}}$ may be explained by the fact that CLFreqNet$_{\text{CC}}$ is the only method of the three that uses two independent networks for the intensity of each claim code, and hence is based on twice the number of parameters. It is not fully clear however why this would deteriorate model quality. Another interesting observation is that despite FreqNet having a worse accident year level error (Fig. 4), its underlying intensity model is comparable in quality to that of CLFreqNet in the baseline scenario.

## 8 Application to real data

In this section, we will apply the different methods to a large real dataset containing motor legal insurance claims provided by a German insurance company. The dataset is described in Sect. 8.1. More detail on the prediction methods and estimation procedure can be found in Sect. 8.2. Due to the nature of real world data, observations are only available for a limited time frame. Therefore, model performance metrics cannot use $\infty$ as the time of evaluation, but must instead use a finite cutoff date. We examined two artificial truncation points, $\tau = 31\text{stDecember}2017 \, and \, \tau = 31\text{stDecember}2018$ and evaluate predictions for one year into the future, i.e. $\text{RMSE}_{\tau:\tau+365}(\mathfrak{Y}, 365)$ and $\text{RMSE}_{\tau:\tau+365}^{\text{expo}}(\mathfrak{Y}, 365)$. Results of this examination are presented and discussed in Sect. 8.3.

### 8.1 The dataset

The dataset is the same as [7]. It contains a portfolio of about 250,000 motor legal insurance contracts and 65,000 corresponding claims with exposure and claims information available monthly from 31st December 2014 to 31st December 2020. Due to the extreme shock the COVID-19 pandemic had on the dataset, we chose to only consider data available up to 31st December 2019 for model evaluation. For a more detailed description of the data, refer to [7].

### 8.2 Predictors

In this section, we provide a detailed overview of the predictors that we apply to the dataset described in Sect. 8.1. For the macro-level error measure, $\text{RMSE}_{\tau:\tau+365}(\mathfrak{Y}, 365)$, we will compare a total of six predictors, three of which can provide viable micro-level predictors. The micro-level predictors are compared using $\text{RMSE}_{\tau:\tau+365}^{\text{expo}}(\mathfrak{Y}, 365)$. Most of the predictors are defined analogously to those in Sect. 7.2 and we will reuse the notation defined there.

- **Predictor NNet**. The original method from [7], formula (19).
- **Predictor NNet$_{\mathfrak{Y}}$**. This predictor is based on the estimated claim feature model. Since all claim features are discrete, this modelling step was done using a single discrete distribution with 360 different possible outcomes. More precisely,

recalling the notation $\hat{N}^{\mathrm{cw}}$ from Sect. 7.2, we define $\hat{N}_{\mathrm{NNet}_{\mathfrak{Y}}} := \hat{N}_{\mathrm{NNet}_{\mathfrak{Y}}}^{\mathrm{cw}}$ with

$$\hat{c}_{\mathrm{NNet}_{\mathrm{cc}}}(x, t, y, d) := \frac{\hat{p}_{\mathrm{rep};\mathfrak{Y}}(\tau_0, \tau_1, x, k_t)}{\hat{p}_{\mathrm{rep};\mathfrak{Y}}(0, \tau, x, k_t)}$$

and $\hat{p}_{\mathrm{rep};\mathfrak{Y}}(\tau_0, \tau_1, x, k_t)$ defined as

$$\sum_{y \in \mathfrak{Y}'} \int_{I_p(x,k_t)} \hat{P}_{D|X=x, T=t_{k_t}, Y=y}([\tau_0 - s, \tau_1 - s]) \, \mathrm{d}s \, \hat{P}_{Y|X=x, T=t_{k_t}}(y).$$

- **Predictor FreqNNet**. The predictor from (12).
- **Predictor CL**. This predictor is the basic chain ladder predictor.
- **Predictor CL$_{\mathrm{cc}}$**. This predictor is the chain ladder predictor based on splitting by claims code $\mathfrak{Y} = \mathfrak{Y}_0^{\mathrm{cc}} \cup \mathfrak{Y}_1^{\mathrm{cc}} \cup \mathfrak{Y}_2^{\mathrm{cc}} \cup \mathfrak{Y}_3^{\mathrm{cc}} \cup \mathfrak{Y}_4^{\mathrm{cc}}$. More precisely, recalling the notation $\hat{N}^{\mathrm{cw}}$ from Sect. 7.2, $\hat{N}_{\mathrm{CL}_{\mathrm{cc}}} = \hat{N}_{\mathrm{CL}_{\mathrm{cc}}}^{\mathrm{cw}}$ with

$$\hat{c}_{\mathrm{CL}_{\mathrm{cc}}}(x, t, y, d) = \mathbf{1}(\bar{\tau}_0 \leq \lfloor (T + D)/p \rfloor \leq \bar{\tau}_1)$$
$$+ \sum_{m=0}^{4} (\mathrm{FtU}_{k_t + \bar{\tau} - \bar{\tau}_0}^{\mathfrak{Y}_m^{\mathrm{cc}}} - \mathrm{FtU}_{k_t + \bar{\tau} - \bar{\tau}_1}^{\mathfrak{Y}_m^{\mathrm{cc}}}) \mathbf{1}(y_1 = m),$$

where $y_1 = \mathrm{cc}$.

- **Predictor CLFreqNNet$_{\mathrm{cc}}$**. This is the chain ladder network based predictor for the partition $\mathfrak{Y} = \mathfrak{Y}_0^{\mathrm{cc}} \cup \mathfrak{Y}_1^{\mathrm{cc}} \cup \mathfrak{Y}_2^{\mathrm{cc}} \cup \mathfrak{Y}_3^{\mathrm{cc}} \cup \mathfrak{Y}_4^{\mathrm{cc}}$ defined in the description of CL$_{\mathrm{cc}}$. The formula stems from applying (14) with $M = 5$ and the partition by claim code.

## 8.3 Results

The seven available predictors are evaluated for one year ahead and on an evaluation period of $q = 365$. In comparison to the simulation study, the **cheating** predictor is missing and the two plain chain ladder predictors have no uncertainty due to their deterministic algorithm. Also, because stepwise estimation of the claim feature distribution was not necessary, there is only one predictor **NNet$_{\mathfrak{Y}}$** based on the full claim feature distribution instead of the two predictors **NNet$_{\mathrm{severity}}$** and **NNet$_{\mathrm{cc}}$**.

Summarily, despite the fact that the model selection strategy has not been fine-tuned to the problem at hand and shows a rather unreliable performance overall, **FreqNet** shows promising results on a micro-level at an acceptable cost on the macro-level.

Figure 7 shows the accident year level prediction error $\mathrm{RMSE}_{\tau:\tau+365}(\mathfrak{Y}, 365)$ for one year ahead across the seven methods for the two artificial truncation points. In contrast to most simulation results on the ultimate accident year level prediction error, we see an increase in error of **NNet$_{\mathfrak{Y}}$** compared to **NNet**. This loss could possibly be overcome by optimizing the claim feature model architecture, which was fixed as a $(10, 5)$ feed-forward network for simplicity. For $\tau = 31\mathrm{stDecember}2017$ the distribution of errors for **FreqNet** also seems to deteriorate, whereas $\tau = 31\mathrm{st}$ December 2018
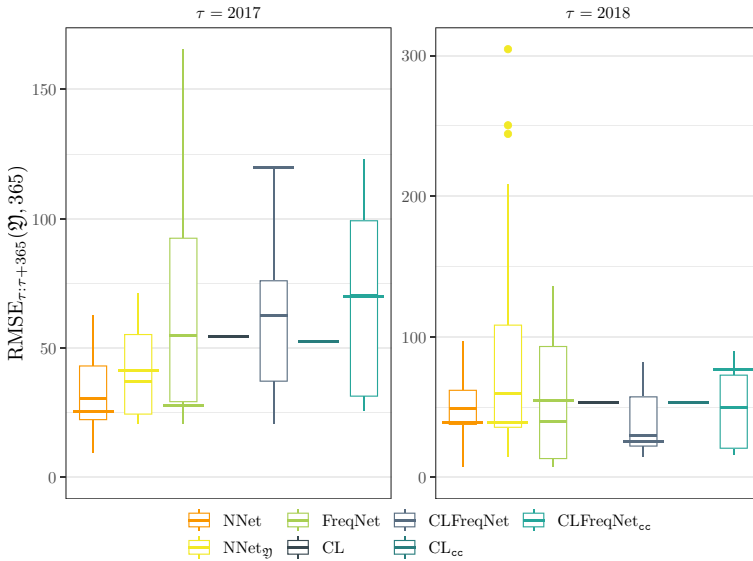
**Fig. 7** Boxplot comparison of $\mathrm{RMSE}_{\tau:\tau+365}(\mathfrak{Y}, 365)$ for the different methods. Final selected models shown as wide horizontal line

exhibits behavior more consistent with the simulation study, decreasing the error while maintaining a similar variance. While for $\tau = $ 31st December 2017, the selected model has a very low error compared to all candidates, the model selected for $\tau = $ 31st December 2018 exhibits a worse accident year level error than the chain ladder methods, even though the median error among all candidate models was lower than that of chain ladder. Regarding **CLFreqNet** and **CLFreqNet**$_{\mathrm{CC}}$, one can see the impact of the training procedure (holding out 25% of the data for validation) increasing the overall variance in error compared to their Chain Ladder counterparts.

In summary, the new methods seem to provide similar accuracy on an accident year level when compared to the underlying methods **NNet**, **CL** or **CL**$_{\mathrm{CC}}$.

The new methods **FreqNet**, **CLFreqNet** and **CLFreqNet**$_{\mathrm{CC}}$ provide exposure-level IBNR predictions, which can be compared in Fig. 8. As with the simulation study, plain triangle based methods can not be used to obtain viable predictors for micro-level claim counts, so the trivial **no IBNR** is used as a basic reference. Unfortunately it is not possible to also provide a theoretical best prediction on real data, so there is no **cheating** benchmark with which the results could be compared. We refer to Fig. 5 for the corresponding simulation study results which do have this benchmark. The micro-level results are more comparable across the different truncation times, showing a similar pattern to that of the simulation study with one exception: There is a larger separation between the errors of **FreqNet** and those of **CLFreqNet** and **CLFreqNet**$_{\mathrm{CC}}$.
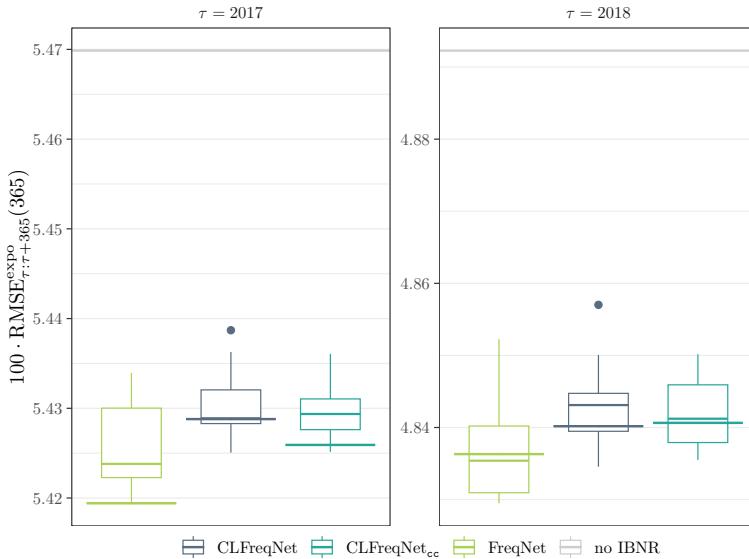
**Fig. 8** Comparison of $\mathrm{RMSE}^{\mathrm{expo}}_{\tau:\tau+365}(\mathfrak{Y}, 365)$ for the different methods

## 9 Conclusion

Two new methods for joint prediction of micro-level IBNR claim counts and claim frequencies have been developed and applied to real and simulated data. Results show promising accuracy on an exposure level compared to the theoretical optimum under laboratory conditions. The new methods also permit assigning IBNR claim count predictions on a policy level such that policies without claims can receive a non-zero IBNR prediction, which is an advantage for analysis of small sub-portfolios where applying chain ladder estimation factors—even with parameters estimated on a larger dataset - produce highly volatile estimates. The presented case studies uncover several opportunities for further research:

1. The distributional assumption of a Blended Dirac-Erlang-Generalized Pareto family for reporting delays in Model 3.1 might not be suitable for all applications. Future work could examine results with other reporting delay distribution families.
2. The functional relationships in Model 3.1, Model 3.2 and Model 3.5 have all been chosen as MLPs. All of these relationships could be chosen from a different function family, e.g., other families used in machine learning such as regression trees.
3. The architecture of all MLPs was simply chosen and no hyperparameter-optimization was performed. Strategies for architecture selection or different architectures could be examined.
4. Different strategies for model selection of a final model among candidate models could be explored.
5. Definition 2.1 can be extended by a claim settlement process for each claim, such as the one presented in [3] but on a policy level, to allow joint modelling of IBNR and RBNS payments.

**Availability of data, materials and code** The scripts used for simulation studies are provided as supplementary material. It requires an R package hosted on CRAN [20]. The real dataset used in Sect. 8 is proprietary.

## Declarations

**Conflict of interest** No conflicts of interest exist.

## Appendix A: Proof of Lemma 3.3

For $x \in \mathfrak{X}, t \geq 0$ and $y = (y_1, \ldots, y_Q) \in \mathfrak{Y}$, write $z^{(q)} = (x, t, y_1, \ldots, y_q)$ and $Z^{(q)} = (X, T, Y_1, \ldots, Y_q)$. Then, in view of the fact that the conditional density of $Y_q$ given $t + d \leq \tau$ and $Z^{(q-1)} = z^{(q-1)}$ may be written as

$$f_{Y_q|t+d\leq\tau, Z^{(q-1)}=z^{(q-1)}}(y_q) = \frac{P(t + d \leq \tau | Z^{(q)} = z^{(q)})}{P(T + D \leq \tau | Z^{(q-1)} = z^{(q-1)})} f_{Y_q|Z^{(q-1)}=z^{(q-1)}}(y_q),$$

we may rewrite each summand in (6) as

$$\mathbb{E}[\tilde{\ell}_{(X,T,Y_1,\ldots,Y_{q-1})}(g|Y_q) \mid T + D \leq \tau, Z^{(q-1)} = z^{(q-1)}]$$

$$= \int \frac{\log f_{g(z^{(q-1)})}(y_q)}{P(T + D \leq \tau | Z^{(q)} = z^{(q)})} f_{Y_q|t+d\leq\tau, Z^{(q-1)}=z^{(q-1)}}(y_q) \, d\mu^{(q)}(y_q)$$

$$= \frac{1}{P(D + T \leq \tau | Z^{(q-1)} = z^{(q-1)})}$$

$$\times \int \log f_{g(z^{(q-1)})}(y_q) f_{Y_q|Z^{(q-1)}=z^{(q-1)}}(y_q) \, d\mu^{(q)}(y_q)$$

$$= \frac{1}{P(D + T \leq \tau | Z^{(q-1)} = z^{(q-1)})} \mathbb{E}[\log f_{g(Z^{(q-1)})}(Y_q) \mid Z^{(q-1)} = z^{(q-1)}].$$

Note that the factor in front of the expectation does not depend on $g$.

Write

$$M(g) = \mathbb{E}\big[ \log f_{g(Z^{(q-1)})}(Y_q) - \log f_{g_0(Z^{(q-1)})}(Y_q) \mid Z^{(q-1)} = z^{(q-1)} \big],$$

and note that $M(g_0) = 0$. Moreover, since $\log(x) \leq 2(\sqrt{x} - 1)$ for $x \geq 0$, we have, for all $g \in \mathcal{G}^{(q)}$,

$$
\begin{aligned}
M(g) &\leq 2\, \mathbb{E}\Big[ \sqrt{f_{g(Z^{(q-1)})}(Y_q)/f_{g_0(Z^{(q-1)})}(Y_q)} - 1 \mid Z^{(q-1)} = z^{(q-1)} \Big] \\
&= 2 \int_{\mathfrak{Y}_q} \Big( \sqrt{f_{g(z^{(q-1)})}(y_q)/f_{g_0(z^{(q-1)})}(y_q)} - 1 \Big) f_{g_0(z^{(q-1)})}(y_q)\, \mathrm{d}\mu^{(q)}(y_q) \\
&= 2 \int_{\mathfrak{Y}_q} \sqrt{f_{g(z^{(q-1)})}(y_q) f_{g_0(z^{(q-1)})}(y_q)}\, \mathrm{d}\mu^{(q)}(y_q) - 2 \\
&= - \int_{\mathfrak{Y}_q} \Big( \sqrt{f_{g(z^{(q-1)})}(y_q)} - \sqrt{f_{g_0(z^{(q-1)})}(y_q)} \Big)^2 \mathrm{d}\mu^{(q)} \leq 0.
\end{aligned}
$$

Hence, $M(g) \leq 0 = M(g_0)$ for all $g \in \mathcal{G}^{(q)}$, which implies the assertion. □

# References

1. Abadi M, Agarwal A, Barham P, Brevdo E, Chen Z, Citro C, Corrado GS, Davis A, Dean J, Devin M, Ghemawat S, Goodfellow I, Harp A, Irving G, Isard M, Jia Y, Jozefowicz R, Kaiser L, Kudlur M, Levenberg J, Mané D, Monga R, Moore S, Murray D, Olah C, Schuster M, Shlens J, Steiner B, Sutskever I, Talwar K, Tucker P, Vanhoucke V, Vasudevan V, Viégas F, Vinyals O, Warden P, Wattenberg M, Wicke M, Yu Y, Zheng X (2015) TensorFlow: large-scale machine learning on heterogeneous systems. https://www.tensorflow.org/. Software available from tensorflow.org
2. Allaire J, Eddelbuettel D, Golding N, Tang Y (2016) tensorflow: R interface to tensorflow. https://github.com/rstudio/tensorflow
3. Antonio K, Plat R (2014) Micro-level stochastic loss reserving for general insurance. Scand Actuar J 2014(7):649–669
4. Antonio K, Godecharle E, Van Oirbeek R (2016) A multi-state approach and flexible payment distributions for micro-level reserving in general insurance. SSRN 20:20
5. Badescu AL, Lin XS, Tang D (2016) A marked Cox model for the number of IBNR claims: theory. Insur Math Econ 69:29–37
6. Baudry M, Robert CY (2019) A machine learning approach for individual claims reserving in insurance. Appl Stoch Model Bus Ind 2019(35):1127–1155
7. Bücher A, Rosenstock A (2023) Micro-level prediction of outstanding claim counts based on novel mixture models and neural networks. Eur Actuar J 13(1):55–90. https://doi.org/10.1007/s13385-022-00314-4
8. Bücher A, Rosenstock A (2023) Supplementary material: micro-level prediction of outstanding claim counts based on novel mixture models and neural networks. https://link.springer.com/article/10.1007/s13385-022-00314-4#Sec27
9. Chaoubi I, Besse C, Cossette H, Côté M-P (2023) Micro-level reserving for general insurance claims using a long short-term memory network. Appl Stoch Model Bus Ind 1:1–26. https://doi.org/10.1002/asmb.2750
10. Chollet F, Allaire J, Kalinowski T, Falbel D, Tang Y, Bijl WVD, Studer M, Keydana S (2017) R interface to keras. https://github.com/rstudio/keras
11. De Felice M, Moriconi F (2019) Claim watching and individual claims reserving using classification and regression trees. Risks 7:4
12. Gabrielli A, Wüthrich MV (2018) An individual claims history simulation machine. Risks 6(2):29

13. Goldburd M, Khare A, Tevet D, Guller D (2016) Generalized linear models for insurance rating. Casual Actuarial Soc CAS Monogr Ser 5:5
14. Goodfellow IJ, Bengio Y, Courville A (2016) *Deep Learning*. MIT Press, Cambridge. http://www.deeplearningbook.org
15. Last G, Penrose M (2018) Lectures on the Poisson process. Cambridge University Press, Cambridge
16. Mikosch T (2009) Non-life insurance mathematics. Universitext, 2nd edn. Springer, Berlin
17. Norberg R (1993) Prediction of outstanding liabilities in non-life insurance. ASTIN Bull 23(1):95–115
18. Norberg R (1999) Prediction of outstanding liabilities ii model variations and extensions. ASTIN Bull 29(1):5–25
19. Okine AN-A, Frees EW, Shi P (2022) Joint model prediction and application to individual-level loss reserving. ASTIN Bull J IAA 52(1):91–116. https://doi.org/10.1017/asb.2021.28
20. Rosenstock A (2023) Reservr: fit distributions and neural networks to censored and truncated data. https://ashesitr.github.io/reservr/. R package version 0.0.1
21. Wang Z, Wu X, Qiu C (2021) The impacts of individual information on loss reserving. ASTIN Bull J IAA 51(1):303–347. https://doi.org/10.1017/asb.2020.42
22. Wüthrich MV (2020) Bias regularization in neural network models for general insurance pricing. Eur Actuar J 10:179–202. https://doi.org/10.1007/s13385-019-00215-z
23. Wüthrich MV, Buser C (2019) Data analytics for non-life insurance pricing. SSRN 20:25
24. Wüthrich MV, Merz M (2015) Stochastic claims reserving manual: advances in dynamic modeling. SSRN 20:20