

# Assignment 1

Jakob og Oskar

2025

## Assignment 1 in interpretable machine learning

We are tasked with building a predictive regression model, with the best possible prediction. This submission will be split up into several parts:

- Initial data preprocessing and overview
- Introduction into the mathematics
- Modelling and justification
- Comparative discussion

### Initial data preprocessing and overview

The given data has the form:

```
## [1] "ID" "Date_start_contract" "Date_last_renewal"
## [4] "Date_next_renewal" "Date_birth" "Date_driving_licence"
## [7] "Distribution_channel" "Seniority" "Policies_in_force"
## [10] "Max_policies" "Max_products" "Lapse"
## [13] "Date_lapse" "Payment" "Premium"
## [16] "Cost_claims_year" "N_claims_year" "N_claims_history"
## [19] "R_Claims_history" "Type_risk" "Area"
## [22] "Second_driver" "Year_matriculation" "Power"
## [25] "Cylinder_capacity" "Value_vehicle" "N_doors"
## [28] "Type_fuel" "Length" "Weight"
```

We are asked to predict the **Cost\_claims\_year** given the rest of the covariate-vector. Initially it is important to note that our data is a classical insurance dataset, where we are given rows corresponding to insurance periods for a given contract. There are several issues with this, since some contracts might overlap into multiple contracts, which can be identified by the ID. It is however very difficult to locate these policies, and we overlook this issue.

There are numerous character vectors in the data, which can be seen here:

```
## Classes 'data.table' and 'data.frame': 105555 obs. of 30 variables:
## $ ID : int 1 1 1 1 2 2 3 3 3 3 ...
## $ Date_start_contract : chr "05/11/2015" "05/11/2015" "05/11/2015" "05/11/2015" ...
## $ Date_last_renewal : chr "05/11/2015" "05/11/2016" "05/11/2017" "05/11/2018" ...
## $ Date_next_renewal : chr "05/11/2016" "05/11/2017" "05/11/2018" "05/11/2019" ...
## $ Date_birth : chr "15/04/1956" "15/04/1956" "15/04/1956" "15/04/1956" ...
## $ Date_driving_licence: chr "20/03/1976" "20/03/1976" "20/03/1976" "20/03/1976" ...
```

```
## $ Distribution_channel: chr "0" "0" "0" "0" ...
## $ Seniority          : int  4 4 4 4 4 4 15 15 15 15 ...
## $ Policies_in_force  : int  1 1 2 2 2 2 1 1 1 1 ...
## $ Max_policies       : int  2 2 2 2 2 2 2 2 2 2 ...
## $ Max_products       : int  1 1 1 1 1 1 1 1 1 1 ...
## $ Lapse              : int  0 0 0 0 0 0 0 0 0 0 ...
## $ Date_lapse         : chr  "" "" "" "" ...
## $ Payment            : int  0 0 0 0 1 1 0 0 0 0 ...
## $ Premium            : num 223 214 215 217 214 ...
## $ Cost_claims_year   : num  0 0 0 0 0 0 0 0 0 0 ...
## $ N_claims_year      : int  0 0 0 0 0 0 0 0 0 0 ...
## $ N_claims_history   : int  0 0 0 0 0 0 0 0 0 0 ...
## $ R_Claims_history   : num  0 0 0 0 0 0 0 0 0 0 ...
## $ Type_risk          : int  1 1 1 1 1 1 3 3 3 3 ...
## $ Area               : int  0 0 0 0 0 0 0 0 0 0 ...
## $ Second_driver      : int  0 0 0 0 0 0 0 0 0 0 ...
## $ Year_matriculation : int 2004 2004 2004 2004 2004 2004 2013 2013 2013 2013 ...
## $ Power              : int  80 80 80 80 80 80 85 85 85 85 ...
## $ Cylinder_capacity  : int  599 599 599 599 599 599 1229 1229 1229 1229 ...
## $ Value_vehicle      : num 7068 7068 7068 7068 7068 ...
## $ N_doors            : int  0 0 0 0 0 0 5 5 5 5 ...
## $ Type_fuel          : chr  "P" "P" "P" "P" ...
## $ Length             : num  NA NA NA NA NA ...
## $ Weight             : int  190 190 190 190 190 190 1105 1105 1105 1105 ...
## - attr(*, ".internal.selfref")=<externalptr>
## NULL
```

One could, model some of the char. vectors like proposed in the lectures, by

$$X_i = E(Y \mid X_i) \\ = \int_Y y \mu(x, dy)$$

where  $\mu$  is the appropriate probability kernel, and we let  $y$  be our response. However, we choose to take the our char. vectors and round them to yearly values, which then has a ordinal ordering and can thus be used as features. Further we take and one-hot encode *Distribution\_channel* by creating three new features which are either 1 or 0. Same for the type of fuel.

Next there are some missing values.

It becomes apparent that there are missing values in the length and type\_fuel variable. We can see that the missingness is overlapping in 1764 rows. However the length variable suffers way heavier from missingness compared to type\_fuel. In order to impute values, we assume the missingness is completely at random for both features. We consider correlated features to do imputation. We see from the correlation plot (fig. 3) that type\_fuel is mostly correlated with Cylinder\_capacity, Value\_vehicle and Weight. The same goes for the feature Length. Therefore, we fit a multivariate linear regression model with these 3 covariates to predict both type\_fuel and Length (using cbind in the response formula). Finally, we predict the missing values using this trained model.

Our response variable **Cost\_claims\_year** suffers from a few extreme values. We decide to remove these values to later on achieve a better model fit. In fig. 2, **N\_claims\_year** is plotted against **Cost\_claims\_year**, where at least 5-10 extreme values of **Cost\_claims\_year** are spotted.

Finally we look at the correlation between the covariates.

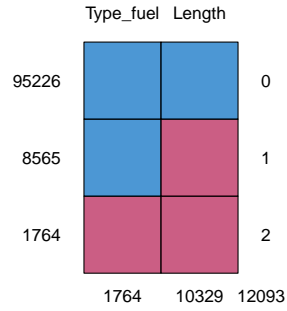


Figure 1: Missing data plot, right axis shows numer of missing columns in that row, and the left axis show how many rows have this missingness pattern

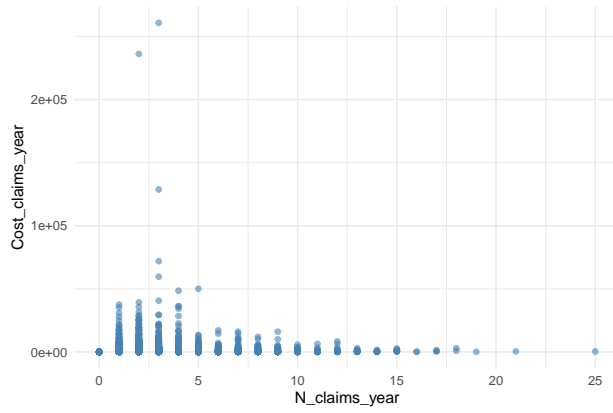


Figure 2: Claim costs over number of claims

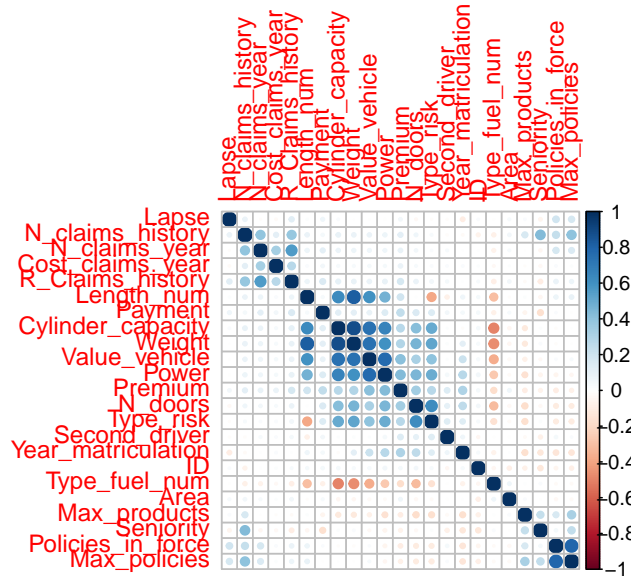


Figure 3: Correlation plot between the continious covariates

We notice some clusters, most meaningfull between *Cylinder\_capacity*, *Value\_vehicle* and *Weight* which is expected. We deem these to have significant predictive ability, and thus we choose to not remove these. The top left cluster, will be ignored for now, since we will later introduce a data-transformation which affects this cluster in a high degree.

**Data Aggregation** For our data aggregation we want to have the ability to assume independence, which we have if we aggregate the rows into being one policy. There are several problems to tackle to achieve data in the desired format. We have to both choose an appropriate aggregations function, and we have to consider events where the insured object change. We start off by identifying some unique-identifiers. If theme columns change, we deem the contract to insure a new car, or at least insure something that is independent of what was previously insured. Amongst these unique identifiers are Length, Weight, Power etc. so things which change, probably mean that the thing that is insured has changes. \ Next, we look at how we have aggregated data. Firstly we sum the exposure for each contract, since we will use this to weight in the models later on. Next we use the max function on the number of policies, lapse, date of birth, date of driving licence, start of contract, products, payments, ratio of claims history, type of risk,. For The mean premium we use the mean of the premium without the forward premium, that is the premium for the time we want to predict for. For the value of the veichle, and number of doors we use mean. For the length we use sum since this is just a scaling. For the variable *cost\_claims\_this\_year* we use the cost of the claims in  $\mathcal{F}_t$  where we are to predict data on  $\mathcal{F}_{t+1}$ .

So we have  $Z_i(\omega_i) := (X_i(\omega_i), Y_i(\omega_i)) : (\Omega_i, \mathcal{F}_i) \rightarrow (\mathcal{X} \times \mathcal{Y}, \mathcal{B}(\mathcal{X} \times \mathcal{Y}))$ , where by aggregating data row-wise we obtain independence on every set in the Borel-set induced. However this is an empirical assumption which can be violated by catastrophe events.

## Introduction into the mathematical framework

Note the classical derivation of how to decompose a total claim sum.

$$\begin{aligned} E(S(t) \mid Z = z) &= E \left( \sum_{i=1}^{N(t)} X_i \mid Z = z \right) \\ &= E \left( \left( \sum_{i=1}^{N(t)} X_i \mid Z = z \cap \mathcal{F}(t) \right) \right) \\ &= E(N E(X_i \mid Z = z) \mid Z = z) \\ &= E(N \mid Z = z) E(X_1 \mid Z = z) \end{aligned}$$

We would however like to modify this slightly, later on. Further the general definition of the Tweedie law is

$$P_{\theta, \sigma^2}(Y \in A) = \int_A \exp \left\{ \frac{\theta z - \kappa_p(\theta)}{\sigma^2} \right\} \nu_y(dz)$$

Where  $A \subset \mathbb{R}$  and measurable.  $\theta$  is the canonical parameter, and  $\kappa_p(\theta)$  is the cumulant function dependent on  $p$  the so called tweedie power parameter.  $\sigma^2$  is the dispersion parameters, and  $\nu_\lambda$  is some sigma-finite base measure depending on the Lebesgue measure.

since it can accommodate a zero-inflated distribution very well, we feel the need to introduce this model here. Lastly we introduce the so-called *bbrier* measure, for classification models.

$$\frac{1}{n} \sum_{i=1}^n w_i (1\{X_i = 1\} - P(X_i = 1))^2$$

Which is the mean weighted square euclidean distance between the probability and the true value.

## Parameter search space

Model	tweedie_power	learning_rate_eta	leaves	iterations	s	alpha
Xgboost	1-1.9	1e-3 - 0.2	16 - 32	200 - 1000		
lgbm	1-1.9	1e-3 - 0.2	3 - 9	200 - 1000		
glmnet					1e-4 - 1	0 - 1

## Modelling and justification

In this section, we train and evaluate different models. For all the proposed models, we have carried out nested crossvalidation. We do stratified resampling in both the inner and outer cv loops, to ensure both zero and non-zero claims over the training and test set. I final note, that goes for all our trained models, is that we introduce offset on the exposure. Meaning: Give more importance to observations with more exposure since they represent more information.

```
#Common trans dataset to skip rerunning trans
```

```
data_trans <- data_trans(data)
```

**Tweedie** Initily we want to fit a Tweedie model on the entire data, since the Tweedie model is known for handling zero-inflated distributions well. We plot the distribution of the claim size.



Figure 4: Distribution of claim size pr. unique contract

We apply both the Xgboost, Ligth-gbm and glmnet. For each of these we supply search-space information in the **Parameter search space** Table.

We define the workflow as shown in 5.

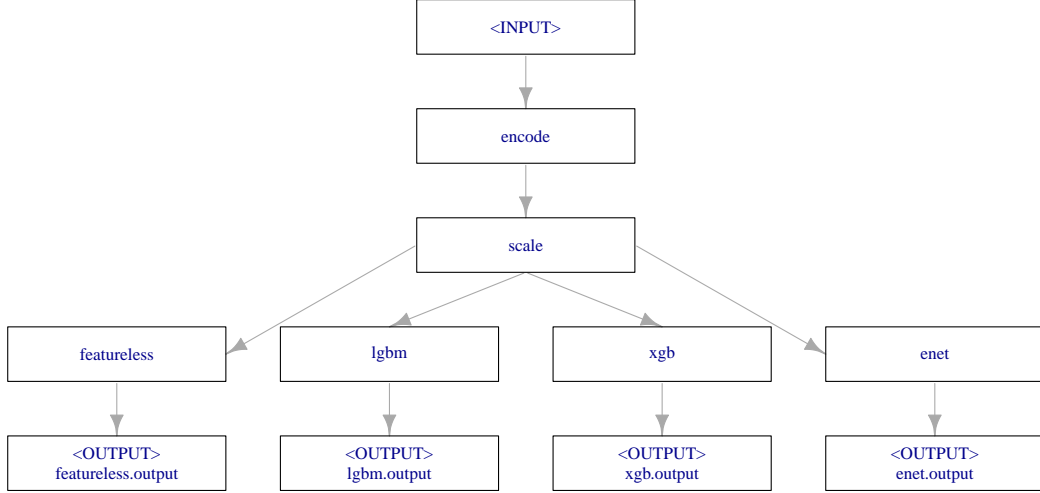


Figure 5: Input and model flow

The results of the nested cross-validation are shown in 6.

And finally the chosen hyperparameters are presented in Table **Best parameters tweedie**.

This tweedie power parameter makes the model a compound gamma.

**Split model** Next we can look at the split model proposed in the mathematical frame, where we aim to create two models, and then predict in the following manner:

$$P(N > 0 \mid Z) \cdot E(X \mid Z)$$

Which is unlike the traditional and proposed  $E(N)E(X)$ , but since we are asked to predict the next year, we feel like we can provide a suitable approximation with the above prediction. This entails fitting both a frequency model and a severity model. Hence we begin with the frequency regression model.

We have fitted the model with nested cross-validation since it provides a more stable estimate of the generalisation error. We use  $K = 2$  folds in both the inner and outer loops. Next we look at some model diagnostics. We can look at the combined *bbrier* score for the aggregated model:

which seems quite small. Further we have some confusion matrices for each of the folds.

Which in general show that the accuracy and precision is very high, meaning that we are usually correct when predicting a claim. We see that the True positive ratio is lower than the  $PPV = \frac{TP}{TP+FP}$ , indicating that we miss a fair share of claims. Overall we note that our model seems stable through the cross-validation and that it performs somewhat well.

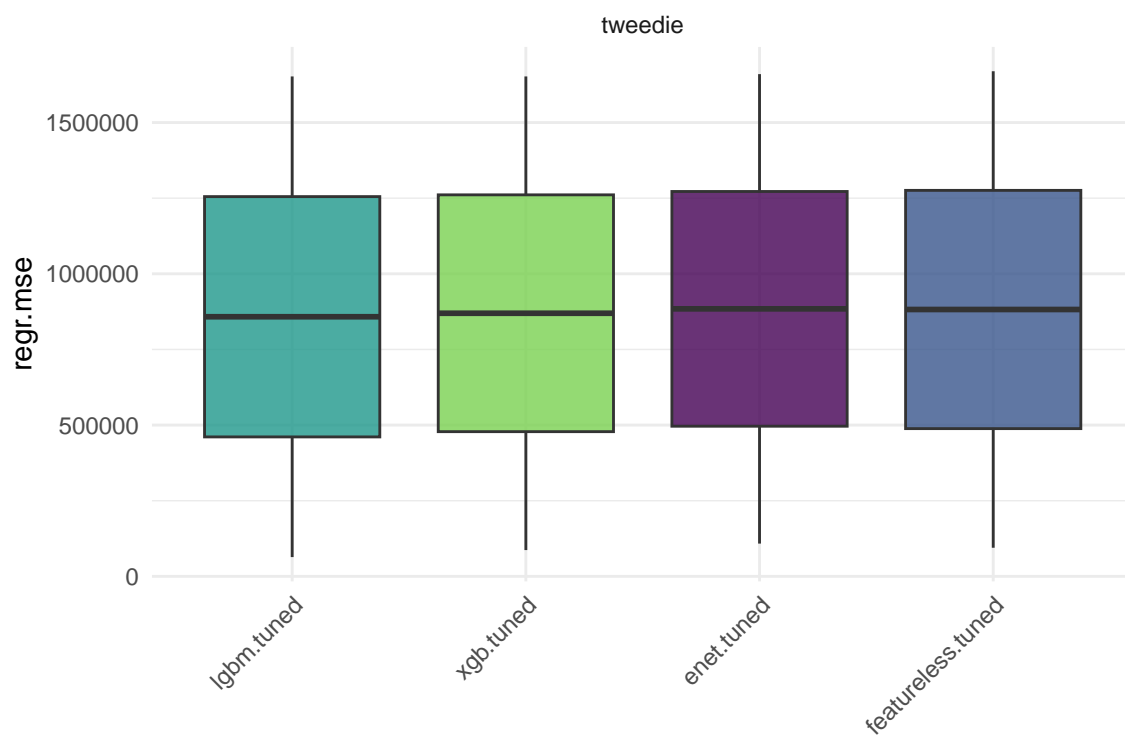


Figure 6: MSE boxplot

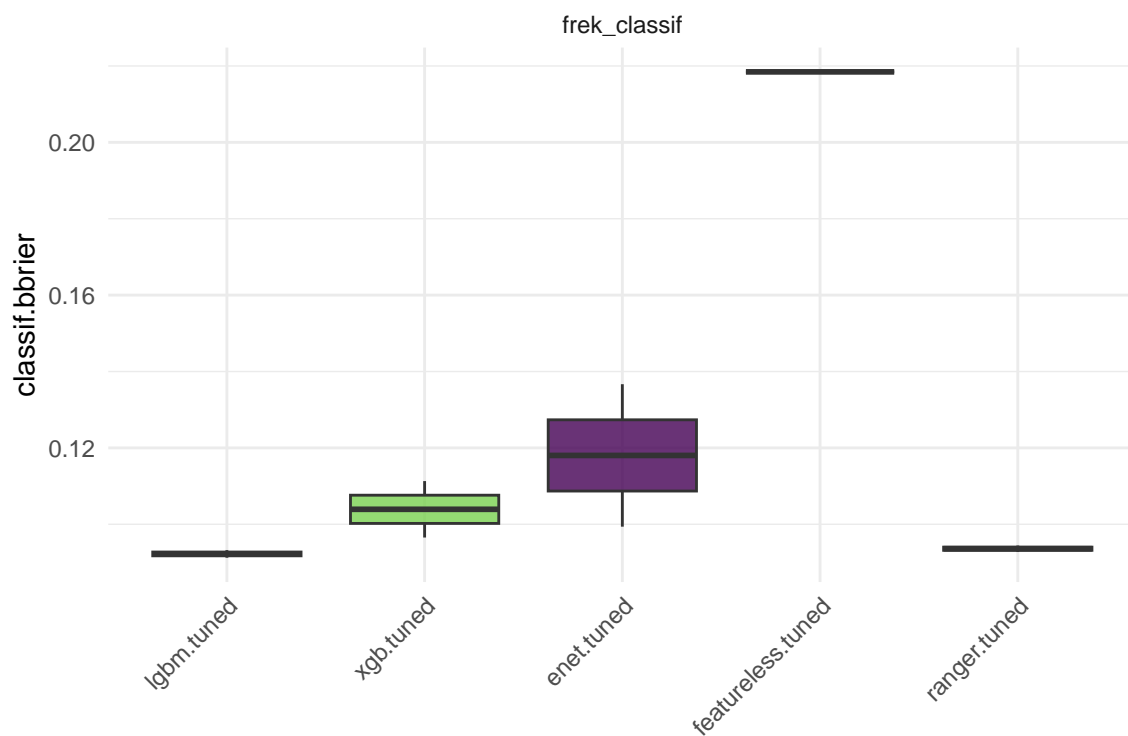


Figure 7: MSE boxplot

## Best parameters tweedie

parameter	value
encode.method	one-hot
scale.robust	FALSE
regr.lightgbm.objective	tweedie
regr.lightgbm.verbose	-1
regr.lightgbm.learning_rate	0.03195343
regr.lightgbm.num_leaves	23
regr.lightgbm.num_threads	1
regr.lightgbm.tweedie_variance_power	1.870088
regr.lightgbm.num_iterations	989

Here we see the *acc* error based on the choosen probability threshold, which as expected shows that for example a threshold at 0.5 results in a pretty good model. By setting the threshold higher we could risk introducing more false negatives.

Next we look at the severity model.

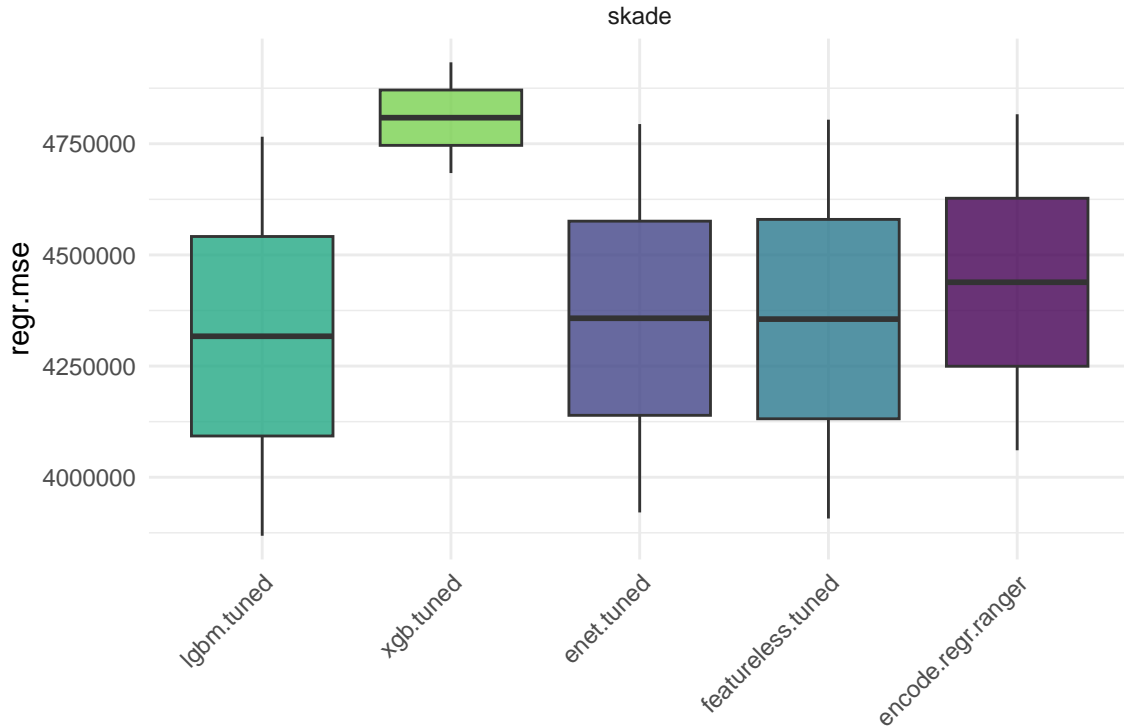


Figure 8: Shows the predicted values along the x-axis vs the real values on the y-axis

Then we can combine the models, and calculate the mean square error for the combined model. We have shown that our models are somewhat acceptable, and we not fit them on the entire dataset and find the out of sample *MSE* for 20 percent of data.

We see here that our final model for the combined model. It is clear that both our models perform extremely poorly, and one should probably have split up the claims into small and large claims so we could model the



large and small claims independently. Since both models are somewhat lacking we stick with the split model, since this is more common in the actuarial business.

## Discussion

Here we touch upon the various modeling problems. ##### Pre-processing We took a somewhat minimalist approach: obvious date parsing, encoding etc. some slight imputation. Since we impute, we impute off random values which adds noise. We could have used better imputation, and our data aggregation could have been slightly better.

**Model choices** In the frequency model, despite the good accuracy, we miss a lot of false negatives. In the severity model, we severely miss the tail of the distribution.

**Evaluation protocol** We did nested cross-validation, but we could have increased the iterations and the number of folds for more stable results. For severity we could have chosen a better metric to optimize for, which will hopefully result in a better model, especially in the tail. We could have fit an EVT distribution on the heavy tailed data, but again this is quite strenuous for this assignment.