Motivation
0000000

Methods
00
00000
0000000000000

Data
000

Results
00000

Summary
0

Anlage
00

# XGBoost in Non-Life Insurance Pricing

Berlin, 17.01.2025

**Hans Alcahya Tjong**
570795

# Inhaltsverzeichnis

# Motivation

## Rating factors

Insurance premium depends on *rating factors*.

The rating factors can be distinguished in three characteristics:

1. Characteristics of the policyholder:
   Gender, age, marital status
2. Characteristics of the insured objects:
   Car manufacturer, engine size, model year
3. Characteristics of the geographical region:
   City, postal code

Motivation
○○○●○○○

Methods
○○
○○○○○
○○○○○○○○○○○○

Data
○○○

Results
○○○○○

Summary
○

Anlage
○○

# Key ratios

**Duration** $n_o$
  Policy years

**Claim number** $N$

**Claim frequency** $SH$
  Number of claims divided by the duration

# Key ratios

**Total claim amount** $S$

**Claim severity** $SD$
   Total claim amount divided by the number of claims (average cost per claim)

**Pure premium** $SB$
   Total claim amount divided by the duration (average cost per policy year)

# Key ratios

| Exposure $[\omega]$ | Target value $[X]$ | Normalized target value $\left[Y = \frac{X}{\omega}\right]$ |
|:---:|:---:|:---:|
| $n_o$ | $N$ | $SH = \frac{N}{n_o}$ |
| $N$ | $S$ | $SD = \frac{S}{N}$ |
| $n_o$ | $S$ | $SB = \frac{S}{n_o}$ |

Tabelle: Connection between key ratio.

Motivation
○○○○○●○

Methods
○○
○○○○○
○○○○○○○○○○○○○

Data
○○○

Results
○○○○○

Summary
○

Anlage
○○

## Multiplicative model

Let $M$ denote number of rating factors and $i_k$ denote the class of the rating factor $k$.

The multiplicative model is defined by:

$$\mu_{i_1, i_2, \ldots, i_M} = \gamma_0 \prod_{k=1}^{M} \gamma_{k, i_k}$$

where $\gamma_0 \in \mathbb{R}$ is basis premium and $\gamma_{k, i_k} \in \mathbb{R}$ are price relativities.

Motivation
○○○○○○○●

Methods
○○
○○○○○
○○○○○○○○○○○○○

Data
○○○

Results
○○○○○

Summary
○

Anlage
○○

## Example

| Age | Location | $i_1$ | $i_2$ | $\mu_{i_1,i_2}$ |
|-------|------------|------|------|---------------------------------|
| 21-40 | Urban | 1 | 1 | $\gamma_0\gamma_{1,1}\gamma_{2,1}$ |
| 21-40 | Small town | 1 | 2 | $\gamma_0\gamma_{1,1}\gamma_{2,2}$ |
| 41-60 | Urban | 2 | 1 | $\gamma_0\gamma_{1,2}\gamma_{2,1}$ |
| 41-60 | Small town | 2 | 2 | $\gamma_0\gamma_{1,2}\gamma_{2,2}$ |
| 61-80 | Urban | 3 | 1 | $\gamma_0\gamma_{1,3}\gamma_{2,1}$ |
| 61-80 | Small town | 3 | 2 | $\gamma_0\gamma_{1,3}\gamma_{2,2}$ |

Tabelle: Example of a multiplicative model.

htw
Hochschule für Technik
und Wirtschaft Berlin
University of Applied Sciences

Motivation
Methods
Data
Results
Summary
Anlage

# Methods

Motivation
0000000

Methods
○●
○○○○○
○○○○○○○○○○○○○

Data
○○○

Results
○○○○○

Summary
○

Anlage
○○

## Methods

1. Generalized Linear Models (GLM)
2. Extreme Gradient Boosting (XGBoost)

Motivation
0000000

Methods
00
●0000
000000000000

Data
000

Results
00000

Summary
0

Anlage
00

Generalized Linear Models (GLM)

# Components of GLM

1. Distribution
2. Systematic component
3. Link-Function

Motivation
0000000

Methods
00
00000
000000000000

Data
000

Results
00000

Summary
0

Anlage
00

Generalized Linear Models (GLM)

# Components of GLM
1. Distribution

Let the response variable $Y_i$ have a distribution from the exponential family:

$$f_{Y_i}(y_i; \theta_i, \phi) = exp\left(\frac{y_i\theta_i - a(\theta_i)}{\phi/\nu_i}\right) c(y_i, \phi, \nu_i)$$

Motivation
0000000

Methods
○○
○○●○○
○○○○○○○○○○○○○

Data
○○○

Results
○○○○○

Summary
○

Anlage
○○

Generalized Linear Models (GLM)

# Components of GLM
## 1. Distribution

where:

$\theta_i$ = canonical parameter (real value)

$\phi$ = Dispersion parameters (positive)

$\nu_i$ = Weight (positive)

$a(\cdot)$ = Cumulant function (bijective and 2 time continuously differentia

$c(\cdot)$ = Positive normalizing function

# Components of GLM
2. Systematic component

The explanatory variables result in a linear predictor:

$$\eta_i = X_i^T \beta = \beta_0 + \sum_{j=1}^{p} x_{ij}\beta_j$$

Motivation
0000000

Methods
00
0000●0000000000

Data
000

Results
00000

Summary
0

Anlage
00

Generalized Linear Models (GLM)

# Components of GLM
## 3. Link-Function

The bijective and two time continuously differentiable link function $g(\cdot)$ links the random component and the systematic component:

$$g(\mu_i) = \eta_i \Leftrightarrow \mu_i = g^{-1}(\eta_i)$$

where $\mu_i = \mathbb{E}(Y_i|X)$

Motivation
0000000

Methods
00
00000
●000000000000

Data
000

Results
00000

Summary
0

Anlage
00

Extreme Gradient Boosting (XGBoost)

# XGBoost: Example

| $i$ | Age | Location | $\hat{\mu}_i$ | $S_i$ | Error | Error$^2$ |
|-----|-----|----------|---------------|-------|-------|-----------|
| 1 | 25 | Urban | 1000 | 1300 | $+300$ | 90000 |
| 2 | 40 | Small town | 1000 | 500 | -500 | 250000 |
| 3 | 32 | Urban | 1000 | 1000 | 0 | 0 |

Motivation
0000000

Methods
00
00000
0●00000000000

Data
000

Results
00000

Summary
0

Anlage
00

Extreme Gradient Boosting (XGBoost)

# XGBoost: Example

- Tree 1: If Age $< 35$, then $+100$. Otherwise $-200$
- Tree 2: If Location $=$ 'Urban', then $+150$. Otherwise $-300$

Motivation
0000000

Methods
OO
OOOOO
OOO●OOOOOOOOO

Data
OOO

Results
OOOOO

Summary
O

Anlage
OO

Extreme Gradient Boosting (XGBoost)

# XGBoost: Example

| $i$ | Age | Location | $\hat{\mu}_i$ | $S_i$ | Error | Error$^2$ |
|---|---|---|---|---|---|---|
| 1 | 25 | Urban | 1250 | 1300 | $+50$ | 2500 |
| 2 | 40 | Small town | 500 | 500 | 0 | 0 |
| 3 | 32 | Urban | 1150 | 1000 | -150 | 22500 |

Motivation
0000000

Methods
00
00000
0000●00000000

Data
000

Results
00000

Summary
0

Anlage
00

Extreme Gradient Boosting (XGBoost)

# XGBoost

- Objective function

$$obj(\Theta) = L(\Theta) + \Omega(\Theta)$$

where:

| | |
|---|---|
| $\Theta$ | Parameters |
| $L(\Theta)$ | loss function |
| $\Omega(\Theta)$ | regularization term |

Motivation
OOOOOOO

Methods
OO
OOOOO
OOOO●OOOOOOOO

Data
OOO

Results
OOOOO

Summary
O

Anlage
OO

Extreme Gradient Boosting (XGBoost)

# XGBoost

- Loss function $L(\Theta)$ measures how predictive the model is with respect to the training data, e.g. *MSE*:

$$L(\Theta) = \sum_i (y_i - \hat{y}_i)^2$$

- Regularization term $\Omega(\Theta)$ measures the complexity of the model

Motivation
0000000

Methods
00
00000
00000●0000000

Data
000

Results
00000

Summary
0

Anlage
00

Extreme Gradient Boosting (XGBoost)

# XGBoost

- The model can be mathematically written as

$$\hat{y}_i^{(t)} = \sum_{k=1}^{t} f_k(x_i) = \hat{y}_i^{(t-1)} + f_t(x_i)$$

where:

$\hat{y}_i^{(t)}$    prediction value at step $t$

$f_t(x_i)$    function of $t$-th tree

# XGBoost: Example



tree1

age < 20

Y          N

+2

-1

tree2

Use Computer Daily

Y          N

+0.9

-0.9

f( ) = 2 + 0.9 = 2.9    f( ) = -1 - 0.9 = -1.9

Abbildung:
https://xgboost.readthedocs.io/en/stable/tutorials/model.html

Motivation
0000000

Methods
00
00000
000000000●00000

Data
000

Results
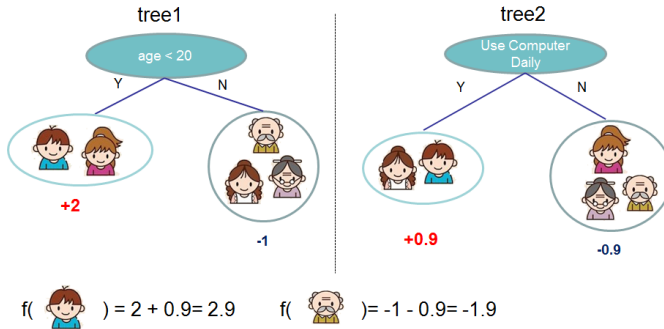00000

Summary
0

Anlage
00

Extreme Gradient Boosting (XGBoost)

# XGBoost

- Objective function of the $t$-th tree to be optimized using MSE as loss function

$$
\begin{aligned}
obj^{(t)} &= \sum_{i=1}^{n} (y_i - (\hat{y}_i^{(t-1)} + f_t(x_i)))^2 + \sum_{k=1}^{t} \omega(f_k) \\
&= \sum_{i=1}^{n} \left[ 2(\hat{y}_i^{(t-1)} - y_i) f_t(x_i) + f_t(x_i)^2 \right] + \omega(f_t) + \epsilon \\
&\overset{T}{=} \sum_{i=1}^{n} \left[ l(y_i, \hat{y}_i^{(t-1)}) + g_i f_t(x_i) + \frac{1}{2} h_i f_t(x_i)^2 \right] + \omega(f_t) + \epsilon
\end{aligned}
$$

Motivation
0000000

Methods
00
00000
000000000●0000

Data
000

Results
00000

Summary
0

Anlage
00

Extreme Gradient Boosting (XGBoost)

# XGBoost

where:

- $\omega(f_k)$   the complexity of the tree $f_k$

- $g_i = \delta_{\hat{y}_i^{(t-1)}} I(y_i, \hat{y}_i^{(t-1)})$

- $h_i = \delta^2_{\hat{y}_i^{(t-1)}} I(y_i, \hat{y}_i^{(t-1)})$

Motivation
0000000

Methods
OO
OOOOO
OOOOOOOOOO●OOO

Data
OOO

Results
OOOOO

Summary
O

Anlage
OO

Extreme Gradient Boosting (XGBoost)

# XGBoost

- The complexity of the tree $f$ is defined by

$$\omega(f) = \gamma T + \frac{1}{2} \lambda \sum_{j=1}^{T} \omega_j^2$$

where:

|   |   |
|---|---|
| $\gamma$ | minimum loss reduction parameter |
| $T$ | number of leaves |
| $\lambda$ | Ridge regularization parameter |
| $\omega_j$ | Score of leaf |

htw.
Hochschule für Technik
und Wirtschaft Berlin
University of Applied Sciences

Motivation
0000000

Methods
00
00000
0000000000000●00

Data
000

Results
00000

Summary
0

Anlage
00

Extreme Gradient Boosting (XGBoost)

# XGBoost

- After re-formulating

$$obj^{(t)} \approx \sum_{j=1}^{T} \left[ G_j \omega_j + \frac{1}{2}(H_j + \lambda)\omega_j^2 \right] + \gamma T$$

where:

$$G_j = \sum_{i \in I_j} g_i$$

$$H_j = \sum_{i \in I_j} h_i$$

$I_j$   set of indices of data points assigned to the j-th leaf

# XGBoost

- The best $\omega_j$ and objective function $obj^*$ are

$$\omega_j^* = -\frac{G_j}{H_j + \lambda}$$

$$obj^* = -\frac{1}{2} \sum_{j=1}^{T} \frac{G_j^2}{H_j + \lambda} + \gamma T$$

# XGBoost: Example

Instance index    gradient statistics

1    g1, h1

2    g2, h2

3    g3, h3

4    g4, h4

5    g5, h5

age < 15

Y          N

is male?

Y          N

$I_1 = \{1\}$          $I_2 = \{4\}$
$G_1 = g_1$          $G_2 = g_4$
$H_1 = h_1$          $H_4 = h_4$

$I_3 = \{2, 3, 5\}$
$G_3 = g_2 + g_3 + g_5$
$H_3 = h_2 + h_3 + h_5$

$$Obj = -\sum_j \frac{G_j^2}{H_j + \lambda} + 3\gamma$$

The smaller the score is, the better the structure is

Abbildung:
https://xgboost.readthedocs.io/en/stable/tutorials/model.html

Motivation
○○○○○○○

Methods
○○
○○○○○
○○○○○○○○○○○○○

Data
●○○

Results
○○○○○

Summary
○

Anlage
○○

# Data

Motivation
0000000

Methods
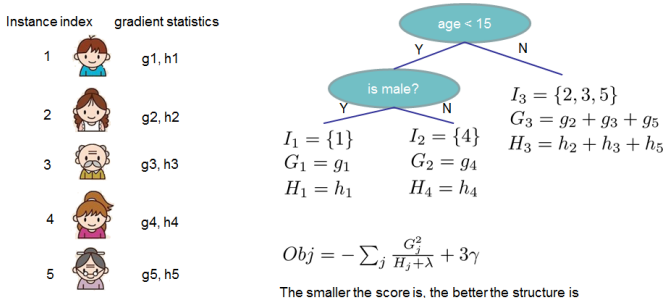00
00000
0000000000000

Data
0●0

Results
00000

Summary
0

Anlage
00

# Data Profile

- Available in R package `dataOhlsson`: `insuranceData`.

- Derived from period 1994-1998

- Partially comprehensive insurance for motorcycles.

- 64548 observations, each corresponding to one insurance policy.

Motivation
0000000

Methods
00
00000
0000000000000

Data
00●

Results
00000

Summary
0

Anlage
00

## Features

| | |
|---|---|
| **agarald** | Age of the policyholder (0-99) |
| **kon** | Gender (K/M) |
| **zon** | Zone (1-7) |
| **mcklass** | Vehicle type (1-7) |
| **fordald** | Vehicle age (0-99) |
| **bonuskl** | Bonus class (1-7) |
| **duration** | Policy year |
| **antskad** | Claim number |
| **skadkost** | Total claim amount |

htw.
Hochschule für Technik
und Wirtschaft Berlin
University of Applied Sciences

Motivation
○○○○○○○

Methods
○○
○○○○○
○○○○○○○○○○○○○○○

Data
○○○

Results
●○○○○

Summary
○

Anlage
○○

# Results

Motivation
0000000

Methods
00
00000
0000000000000

Data
000

Results
0●000

Summary
0

Anlage
00

## Claim Frequency

| Model | Features |
|-------|----------|
| GLM0 | (Intercept model) |
| GLM1 | Alter + Geschlecht + Zone + Fahrzeugtyp + Fahrzeugalter + Bonusklasse |
| GLM2 | Alter + Geschlecht + Zone + Fahrzeugtyp + Fahrzeugalter + Bonusklasse + Bonusklasse * Zone |
| XGB | Alter + Geschlecht + Zone + Fahrzeugtyp + Fahrzeugalter + Bonusklasse + Schadenaufwendungen |

Motivation
0000000

Methods
00
00000
0000000000000

Data
000

Results
00●00

Summary
0

Anlage
00

# Claim Frequency

| Model | Claim Frequency (Test) | | Poisson Deviance | | RMSE |
|-------|-------|-----------|-------|-------|--------|
|       | Actual | Predicted | Train | Test |        |
| GLM0  | 1.00%  | 1.09%     | 8.75% | 7.50% | 0.1066 |
| GLM1  | 1.00%  | 1.10%     | 8.49% | 6.83% | 0.1055 |
| GLM2  | 1.00%  | 1.10%     | 8.38% | 6.76% | 0.1055 |
| XGB   | 1.00%  | 1.01%     | 0.00% | 0.12% | 0.0253 |

Motivation
0000000

Methods
00
00000
0000000000000

Data
000

Results
00000

Summary
0

Anlage
00

## Claim Severity

| Model | Features |
|-------|----------|
| GLM0 | (Intercept model) |
| GLM1 | Alter + Geschlecht + Zone + Fahrzeugtyp + Fahrzeugalter + Bonusklasse |
| XGB | Alter + Geschlecht + Zone + Fahrzeugtyp + Fahrzeugalter + Bonusklasse + Versicherungsjahre |

Motivation
0000000

Methods
00
00000
0000000000000

Data
000

Results
0000●

Summary
0

Anlage
00

# Claim Severity

| Model | Claim Severity (Test) | | RMSE |
|-------|-----------|-----------|----------|
| | Actual | Predicted | |
| GLM0 | 4762611.00 | 5116337.08 | 38224.70 |
| GLM1 | 4762611.00 | 4861973.03 | 38418.00 |
| XGB | 4762611.00 | 4984814.42 | 42570.22 |

htw.
Hochschule für Technik
und Wirtschaft Berlin
University of Applied Sciences

# Summary

- For predicting claim frequency, XGBoost outperformed GLM
- XGBoost's RMSE in predicting claim severity higher than GLM's RMSE
- XGBoost captures interactions among rating variables, GLM needs manual creation of interaction terms to find a better model
- Parameter tuning in XGBoost could be required

Motivation
0000000

Methods
00
00000
0000000000000

Data
000

Results
00000

Summary
0

Anlage
●○

## Anlagen
Literaturverzeichnis

📄 Esbjörn Ohlsson und Björn Johansson. Non-Life Insurance Pricing with Generalized Linear Models. Springer Berlin, Heidelberg, 18. März 2010. ISBN: 978-3-642-10791-7.

📄 J. A. Nelder und R. W. M. Wedderburn. „Generalized Linear Models". In: Journal of the Royal Statistical Society. Series A (General) 135.3 (1972), S. 370–384. ISSN: 00359238. URL: http://www.jstor.org/stable/2344614.

📄 Michel Denuit, Donatien Hainaut und Julien Trufin. Effective Statistical Learning Methods for Actuaries I. Springer Cham, 3. Sep. 2019. ISBN: 978-3-030-25820-7.

Motivation
0000000

Methods
OO
OOOOO
OOOOOOOOOOOOO

Data
OOO

Results
OOOOO

Summary
O

Anlage
O●

## Anlagen
Literaturverzeichnis

📄 Trevor Hastie und Robert Tibshirani. „Generalized Additive Models". In: Statistical Science 1.3 (1986), S. 297–310. ISSN: 08834237. URL: http://www.jstor.org/stable/2245459.

📄 Mario Wuthrich und Christoph Buser. „Data Analytics for Non-Life Insurance Pricing". In: SSRN Electronic Journal (Jan. 2017). DOI: 10.2139/ssrn.2870308

📄 M.V. Wüthrich und M. Merz. Statistical Foundations of Actuarial Learning and its Applications. Springer Actuarial. Springer International Publishing, 2022. ISBN: 978-3-031-12411-2

htw.
Hochschule für Technik
und Wirtschaft Berlin
University of Applied Sciences