

Dynamically updating motor insurance prices with telematics collected driving behavior data

Katrien Antonio

LRisk - KU Leuven and ASE - University of Amsterdam

May 12, 2022





Katrien Antonio



Roel Henckaerts

Paper will appear in **Insurance: Mathematics and Economics**.

Recent work on insurance pricing analytics

[Henckaerts et al., 2018]

Henckaerts, O., Koenig, L., & Schmitt, M. (2018). *Data driven biasing strategy for the construction of insurance tariff classes*. *Journal of Risk and Insurance*, 87(1), 1-18.

Abstract
This paper presents a data-driven strategy for the construction of insurance tariff classes. The strategy is based on the use of machine learning techniques to identify the most relevant features for the construction of tariff classes. The strategy is implemented using a decision tree algorithm. The results of the study show that the proposed strategy outperforms traditional methods in terms of predictive accuracy and stability.

Introduction
The insurance industry is facing a number of challenges, including the need to improve its pricing accuracy and to reduce its operational costs. One of the key challenges is the construction of insurance tariff classes. This paper presents a data-driven strategy for the construction of insurance tariff classes. The strategy is based on the use of machine learning techniques to identify the most relevant features for the construction of tariff classes. The strategy is implemented using a decision tree algorithm. The results of the study show that the proposed strategy outperforms traditional methods in terms of predictive accuracy and stability.

SAJ

[Henckaerts et al., 2021]

Henckaerts, O., Koenig, L., & Schmitt, M. (2021). *Boosting insights in insurance tariff class with tree-based machine learning methods*. *Journal of Risk and Insurance*, 90(1), 1-18.

Abstract
This paper presents a data-driven strategy for the construction of insurance tariff classes. The strategy is based on the use of machine learning techniques to identify the most relevant features for the construction of tariff classes. The strategy is implemented using a decision tree algorithm. The results of the study show that the proposed strategy outperforms traditional methods in terms of predictive accuracy and stability.

Introduction
The insurance industry is facing a number of challenges, including the need to improve its pricing accuracy and to reduce its operational costs. One of the key challenges is the construction of insurance tariff classes. This paper presents a data-driven strategy for the construction of insurance tariff classes. The strategy is based on the use of machine learning techniques to identify the most relevant features for the construction of tariff classes. The strategy is implemented using a decision tree algorithm. The results of the study show that the proposed strategy outperforms traditional methods in terms of predictive accuracy and stability.

NAAJ

[Henckaerts et al., 2022]

Henckaerts, O., Koenig, L., & Schmitt, M. (2022). *When should an high-banking strategy and temporary with Model-Agnostic Interpretable Data-driven riskPricing*. *Journal of Risk and Insurance*, 91(1), 1-18.

Abstract
This paper presents a data-driven strategy for the construction of insurance tariff classes. The strategy is based on the use of machine learning techniques to identify the most relevant features for the construction of tariff classes. The strategy is implemented using a decision tree algorithm. The results of the study show that the proposed strategy outperforms traditional methods in terms of predictive accuracy and stability.

Introduction
The insurance industry is facing a number of challenges, including the need to improve its pricing accuracy and to reduce its operational costs. One of the key challenges is the construction of insurance tariff classes. This paper presents a data-driven strategy for the construction of insurance tariff classes. The strategy is based on the use of machine learning techniques to identify the most relevant features for the construction of tariff classes. The strategy is implemented using a decision tree algorithm. The results of the study show that the proposed strategy outperforms traditional methods in terms of predictive accuracy and stability.

Expert Syst. Appl.

[Henckaerts & Antonio, 2022]

Henckaerts, O., & Antonio, R. (2022). *The added value of dynamically updating reinsurance prices with alternative collected driving behavior data*. *Journal of Risk and Insurance*, 91(1), 1-18.

Abstract
This paper presents a data-driven strategy for the construction of insurance tariff classes. The strategy is based on the use of machine learning techniques to identify the most relevant features for the construction of tariff classes. The strategy is implemented using a decision tree algorithm. The results of the study show that the proposed strategy outperforms traditional methods in terms of predictive accuracy and stability.

Introduction
The insurance industry is facing a number of challenges, including the need to improve its pricing accuracy and to reduce its operational costs. One of the key challenges is the construction of insurance tariff classes. This paper presents a data-driven strategy for the construction of insurance tariff classes. The strategy is based on the use of machine learning techniques to identify the most relevant features for the construction of tariff classes. The strategy is implemented using a decision tree algorithm. The results of the study show that the proposed strategy outperforms traditional methods in terms of predictive accuracy and stability.

IME



github/henckr/distRforest



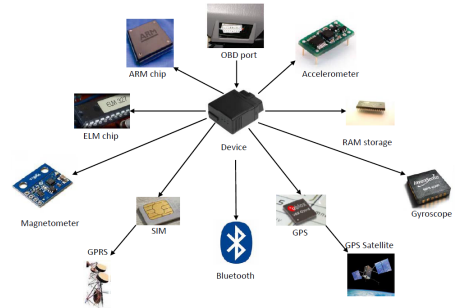
github/henckr/maidrr

- ▶ Denote for policyholder i in a given policy period:
 - e_i : exposure-to-risk
 - N_i : number of claims filed during the exposure period
 - L_i : total loss amount reported during the exposure period.
- ▶ The **pure premium** π_i :

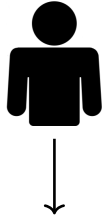
$$\pi_i = \mathbb{E} \left[\frac{L_i}{e_i} \right] \stackrel{\text{indep.}}{=} \mathbb{E} \left[\frac{N_i}{e_i} \right] \times \mathbb{E} \left[\frac{L_i}{N_i} \mid N_i > 0 \right] = \underbrace{\text{Freq}_i}_{\text{frequency}} \times \underbrace{\text{Sev}_i}_{\text{severity}}$$

- ▶ Build $f(\text{risk factors})$ to predict frequency and severity, respectively.

Products: usage-based insurance (UBI)
pay-as-you-drive (PAYD)
pay-how-you-drive (PHYD)



- ▶ **Telematics** is the integrated use of **telecommunications** and **informatics**.
- ▶ Black-box device is installed in the vehicle.
- ▶ **Real driving behavior** is monitored.
- ▶ Very often targets **young drivers**.



Static, demographic data



License age



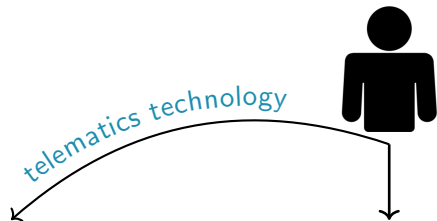
Car make/model



Type of fuel



Postal code



Driving habits

Static, demographic data



Mileage



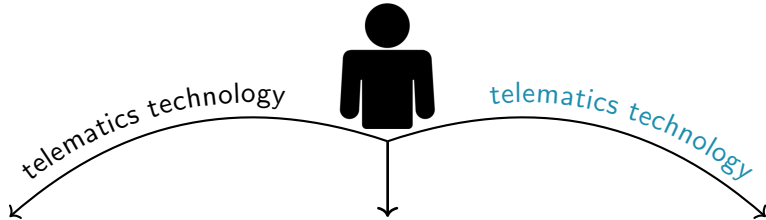
Travel time



Time slot



Road type



Driving habits

Static, demographic data

Driving style



Speed



Acceleration

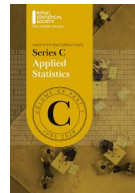


Attention



Weather

Insurance analytics literature on telematics



- ▶ Verbelen, Antonio & Claeskens (2018, JRSS C):
 - claim frequency models with classic, static features and driving habit information
 - compositional data and their use in GAMs.

- ▶ Wüthrich (2017, EAJ), Gao & Wüthrich (2018, EAJ), Gao et al. (2019, SAJ) and more papers:
 - the construction of $v - a$ **heatmaps from GPS signals**
 - feature-engineering on these heatmaps
 - use of these features in claim frequency models.



[Open Access](#) [Printable Paper](#) [Article](#)

Address Identification Using Telematics: An Algorithm to Identify Dwell Locations

by [Christopher Grumiau](#)¹ , [Mina Mostoufi](#)^{1*} , [Solon Pavlioglou](#)^{1*}  and [Tim Verdonck](#)^{2,3} 

¹ Allianz Benelux, 1000 Brussels, Belgium

² Department of Mathematics (Faculty of Science), University of Antwerp, 2000 Antwerpen, Belgium

³ Department of Mathematics (Faculty of Science), Katholieke Universiteit Leuven, 3000 Leuven, Belgium

* Authors to whom correspondence should be addressed

Risks 2020, 8(3), 92; <https://doi.org/10.3390/Risks8030092>

Received: 16 June 2020 / Revised: 7 August 2020 / Accepted: 21 August 2020 / Published: 1 September 2020

(This article belongs to the Special Issue Data Mining in Actuarial Science: Theory and Applications)

- ▶ Denuit, Guillen & Trufin (2019, Annals of Actuarial Science) on **Multivariate credibility modelling for usage-based motor insurance pricing with behavioural data.**
- ▶ Grumiau, Mostoufi, Pavlioglou & Verdonck (2020, Risks) on **Address identification using telematics: an algorithm to identify dwell locations.**
- ▶ Banghee So, J.-P. Boucher & E. Valdez (2021, Risks) on **Synthetic dataset generation of driver telematics.**

Managerial insights, based on Carbone & Taub (2018) **UBI insurance is not usage-based. Sorry, not sorry!**

- In 2017, 14 million policies sent telematics data to insurers around the world.
- However, less than 9 percent of the global insurance telematics policies were characterized by **usage-based pricing**.
- **Use of driving data in pricing:**
 - * use driving score at underwriting stage
 - * propose tailored renewal price (with discounts, or discounts + surcharges)
 - * usage-based, i.e. charge price for period of coverage based on how policyholder behaves during this period, and avoid **premium leakage**.

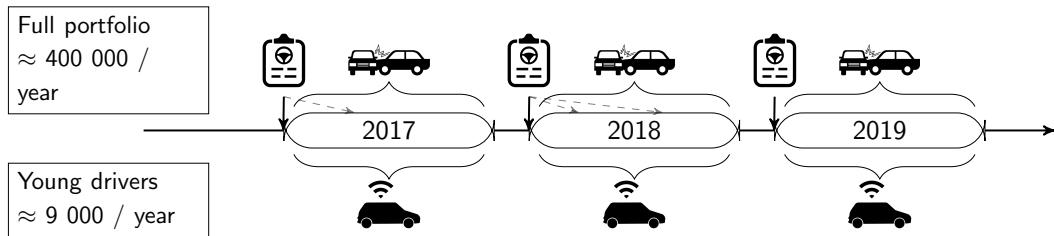
Our focus in this talk:

- How to use driving behavior (i.e. habits + style) to **update** a baseline **tariff** (with only self-reported characteristics)?
- What is the **added value** of telematics for pricing via risk classification?
- Managerial insights? Impact on retention rates, profit?

Focus on **frequency, severity and churn models** in the presence of static self-reported characteristics as well as telematics collected data.

Aim for an **explainable** updating mechanism.

Data and methodology



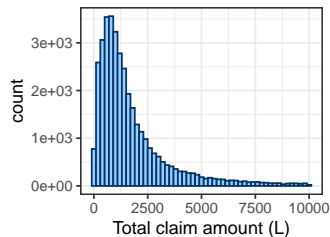
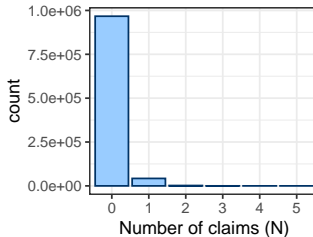
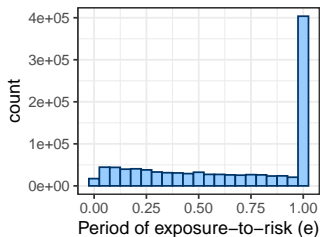
Policy information at the start of the policy period, subject to possible changes during the policy period (e.g., new vehicle).



Claims reported (68 196 in total).

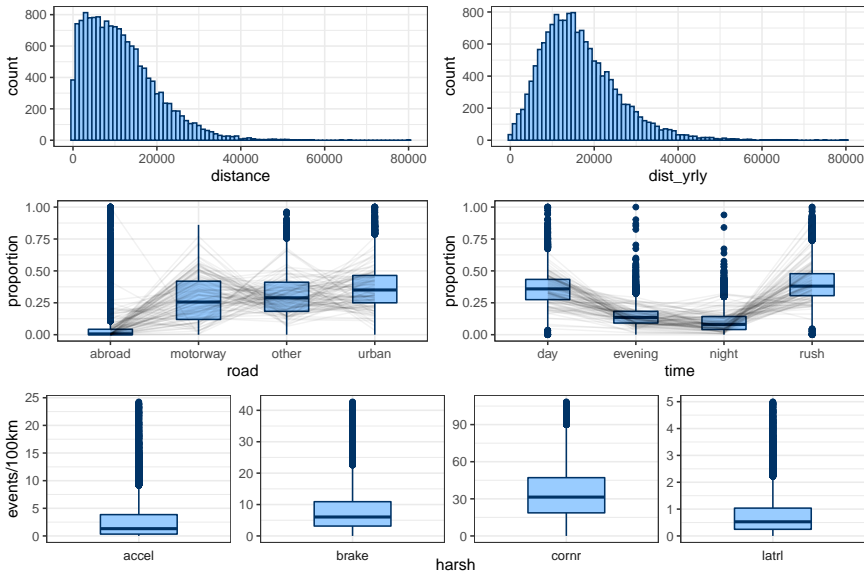


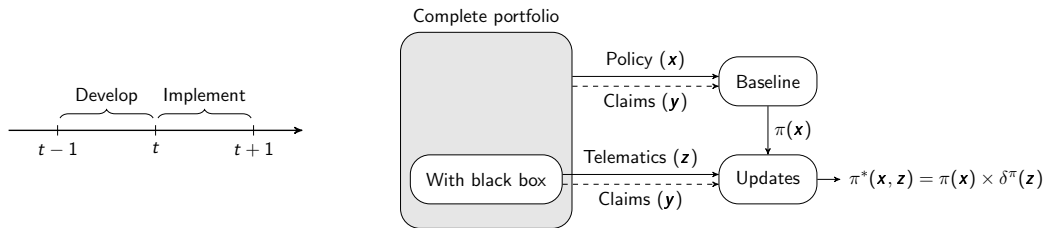
Driving behavior (only for drivers < 26 years at underwriting time, 308M kilometers driven in total).



Policy information with **self-reported** risk characteristics:

- **driver**: age, experience, additional young drivers, etc.
- **payments**: frequency and SEPA indicator
- **geographical**: postal code and mosaic segment
- **vehicle**: age, weight, value, power, fuel, make, etc.





The idea:

- charge baseline tariff $\pi(\mathbf{x})$ at t
- ex post, multiplicative update $\delta^\pi(\mathbf{z})$ at $t + 1$, based on driving data in $[t, t + 1]$.

Baseline pricing and churn models

- ▶ Predictive models for the **complete portfolio** using traditional features \mathbf{x}
 - claim frequency *and* severity \rightarrow **tariff**
 - customer churn prediction \rightarrow client **retention** analysis.

- ▶ Stochastic gradient boosting (Friedman, 2002) with the following **assumptions**:

| | Distribution | Prediction $f(\mathbf{x})$ | Loss function $D(y, f(\mathbf{x}))$ |
|-----------------|---------------------------|---------------------------------|---|
| Claim frequency | $N \sim \text{Poisson}$ | $\mathbb{E}(N \mathbf{x}, e)$ | $\frac{2}{n} \sum_{i=1}^n \left[y_i \ln \left\{ \frac{y_i}{f_i} \right\} - \{y_i - f_i\} \right]$ |
| Claim severity | $L/N \sim \text{gamma}$ | $\mathbb{E}(L/N \mathbf{x})$ | $\frac{2}{\sum_i N_i} \sum_{i=1}^n N_i \left[\frac{y_i - f_i}{f_i} - \ln \left\{ \frac{y_i}{f_i} \right\} \right]$ |
| Customer churn | $C \sim \text{Bernoulli}$ | $\mathbb{E}(C \mathbf{x})$ | $-\frac{1}{n} \sum_{i=1}^n [y_i \ln \{f_i\} + (1 - y_i) \ln \{f_i\}]$ |

▶ Parameter **tuning**:

H2O random grid search + 5-fold cross-validation (LeDell et al., 2020).

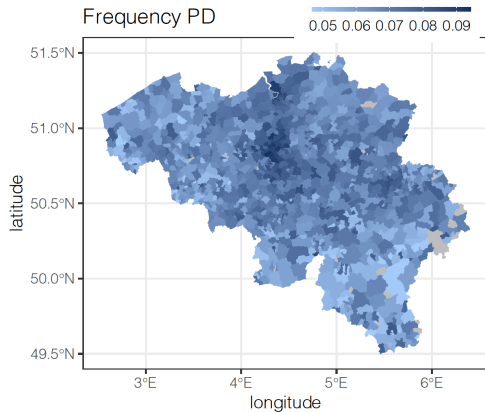
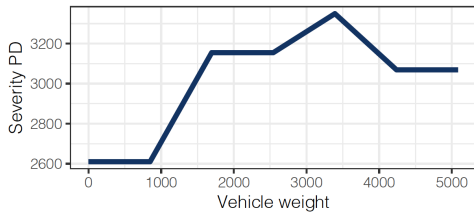
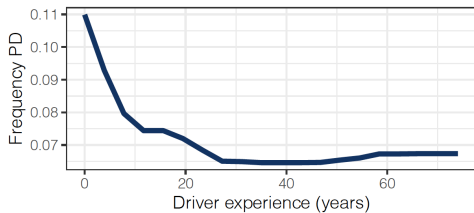
▶ Enforce the **balance property** by scaling predictions (for the young drivers):

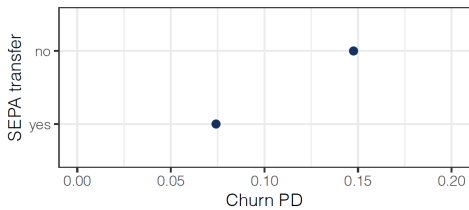
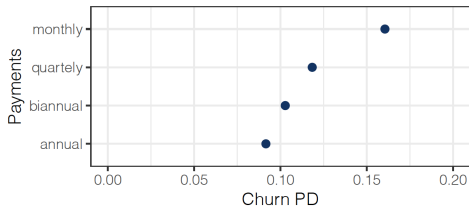
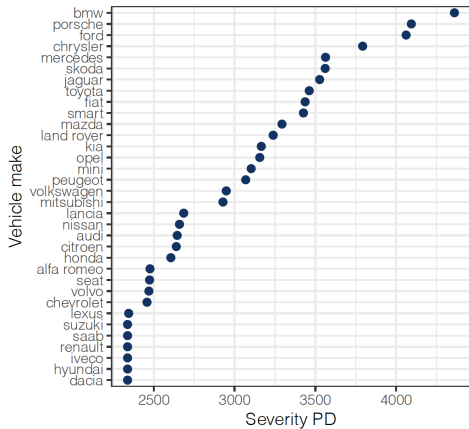
$$\sum_{i=1}^n f_i = \sum_{i=1}^n y_i$$

automatically fulfilled for GLMs with canonical link

see e.g. Wüthrich (2020).

| Rank | Claim frequency | | Claim severity | | Customer churn | |
|----------|--------------------|-------|--------------------|-------|-----------------|-------|
| | Feature | % | Feature | % | Feature | % |
| 1 | geo_postcode | 34.72 | veh_weight | 23.21 | paym_split | 43.48 |
| 2 | driv_experience | 14.08 | veh_make | 21.37 | geo_postcode | 11.67 |
| 3 | driv_seniority | 8.52 | geo_postcode | 10.54 | veh_age | 9.85 |
| 4 | veh_make | 6.25 | veh_segment | 10.48 | paym_sepa | 9.44 |
| 5 | geo_mosaic | 5.85 | geo_mosaic | 6.59 | driv_seniority | 6.90 |
| 6 | veh_fuel | 5.09 | driv_seniority | 5.83 | veh_make | 3.43 |
| 7 | veh_segment | 4.66 | veh_value | 3.50 | driv_experience | 2.85 |
| 8 | paym_split | 3.91 | veh_age | 3.44 | geo_mosaic | 2.45 |
| 9 | driv_add_younger26 | 3.29 | driv_experience | 2.98 | driv_age | 2.43 |
| 10 | driv_age | 2.75 | driv_add_younger26 | 2.91 | veh_use | 1.99 |
| Σ | | 89.12 | | 90.86 | | 94.48 |





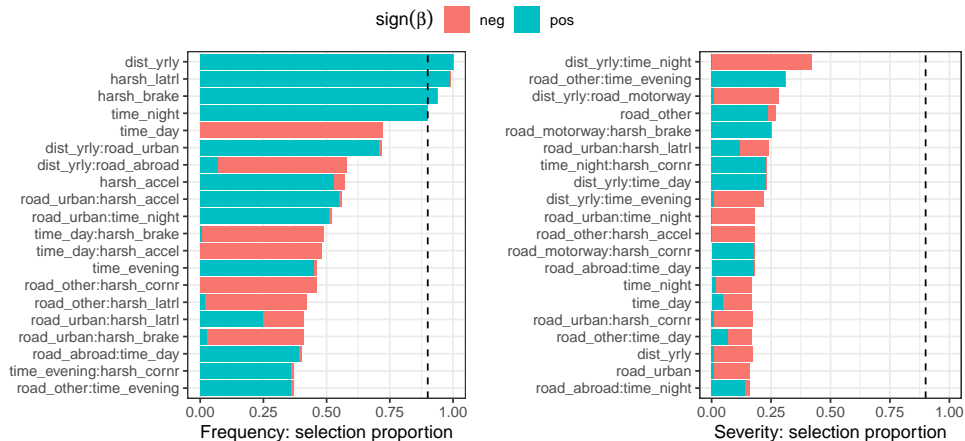
Updating pricing with driving behavior

- ▶ Aim is to update premiums for the drivers with telematics features \mathbf{z} .
- ▶ Log-link GLM with the baseline prediction $\ln[f(\mathbf{x})]$ as an offset:

$$\ln[\mathbb{E}(y \mid \mathbf{x}, \mathbf{z})] = \ln[f(\mathbf{x})] + \beta_0 + \sum_{j=1}^p \beta_j z_j$$

$$\mathbb{E}(y \mid \mathbf{x}, \mathbf{z}) = f(\mathbf{x}) \times \exp(\beta_0) \times \prod_{j=1}^p \exp(\beta_j z_j).$$

- ▶ Updated prediction is then **multiplicative**:
 - baseline GBM prediction $f(\mathbf{x})$ for a policyholder with risk characteristics \mathbf{x}
 - overall update factor $\exp(\beta_0)$ via the intercept
 - update $\exp(\beta_j z_j)$ from each telematics feature z_j .

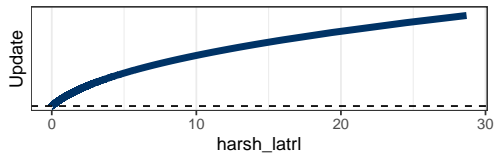
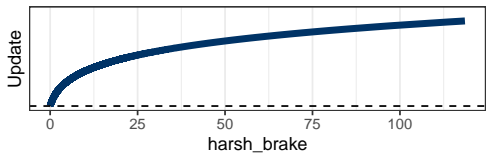
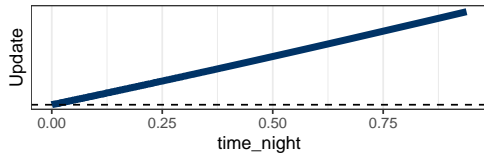
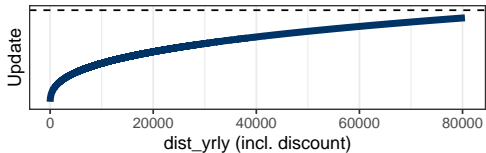


- ▶ Let $\mathbf{z}^* \in \mathbb{R}^4$ denote $\{\text{dist_yrly}, \text{harsh_latrl}, \text{harsh_brake}, \text{time_night}\}$.
- ▶ Log-link Poisson GLM with offset for **claim frequency**:

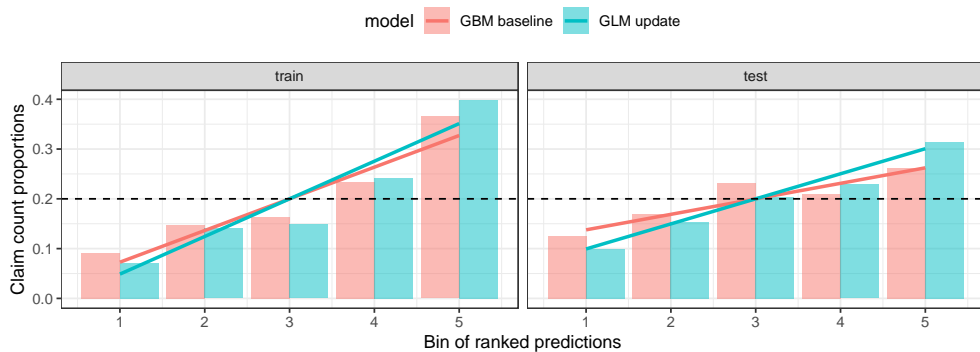
$$\ln[\mathbb{E}(N | \mathbf{x}, \mathbf{z}^*)] = \ln[\mathbb{E}(N | \mathbf{x}, e)] + \beta_0 + \sum_{j=1}^4 \beta_j \log(z_j^* + 1)$$

$$\mathbb{E}(N | \mathbf{x}, \mathbf{z}^*) = \mathbb{E}(N | \mathbf{x}, e) \times \exp(\beta_0) \times \prod_{j=1}^4 (z_j^* + 1)^{\beta_j}.$$

- ▶ Updated prediction is **multiplicative** in the following terms:
 - baseline GBM prediction $\mathbb{E}(N | \mathbf{x}, e)$ for a policyholder with risk characteristics \mathbf{x}
 - overall discount factor $\exp(\beta_0) \approx 2\%$
 - update $(z_j^* + 1)^{\beta_j}$ from each telematics feature z_j^* .

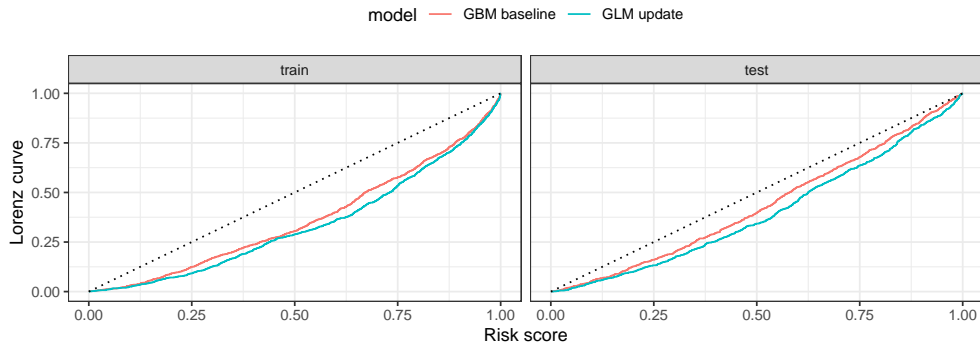


- ▶ Mileage + discount remains < 1 .
- ▶ Penalty once night driving, harsh braking or lateral events are registered.
- ▶ Safe driving is the key to earn discounts!



Here we consider:

- $r_i^m = F_n(f^m(\mathbf{x}_i, \mathbf{z}_i^*))$ with $F_n(\cdot)$ the ecdf
- $PC^m(s) = \frac{\sum_{i=1}^n N_i \mathbb{1}\{\frac{s-1}{5} < r_i^m \leq \frac{s}{5}\}}{\sum_{i=1}^n N_i}$ for $s \in \{1, \dots, 5\}$.



Here we consider:

- $LC^m(s) = \frac{\sum_{i=1}^n N_i \mathbb{1}\{r_i^m \leq s\}}{\sum_{i=1}^n N_i}$ for $s \in [0, 1]$.

Managerial insights

For the discussion of managerial insights, I refer to our paper:

- **adjust baseline churn** $\rho(\mathbf{x})$ to $\rho^*(\mathbf{x}, \delta^\pi) = \rho(\mathbf{x}) + \epsilon_p \cdot (\delta^\pi - 1) = \rho(\mathbf{x}) + \delta^\rho$, with ϵ_p the price elasticity
- study expected **profit and retention rate**

$$P = \frac{1}{n} \sum_{i=1}^n (1 - (\rho_i + \delta_i^\rho)) \cdot (\delta_i^\pi \pi_i - L_i) \quad R = \frac{1}{n} \sum_{i=1}^n 1 - (\rho_i + \delta_i^\rho)$$

- restrict penalties/discounts + redistribute \Rightarrow **fairness, solidarity, commercially appealing**

$$\delta_{lo}^\pi \leq \delta^\pi \leq \delta_{hi}^\pi$$

$$\sum_{i=1}^n (1 - \rho_i) \cdot \pi_i = \sum_{i=1}^n (1 - \rho_i) \cdot \alpha \cdot \delta_i^\pi \cdot \pi_i$$

Our paper puts focus on:

- a baseline pricing model with self-reported characteristics
- an explainable updating mechanism to incorporate driving behavioral information.

Added value of telematics for insurance pricing is studied from both a statistical and managerial perspective.

For more information, please visit:

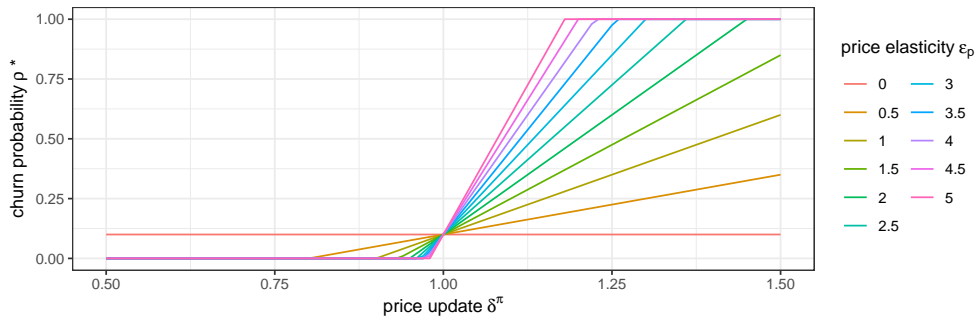
- draft of the paper, including complete set of references
- LRisk website, www.lrisk.be
- my homepage <https://katrienantonio.github.io>.

Special thanks to

- the organizers
- the companies and funding agencies supporting/having supported my research lab: Ageas, Argenta, Atlas Copco, CNP Assurances, FWO, KU Leuven internal funds.

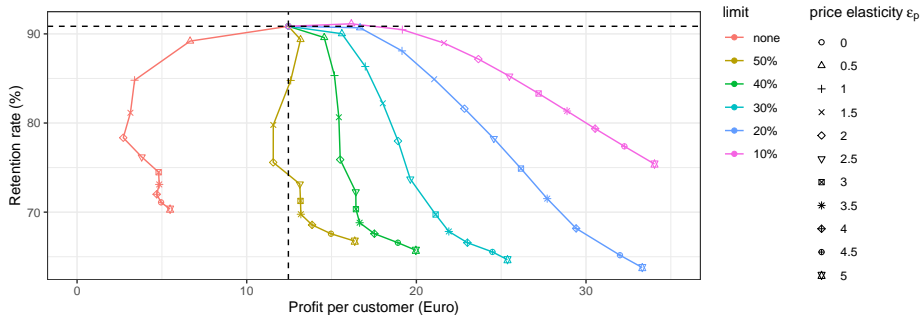
Extra sheets

- ▶ Adjust baseline churn via price **update** δ^π and **elasticity** of demand ϵ_p .
- ▶ $\rho^*(\mathbf{x}, \delta^\pi) = \rho(\mathbf{x}) + \epsilon_p \cdot (\delta^\pi - 1)$ for $\rho(\mathbf{x}) = 0.1$ and $\epsilon_p \in [0, 5]$:



- ▶ $\rho^* = \rho$ when $\delta^\pi = \pi^*/\pi = 1$ (no price change).
- ▶ Linear increases/decreases ($\delta^\pi > 1$ / $\delta^\pi < 1$) with slope ϵ_p .

- ▶ Expected profits and retention rates under different scenarios:



- ▶ Stricter limits result in higher profits
- ▶ Profits increase with ϵ_p at the cost of lower retention
- ▶ No limit results in lower profits than baseline (driven by low premiums on average)

- ▶ Maximize expected profit P while retaining a minimum proportion of the portfolio R^* :

$$\begin{aligned} \max_{\alpha} P(\alpha) &= \frac{1}{n} \sum_{i=1}^n (1 - (\rho_i + \delta_i^{\rho})) \cdot (\alpha \delta_i^{\pi} \pi_i - L_i) \\ \text{subject to } R(\alpha) &= \frac{1}{n} \sum_{i=1}^n 1 - (\rho_i + \delta_i^{\rho}) \geq R^* \\ \delta_{lo}^{\pi} &\leq \delta^{\pi} \leq \delta_{hi}^{\pi}. \end{aligned}$$

- ▶ Implicit dependence of R on α as $\delta^{\rho} = \epsilon_{\rho} \cdot (\alpha \delta^{\pi} - 1)$.
- ▶ Efficient frontier by varying R^* over a range of values and maximizing $P(R^*)$ via α .

Efficient frontiers for $R^* \in [0.75, 0.9]$ under different scenarios

