

Longitudinal analysis of distance traveled

Longitudinal analysis of distance traveled

Context

- ▶ New technologies such as **GPS-collected data** have emerged, which offer new ways to approach car insurance pricing.
 - ▶ Processing these data provides **reliable information** about drivers' behavior.
- One piece of GPS-collected information that is directly related to the risk insured is distance driven.

Relevance

- Covariates such as territory, gender and age only describe the **general behavior** of insured in those groups.
- ▶ Ayuso et al. (2016b) shows that the **differences** observed in claims frequency between men and women are largely attributable to **vehicle use** ;
 - ▶ Verbelen et al. (2018) reached a similar conclusion
- In a social-political context where the use of gender in ratemaking is restricted or criticize, calculating premiums on **more objective information** is of interest.

Overview

Objective

Using telematics data, we study the relationship between **claim frequency** and **distance driven** through different models by observing **smooth functions**.

- 1 Generalized Additive Models (**GAM**) for a Poisson distribution (fixed effects),
- 2 Generalized Additive Models for Location, Scale, and Shape (**GAMLSS**) that we generalize for panel count data (random effects).

Why GPS-collected data ?

- ▶ As shown by many authors, such as Lemaire et al. (2016), the **self-reported** approximation of the distance driven is **not reliable** and is often very different from the exact distance driven.
- ▶ There are **important differences between driving uses and driving habits**, which justifies consideration of other measures than exposure time in the modeling.

A First Model

Starting Point

Boucher et al. (2017), by using a **GAM Poisson model**, analyzed the influence of duration and distance driven on the number of claims with **independent cubic splines** : $\log(\mu_i) = \beta_0 + s_1(km_i) + s_2(d_i)$.

$$\begin{aligned} \mu_{i,t} &= \exp(\mathbf{X}_{i,t}\beta + s_1(km) + s_2(d)) \\ &= \exp(s_1(km)) \exp(s_2(d)) \exp(\mathbf{X}_{i,t}\beta) \\ &= \exp(s_1(km)) \exp(s_2(d)) \lambda_{i,t}. \end{aligned} \quad (1)$$

GAM

- ▶ GAMs : introduced by Hastie and Tibshirani (1986).
- ▶ Extension of the generalized linear models (GLM) theory : relax the hypothesis of linearity, and smoothing functions s of the covariates could be included in the predictor.
- ▶ Example : the mean for an individual i could be given by $g(\mu_i) = s_0 + s_1(x_{1,i}) + s_2(x_{2,i}) + s_3(x_{3,i})$.

A First Model

What do you think ?

We model $N_{it} \sim \text{Pois}(\mu_{it})$, where $\mu_{it} = \exp(s_1(km)) \exp(s_2(d)) \lambda_{i,t}$ with real canadian insurance data.

Questions :

- 1 What the relation between $\exp(s_1(km))$ and claim frequency would look like when a **linear trend is not imposed** by the model structure ?
- 2 And $\exp(s_2(d))$?

To help you :

- ▶ Would it be nonetheless nearly linear ?
- ▶ Would it stop increasing at some point ?
- ▶ Would it start to decline at some point ? Would it go up again ?
- ▶ Any other intuition ?

A First Model

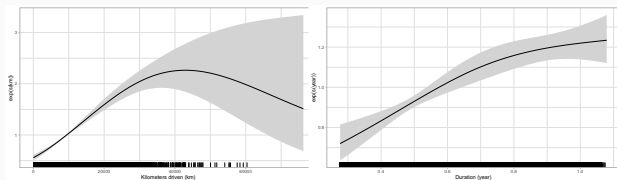


Figure 1: $\exp(\hat{s}_1(km))$ and $\exp(\hat{s}_2(year))$ from the Poisson GAM

Case Study

- 1 All models are illustrated using data from a **major Canadian insurance company**.
- 2 The model $\log(\mu_i) = \beta_0 + s_1(km_i) + s_2(d_i)$ yields **similar results** to those obtained by Boucher et al. (2017) (**Spanish data**).
- 3 In the study by Boucher et al. (2017), a "learning effect" is advanced to justify the look of $\exp(\hat{s}_1(km))$.

A First Model

Consistency problem

The **slope** could change as distance increases, but it should always be **strictly positive** since the **risk is greater**, meaning that the smoothing function should always be increasing.

- ▶ **One explanation** comes from the fact that GAM **supposes independence** between all contracts of the same insured.

Results Analysis

One can argue that distance driven is **correlated** with **other driving habits** resulting from driving experience, (Ferreira and Minikel (2010)).

- 1 The model **does not** take this correlation **into account**.
- 2 The resulting relationship between claim frequency and the distance driven do **not** give an **appropriate representation** of how the claim frequency could change **when insureds change their driving habits**.

We think that the shape of the smoothing function comes from the **driver profiles** : the lower quantiles of the distribution of the distance driven does not come from the same (type of) drivers as the higher quantiles.

A Longitudinal Analysis

Search for a "marginal" effect

- 1 The objective is not to compute a premium.
- 2 The objective is mainly to understand **how the distance impacts** the claim frequency when **all individual characteristics** of policyholders have been **considered**.

Panel Data Modeling

In non-life insurance, however, we can observe the same insured over **many** contracts.

- ▶ Instead of modeling the marginal distribution of each $N_{i,t}$ for $t = 1, \dots, T$, we are now looking for the **joint distribution** :

$$\Pr(N_1 = n_1, N_2 = n_2, \dots, N_T = n_T) = \Pr(N_1 = n_1) \times \Pr(N_2 = n_2 | N_1 = n_1) \times \dots \times \Pr(N_T = n_T | N_1 = n_1, \dots, N_{T-1} = n_{T-1}),$$

A Longitudinal Analysis

Construct Multivariate Count Models

- ▶ One popular way, is to **include an individual parameter** α in the mean parameter of the count distribution of each contract t :

$$N_{i,t} | \alpha_j \sim \text{Poisson}(\mu_{i,t} = \alpha_j \lambda_{i,t}), \quad (2)$$

Random vs Fixed effects

We can consider two different situations regarding this parameter :

- 1 All $\alpha_j, j = 1, \dots, n$ are i.i.d. **random variables** that come from a selected **prior distribution** (we call this the random effects model) ;
- 2 All $\alpha_j, j = 1, \dots, n$ are **unknown parameters** that need to be **estimated** (we call this the fixed effects model).

Random Effects Model

Model Specification

In random effects models, we suppose that α_i , $i = 1, \dots, n$, are **random variables**, with prior density $f(\cdot)$.

- Conditionally on the random effects α_i^{RE} , all numbers of claims $N_{i,1}, N_{i,2}, \dots, N_{i,T}$ from insured i are independent.

$$\Pr[N_{i,1} = n_{i,1}, \dots, N_{i,T} = n_{i,T}] = \int_0^\infty \left(\prod_{t=1}^T \exp(-\alpha_i^{RE} \lambda_{i,t}^{RE}) \frac{(\alpha_i^{RE} \lambda_{i,t}^{RE})^{n_{i,t}}}{n_{i,t}!} \right) f(\alpha_i^{RE}) d\alpha_i^{RE} \quad (3)$$

- Many distributions can be used for α_i^{RE} , such as the **gamma** or the inverse Gaussian.

Gamma Distribution

If we suppose that α_i^{RE} follows a **gamma distribution of mean 1 and variance $\frac{1}{\nu}$** , the joint distribution can be expressed as :

$$\Pr[N_{i,1} = n_{i,1}, \dots, N_{i,T} = n_{i,T}] = \left(\prod_{t=1}^T \frac{(\lambda_{i,t}^{RE})^{n_{i,t}}}{n_{i,t}!} \right) \frac{\Gamma(n_{i,\bullet} + \nu)}{\Gamma(\nu)} \left(\frac{\nu}{\lambda_{i,\bullet}^{RE} + \nu} \right)^\nu (\lambda_{i,\bullet}^{RE} + \nu)^{-n_{i,\bullet}} \quad (4)$$

$$(n_{i,\bullet} = \sum_{t=1}^T n_{i,t} \text{ and } \lambda_{i,\bullet}^{RE} = \sum_{t=1}^T \lambda_{i,t}^{RE})$$

Random Effects Model

MVNB

This well-known distribution is the multivariate negative binomial distribution.

- This distribution is a **generalization** of the **negative binomial distribution**.
- It is a basic distribution for panel count data modeling with overdispersion ($E[N_{i,t}] = \lambda_{i,t}^{RE} < V[N_{i,t}] = \lambda_{i,t}^{RE} + (\lambda_{i,t}^{RE})^2 / \nu$).
- It is **not** a member of the **linear exponential family**.
- GAM theory **cannot be used** to include smoothing functions.

It can be shown that the first-order condition to obtain $\hat{\beta}_{MLE}$ is :

$$\sum_{i=1}^n \sum_{t=1}^T x_{i,t} \left(n_{i,t} - \lambda_{i,t}^{RE} \frac{n_{i,\bullet} + \nu}{\lambda_{i,\bullet}^{RE} + \nu} \right) = 0. \quad (5)$$

Random Effects Model

GAMLSS

Instead, we use **Generalized Additive Models for Location, Scale and Shape** theory, that can be used for other distributions than the members of the linear exponential family of distribution.

- More flexible** : can model a location parameter μ_i , a variance parameter σ_i (scale), a skewness parameter ν_i and a kurtosis parameter τ_i as additive functions of the covariates.

$$g_k(\theta_k) = X_k \beta_k + \sum_{j=1}^{J_k} Z_{j,k} \gamma_{j,k} \quad (6)$$

- $\theta = [\mu, \sigma, \nu, \tau]$. μ, σ, ν and τ are vectors with n elements
- If a smooth function can be expressed in linear form, Equation (6) can be rewritten as

$$g_k(\theta_k) = X_k \beta_k + \sum_{j=1}^{J_k} h_{j,k}(x_{j,k}),$$

where $h_{j,k}$ is a smooth non-parametric function.

Random Effects Model

Model Specification

It is possible to use a GAMLSS that specify **only the location parameter**. In this case, θ would simply become $\theta = \{\mu\}$.

- 1 We choose to model the parameter $\lambda_{i,t}$ with smoothing function ;
- 2 v is kept **constant** for all individuals.

R package

- 1 To use GAMLSS, many distributions are available in the R package *gamlss*.
 - 2 Unfortunately, the MVNB distribution is not one of them.
 - 3 The distribution is however implemented by itself in the package *multinbmod*.
- Consequently, we have to **write our own code** for convenience.

Random Effects Model

What do you think ?

We model $N \sim MVNB(\mu, v)$, where $\mu = \exp(s_1(\text{km})) \exp(s_2(d)) \lambda$ with real canadian insurance data.

Questions :

- 1 What the **relation** between $\exp(s_1(\text{km}))$ [$\exp(s_2(d))$] and **claim frequency** would look like ?
- 2 **How** would the results **differ** from the **previous model** ?

To help you :

- ▶ Would it be nonetheless nearly linear ?
- ▶ Would it stop increasing at some point ?
- ▶ Would it start to decline at some point ? Would it go up again ?
- ▶ Any other intuition ?

Random Effects Model

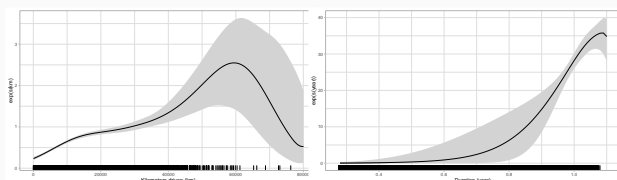


Figure 2: $\exp(\hat{s}_1(\text{km}))$ and $\exp(\hat{s}_2(\text{year}))$ from the GAMLSS with random effects model

Model Fitting

- 1 To fit the model, we maximize a **penalized log-likelihood** function l_p , integrating a quadratic penalty $\gamma^T G \gamma$.
- 2 **Penalty matrix** G : very often define as $\Lambda D_r^T D_r$ (different formulations possible).
- 3 A hyper-parameter, noted here $\Lambda \in \mathbb{R}^+$, controls the **weight given to the penalty**. The **greater** its value, the **smoother** the resulting estimated function.
- 4 To select the penalty parameters in $G(\Lambda)$ associated with both p-splines, we **test out** multiple **combinations** of values of $\Lambda = \{\Lambda_1, \Lambda_2\}$.

A Fixed Effects Approach

The model

Poisson fixed effects model can be seen as a **basic Poisson** regression model without an intercept. Being part of the linear exponential family of distribution, **GAM theory** can then be used when smoothing functions are added to the mean parameter of the distribution.

In practice, as mentioned, it is **relatively easy** to implement the fixed effects model with R; we simply used the *gam* function from the package *mgcv*.

- 1 To include fixed effects in the model the **intercept** of the model is **dropped**.
- 2 We include a **unique identifier** variable for each policyholder as a factor variable and we include the **distance driven** in the model using a **cubic spline s**.

A Fixed Effects Approach

Parameters estimation

In the fixed effects model, we consider each α_i , $i \in \{1, \dots, n\}$ as an **unknown parameter**.

- 1 At least $n + p + 1$ parameters should be estimated, which is quite a high number of parameters given that T_i is usually small for insurance datasets.
- 2 The large number of parameters in the model causes what is called **incidental problem**, which means that an incorrect estimation of the fixed effects α generates **incorrect estimates** of β associated with covariates in the mean.
- 3 It has been shown that a fixed effects model based on a Poisson distribution **does not have this problem** (see (Cameron and Trivedi, 2013)) for a detailed explanation).

First-order condition equation

For the β parameters, the first condition by MLE can be shown to be equal to :

$$\sum_{i=1}^n \sum_{t=1}^{T_i} x_{i,t} \left(n_{i,t} - \lambda_{i,t}^{FE} \frac{n_{i,t}}{\lambda_{i,t}^{FE}} \right) = 0. \quad (7)$$

When we compare the first-order condition equation of the random effects model and (7), we see that when T is large, or when $v \rightarrow 0$, **random and fixed** effects models are **equivalent**.

A Fixed Effects Approach

What do you think ?

We model $N_{i,t} \sim \text{Pois}(\mu_{i,t})$, where $\mu_{i,t} = \exp(a_i) \exp(s(km))$.

Questions :

- 1 What the **relation** between $\exp(s(km))$ and claim frequency would look like ?
- 2 Will the **"learning effect"** be there again ?

Rating structure based on distance driven

We decided to model the Poisson fixed effects by **not including** a smoothing function for the duration.

- 1 Our objective is to measure the **marginal effect** of the distance on the claim frequency. If we want to measure the risk of each additional kilometer the insured decides to drive, the **duration** of the contract is **not important**.
- 2 We want to construct a rating structure based solely on the distance driven as a risk measure.

A Fixed Effects Approach

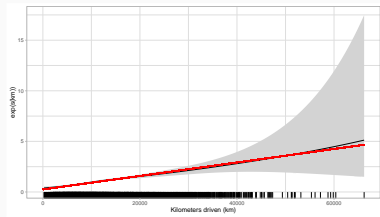


Figure 3: GAM with fixed effects estimated with Canadian data

Results Analysis

- 1 The relationship between distance traveled and claim frequency is **always increasing**, and is even **almost linear**.
- 2 What has been called the “learning effect” has **disappeared**.
- 3 We observe a **much more logical** and **coherent** relationship between distance traveled and frequency than before.

Longitudinal analysis of distance traveled
 ○○○○○○○○○○○○○○○○○●○○○○○

Claim Classification Using Partial Telematics Information
 ○○○○○○○○○○○○○○○○

A Fixed Effects Approach

Marginal impact of each additional kilometer

- 1 The relationship between claim frequency and the distance driven should be understood as the **marginal impact** of **each** additional kilometer driven or not-driven.
- 2 Explicitly, as we approximated $\exp(s(km))$ by $0.25 + \frac{1}{15000} km_{i,t}$ (the red line), we then have

$$\begin{aligned}
 N_{it} &\sim \text{Poisson}(\exp(\alpha_i) \exp(s(km))) \\
 &\sim \text{Poisson}(\exp(\alpha_i)(a + b km_{i,t})) \\
 &\sim \text{Poisson}\left(0.25 \exp(\alpha_i) + \frac{1}{15000} \exp(\alpha_i) km_{i,t}\right).
 \end{aligned}$$

- 3 We see that the **slope**, i.e., the marginal impact of each additional kilometer driven or not-driven, is **not the same** for each insured because it **depends on α_i** .

Longitudinal analysis of distance traveled
 ○○○○○○○○○○○○○○○○○●○○○○○

Claim Classification Using Partial Telematics Information
 ○○○○○○○○○○○○○○○○

A Fixed Effects Approach

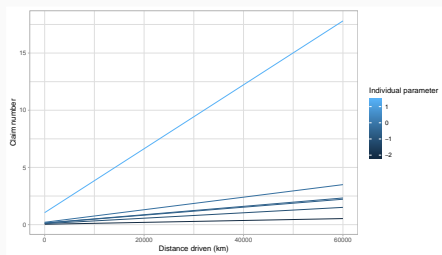


Figure 4: Exposure measure for different individual parameters.

Results Analysis II

- With this model, we then **reconcile** the intuition that **each kilometer** should increase the risk for an individual, but that this increase could be **different for each driver**.

Longitudinal analysis of distance traveled
 ○○○○○○○○○○○○○○○○○●○○○○○

Claim Classification Using Partial Telematics Information
 ○○○○○○○○○○○○○○○○

A Fixed Effects Approach

“Learning effect”

In summary, **instead** of referring to the “learning effect” to understand the left-hand graph of Cross-sectional data model, we **should understand** instead that

- 1 Typical insureds who drive more than 60,000 km per year are **better risks per kilometer** than insureds who drive approximately 40,000 km per year.
- 2 However, for each driver, independently of their driving risk *per kilometer*, the risk of an accident will always **increase for each additional kilometer driven** (by approximately $\frac{1}{15,000}$).

Comparative Analysis

Which Effect Should Be Used in Practice ?

The fixed effects model is **more general** than the random effects model, which means that in case of contradictory results, **fixed effects** should always be **preferred**.

$$\begin{aligned} \Pr[N_{i,1} = n_{i,1}, \dots, N_{i,T} = n_{i,T}] &= \int_0^\infty \Pr[N_{i,1} = n_{i,1}, \dots, N_{i,T} = n_{i,T} | x_{i,1}, \dots, x_{i,T}, a_i^{RE}] f(a_i^{RE} | x_{i,1}, \dots, x_{i,T}) da_i^{RE} \\ &= \int_0^\infty \left(\prod_{t=1}^T \Pr[N_{i,t} = n_{i,t} | x_{i,1}, \dots, x_{i,T}, a_i^{RE}] \right) f(a_i^{RE}) da_i^{RE} \\ &= \int_0^\infty \left(\prod_{t=1}^T \exp(-a_i^{RE} \lambda_{i,t}^{RE}) \frac{(a_i^{RE} \lambda_{i,t}^{RE})^{n_{i,t}}}{n_{i,t}!} \right) f(a_i^{RE}) da_i^{RE} \end{aligned}$$

We can see that we have to suppose an **additional assumption** : from the first to the second line of development, $f(a_i^{RE} | x_{i,1}, \dots, x_{i,T})$ becomes $f(a_i^{RE})$. **The interpretation of random effects results are tricky.**

Comparative Analysis

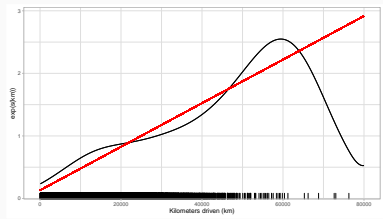


Figure 5: Comparison between the random effect approach and the fixed-effect approach for the median value of the individual parameter

Which Effect Should Be Used in Practice ?

- 1 Fixed effects modeling, even if theoretically better, is **not amenable** to ratemaking.
- 2 On the other hand, the **MVNB** can be used for **predictive rating**, where it can be shown that the predictive distribution of $N_{i,T}$ depends on past values of $\lambda_{i,t}$ and $n_{i,t}$, for $t = 1, \dots, T - 1$.

Take-home points

- 1 **Fixed effects** should be used to understand the **"true" relationship** between covariates and claims experience.
- 2 For ratemaking, fixed effects should be used to compute the **premium surcharge** for each additional kilometer the insureds drive.
- 3 In our case, it represents an increase of $\hat{\alpha}_i \frac{1}{15,000}$ per km, for claim frequency.
 - ▶ Using this approach, insurers will **avoid** the situation where an insured could see a **premium reduction** if, for example, he decides to drive 50,000 km instead of 40,000 km, as we saw with a **basic GAM approach**.
- 4 **Fixed effects** can be used to construct PAYD insurance solely based on kilometers driven for **self-service vehicles**, where drivers' profile cannot be directly used for ratemaking.
- 5 Research is required in this area.

Claim Classification Using Partial Telematics Information

Research question

When has an insurer collected enough information about an insured's driving habits?

General idea

- ▶ Supervised classification with **classic** and **telematics** covariates.
- ▶ Modeling the indicator of one or more claims.
- ▶ Calculation of telematics covariates at different stages of the contract : after 1 month, 2 months, ..., 12 months. Then, comparison of the performance.

Motivations

- ▶ An insurer wishes to keep a minimum of telematic information on its policyholders for reasons of :
 - Confidentiality
 - Data storage
- ▶ But still wants to take advantage of this information, for instance, to avoid adverse selection.

Trip data

Extract from the trip database

| VIN | Trip ID | Starting time | Arrival time | Distance | Maximum speed |
|-----|---------|---------------------|---------------------|----------|---------------|
| A | 1 | 2016-04-09 15:23:55 | 2016-04-09 15:40:05 | 10.0 | 72 |
| A | 2 | 2016-04-09 17:49:33 | 2016-04-09 17:57:44 | 4.5 | 68 |
| ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ |
| A | 3312 | 2019-02-11 18:33:07 | 2019-02-11 18:54:10 | 9.6 | 65 |
| B | 1 | 2016-04-04 06:54:00 | 2016-04-04 07:11:37 | 14.0 | 112 |
| B | 2 | 2016-04-04 15:20:19 | 2016-04-04 15:34:38 | 13.5 | 124 |
| ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ |
| B | 2505 | 2019-02-11 17:46:47 | 2019-02-11 18:19:22 | 39.0 | 130 |
| C | 1 | 2016-01-16 15:41:59 | 2016-01-16 15:51:35 | 3.3 | 65 |
| ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ |

- ▶ These are the only telematic data we have. All telematic covariates are derived from these **4 measurements**.

Contract data

Extract from the contract database

| VIN | Contract start date | Contract end date | Classic covariate #1 | ... | Claim(s) indicator |
|-----|---------------------|-------------------|----------------------|-----|--------------------|
| A | 2015-01-09 | 2016-01-09 | F | ... | 0 |
| A | 2016-01-09 | 2017-01-09 | F | ... | 1 |
| A | 2017-01-09 | 2018-01-09 | F | ... | 0 |
| B | 2015-12-14 | 2016-12-14 | M | ... | 0 |
| B | 2016-12-14 | 2017-12-14 | M | ... | 0 |
| C | 2015-04-26 | 2016-04-26 | F | ... | 1 |
| C | 2016-04-26 | 2017-04-26 | F | ... | 0 |
| C | 2017-04-26 | 2018-04-26 | F | ... | 0 |
| ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ |

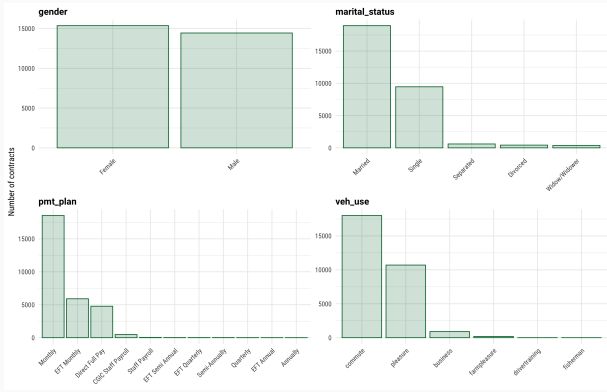
- ▶ Linking of the 2 datasets on the basis of the VIN and the start/end dates of the contract.
- ▶ Expansion of the contract database with 14 telematics variables calculated using the trip dataset.
- ▶ We only consider one-year contracts.

Creation of the classification datasets

- 1 For each row (contract) in the contract dataset, associate the right trips from the trip dataset.
- 2 Compute the 14 telematics variables with different levels of information : 1 months of telematics data, 2 months, 3 months, etc. until 12 months.

- ▶ We end up with **13 tables** 29799×25 (10 classic covariates, 14 telematic covariates and 1 target).
 - 1 table with only classic covariates
 - 12 tables with classic and telematic covariates, respectively calculated with 1, 2, ..., 12 months of data.
- ▶ We keep 70% of the rows (contracts) for the **training** set and 30% for the **test** set.

Classic covariates – Categorical



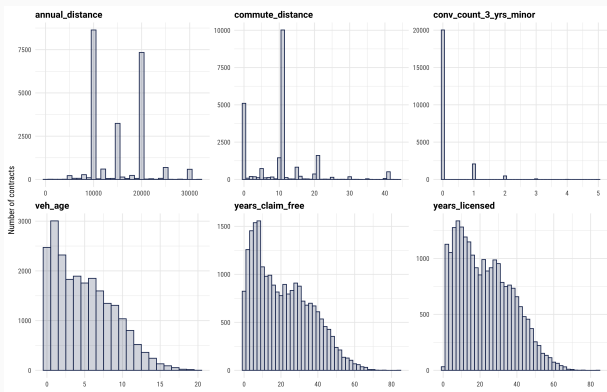
Preprocessing :

Lump rare categories → target encode → normalize → Yeo-Johnson transform

Longitudinal analysis of distance traveled
○○○○○○○○○○○○○○○○○○○○○○○○

Claim Classification Using Partial Telematics Information
○○○○●○○○○○○○○

Classic covariates – Numeric



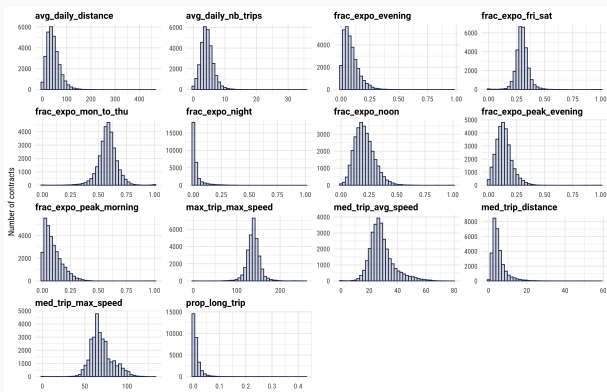
Preprocessing :

Normalize → Yeo-Johnson transform

Longitudinal analysis of distance traveled
○○○○○○○○○○○○○○○○○○○○○○○○

Claim Classification Using Partial Telematics Information
○○○○●○○○○○○○○

Telematic covariates



- ▶ Here, the distributions of the covariates calculated with full information (12 months) are shown.
- ▶ Preprocessing is the same as classic numeric covariates.

Longitudinal analysis of distance traveled
○○○○○○○○○○○○○○○○○○○○○○○○

Claim Classification Using Partial Telematics Information
○○○○●○○○○○○○○

Choice of the classification algorithm

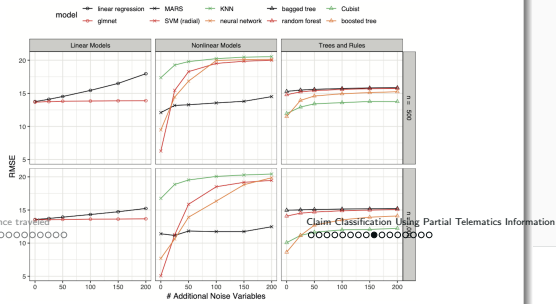
- ▶ Need to be an **off-the-shelf** algorithm.
- ▶ Need to be robust to **redundant** or **unnecessary** predictors.

Question

- ▶ Which classification algorithms do you think are good candidates ?

Answer

- ▶ 2 good candidates are **penalized logistic regression** and **random forest**.
 - Easy to tune hyperparameters
 - Robust to "noise" predictors



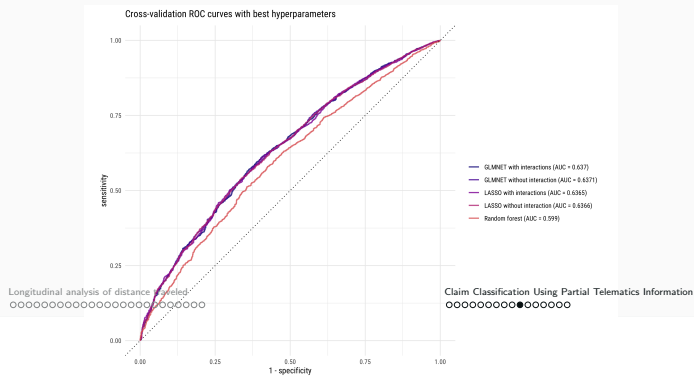
*Figure taken from *Feature Engineering and Selection: A Practical Approach for Predictive Models*, by Max Kuhn and Kjell Johnson.

Choice of the classification algorithm

- ▶ In order to choose the classification algorithm, I use the database with complete information and compare the performance of an **elastic net** logistic regression, a **LASSO** logistic regression and a **random forest**.

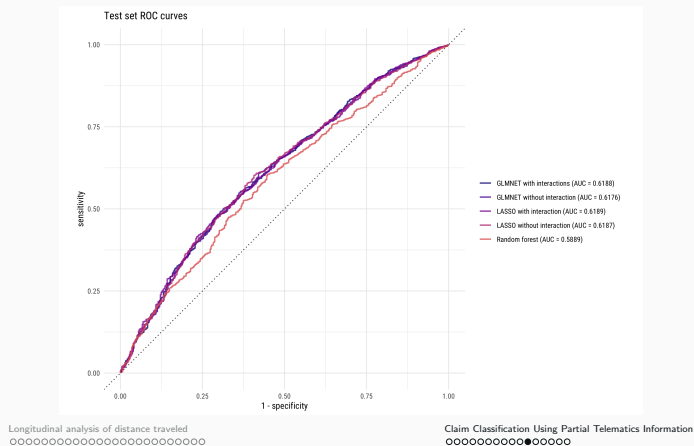
Question

- ▶ Which classification algorithms do you think will perform the best between LASSO, elastic net and random forest ?



Choice of the classification algorithm

- ▶ Same plot, but on the test set instead of cross-validation.



A glimpse at LASSO logistic regression

Loss function

$$L(\beta, \mathbf{y}) = -\frac{1}{n} \sum_{i=1}^n \{y_i \ln(\rho_i) + (1 - y_i) \ln(1 - \rho_i)\} + \lambda \sum_{j=1}^p |\beta_j|, \quad \text{où } \rho_i = \frac{1}{1 + e^{-\mathbf{x}_i^T \beta}}$$

Estimation

- We find the β coefficients that minimize the loss function, which is equivalent to minimizing the negative of the log-likelihood with a constraint on the sum of the absolute values of the coefficients :

$$\hat{\beta}^{\text{LASSO}} = \arg \min_{\beta} \left\{ -\frac{1}{n} \sum_{i=1}^n y_i \ln(\rho_i) + (1 - y_i) \ln(1 - \rho_i) \right\} \quad \text{s.c.} \quad \sum_{j=1}^p |\beta_j| \leq s$$

Prediction

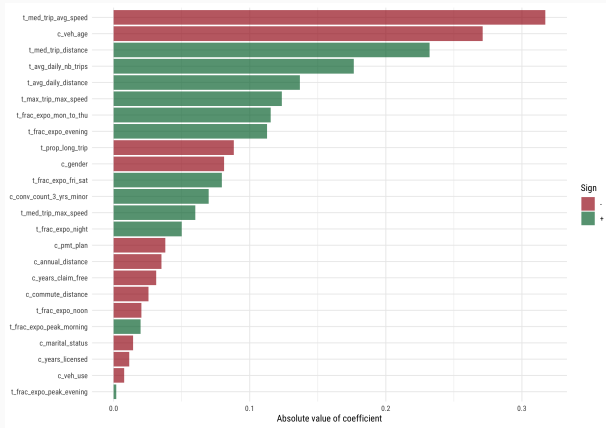
- Same prediction formula as a non-penalized logistic regression, but using LASSO coefficients $\hat{\beta}^{\text{LASSO}}$:

$$\hat{y}_i = \frac{1}{1 + e^{-\mathbf{x}_i^T \hat{\beta}^{\text{LASSO}}}}$$

Longitudinal analysis of distance traveled
○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○

Claim Classification Using Partial Telematics Information
○○○○○○○○○○○○○○○○○○●○○○○

LASSO logistic regression coefficients



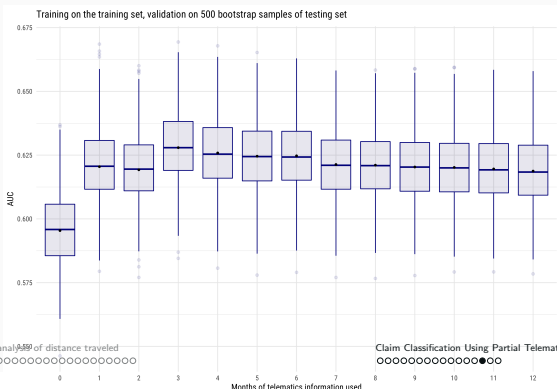
Longitudinal analysis of distance traveled
○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○

Claim Classification Using Partial Telematics Information
○○○○○○○○○○○○○○○○○○●○○○○

Logistic LASSO performance on the 13 datasets

Question

- I'm about to show you the performance of the model on the 13 datasets. After how many months do you think that telematics information no longer significantly improves performance ?

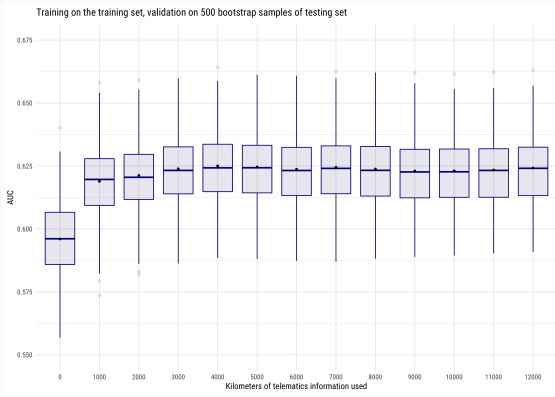


Longitudinal analysis of distance traveled
○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○

Claim Classification Using Partial Telematics Information
○○○○○○○○○○○○○○○○○○●○○○○

- The AUC has improved significantly with the 4-measure trip summaries!
- Telematics information becomes redundant after about 3 months.

Logistic LASSO performance on the 13 datasets – km



► Telematics information becomes redundant after about 4000 km.

Longitudinal analysis of distance traveled
 ○○○○○○○○○○○○○○○○○○○○○

Claim Classification Using Partial Telematics Information
 ○○○○○○○○○○○○○○○○○●

Future considerations

Conclusion

► We found out that telematics information we have at our disposal becomes redundant after about 3 months or 4000 km, at least in the collision claim classification framework.

Integration of contracts of less than one year

► Here, only one-year contracts were used.

Test on other insurance coverage

► In our analysis, only collision type coverages are considered.
 ► Do we come to the same conclusion if we use, for instance, comprehensive coverage claims (theft, hail, etc.)?

Longitudinal analysis of distance traveled
 ○○○○○○○○○○○○○○○○○○○○○

Claim Classification Using Partial Telematics Information
 ○○○○○○○○○○○○○○○○○●