



# AKUR8

DAV Herbsttagung

---

November 15th 2022



**Jan Kütke**

**Aktuar (DAV) / Actuarial Data Scientist**

## **Biography**

Jan ist ein Aktuar (DAV) aus Siegen und arbeitet bei Akur8 als Actuarial Data Scientist, um Versicherern zu helfen, die Potenziale der Pricing-Methoden des 21. Jahrhunderts nutzbar zu machen.

Zuvor hat er drei Jahre bei einer weltweit tätigen Aktuariellen Beratung gearbeitet. Er hat einen Master of Mathematics an der Universität Bonn abgeschlossen und wohnt in Köln.

Daneben ist er ein begeisterter Leser der Werke von Anna Seghers und Dietmar Dath.

# Actuarial Modeling

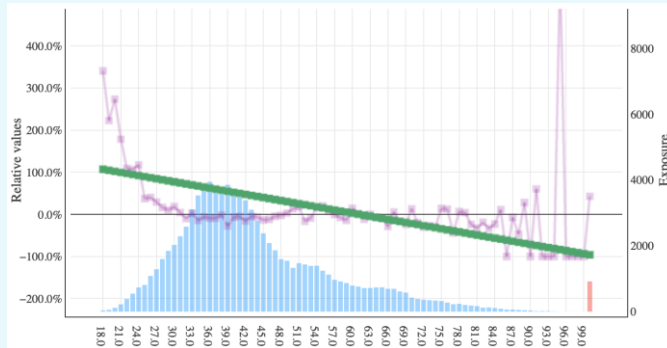
With GLMs, Credibility and Penalized Regressions

# Actuarial Modeling: Capturing Non-Linearities

## What GLMs Offer...

Generalized Linear Models (“GLMs”) are, by definition, linear.

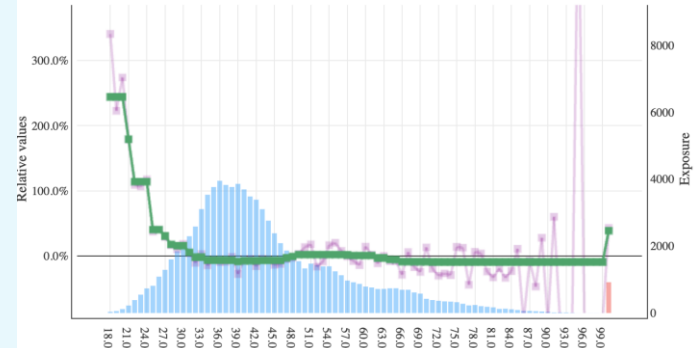
They are easy to fit (as only one parameter has to be found for every variable).



## ...What we want

We want to capture the non-linear relations between the explanatory and predicted variables.

They are hard to fit because, for every variable, a large number of parameters has to be found.



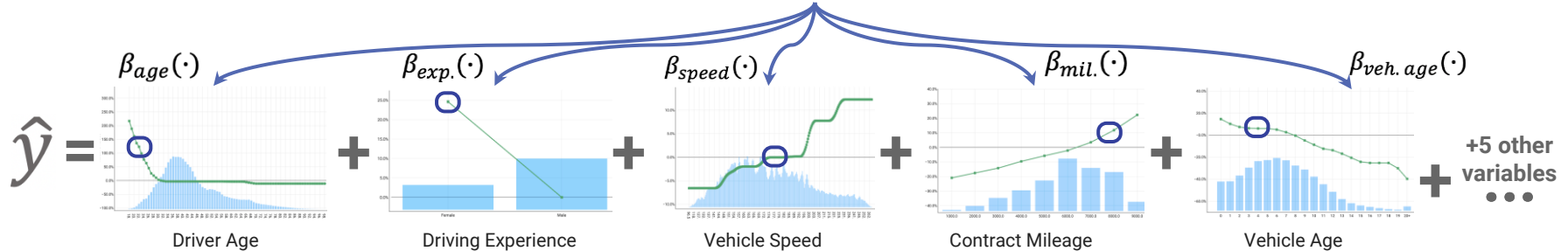
# Additives Models: Easy to visualize

Additive models are a great balance between predictive power and adverse-selection management.

For this reason, they are **currently widely used in the actuarial community**.

Additive models can be visualized as rating tables, but most remarkably, the **human visualization is convenient** for model review and modification as it displays one function per variable.

$$\hat{y}(X) = g^{-1} \left( \sum_d \beta_d(X_d) \right)$$

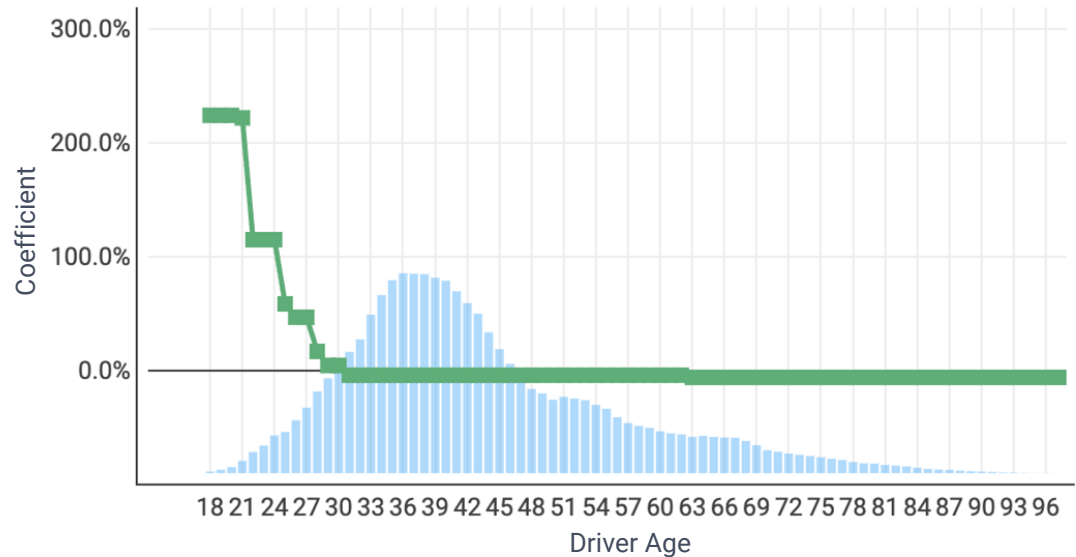


# Representing Additive Models with Tables

Driver Age	Coefficient
18	224%
19	224%
20	224%
21	222%
22	115%
23	115%
24	115%
25	59%
26	47%
27	47%
28	17%
29	4%
30	4%
31	-4%
32	-4%
33	-4%
34	-4%
35	-4%
36	-4%



The function  $\beta_{age}(\cdot)$  is represented by all its values, for all values of the driver age..



# GLM and GAM equivalence

Linear Models

$$\hat{y}(X) = g^{-1} \left( \sum_{i,j} \beta_{i,j} \times I_{X_i=j} \right)$$

Variables Transformations

Driver Age=16
Driver Age=17
Driver Age=18
Driver Age=19
Driver Age=20
Driver Age=21
Driver Age=22
Driver Age=23

Non-Linear Models

$$\hat{y}(X) = g^{-1} \left( \sum_j \beta_j(X_j) \right)$$

$$I_{X_{age}=j} = \begin{cases} 1 & \text{if } X_{age} = j \\ 0 & \text{if } X_{age} \neq j \end{cases}$$

**GLMs and GAMs are equivalent:** coefficients are built for different values of the explanatory variables.

# Controlling the Transformations

# Creating a GLM to capture non-linear relationships

All regression models are built around the same main principle: to **maximize the probability of the observations**.

$$\beta^* = \text{Argmax}_{\beta} p(y|\hat{y})$$

However, **maximizing the likelihood** on hundreds of parameters would lead to overfitting, which needs to be controlled.

Two main approaches are used by the actuarial community:

Manage the number of parameters by carefully **selecting which transformations**

**are used:**

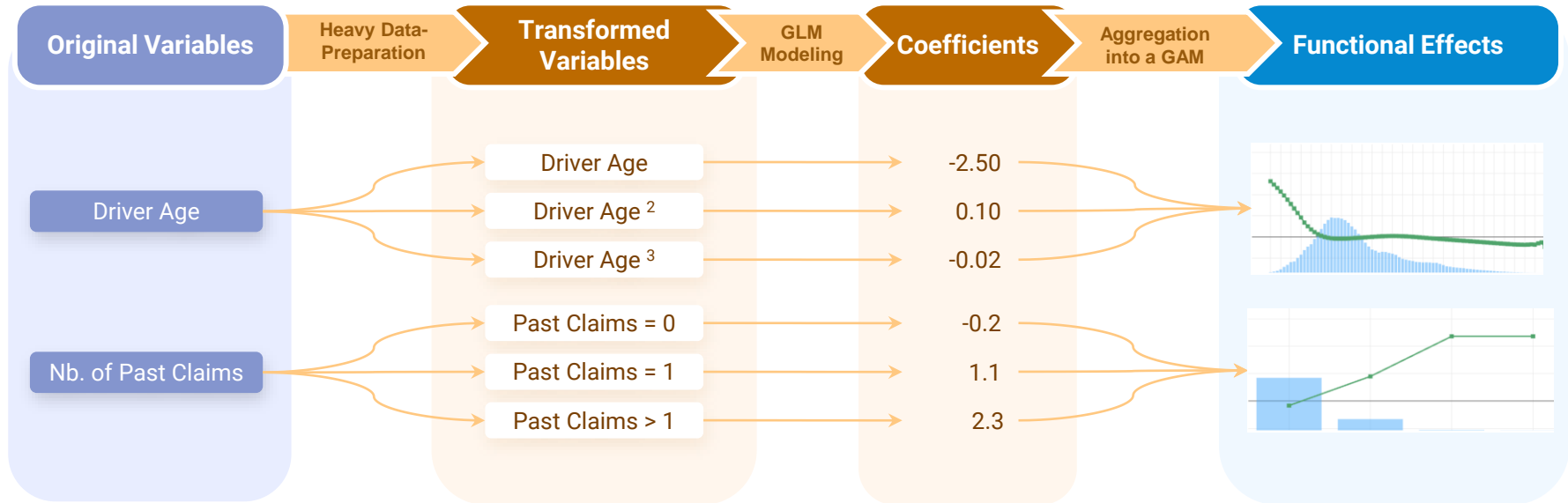
- Polynomials
- Groupings
- ...

Integrate priors on the coefficients into the model creation:

- The priors will be directly included into the likelihood optimization.
- They will reduce the complexity of the models created.

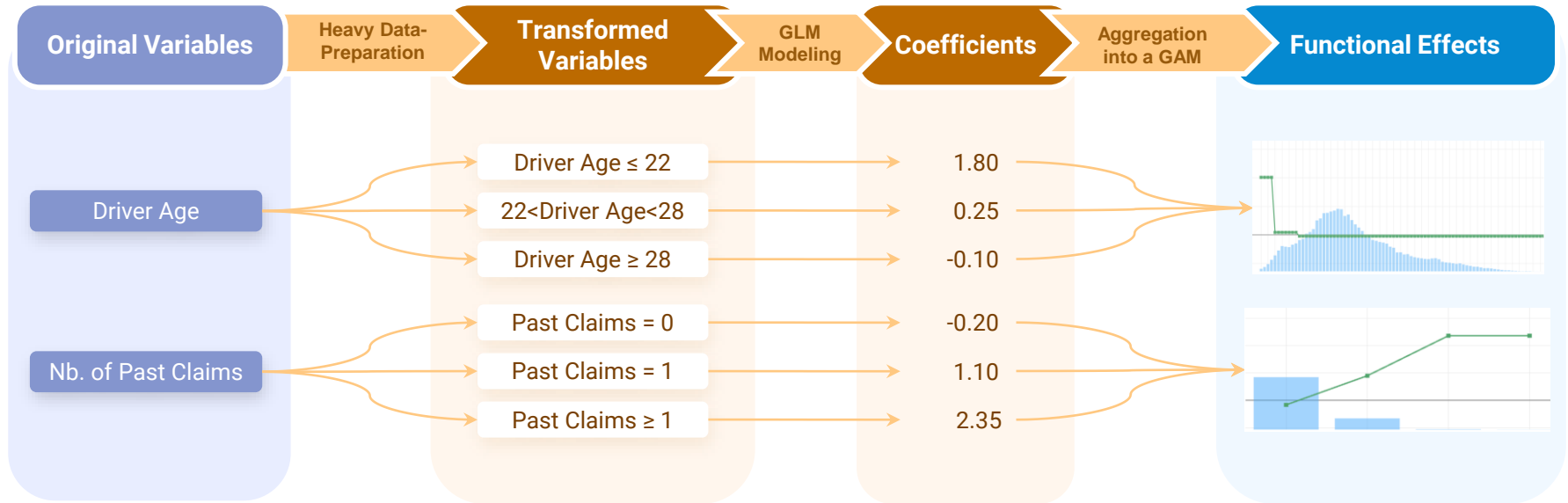
# Modeling with variable transformations

Building the model  $\leftrightarrow$  Selecting the transformation



# Modeling with variable transformations

Different transformations lead to a different model



# Leveraging Credibility

# Creating a GLM to capture non-linear relationships

All regression models are built around the same main principle:

$$\beta^* = \text{Argmax}_{\beta} \text{Likelihood}(y, x, \beta)$$

However, **maximizing the likelihood** on hundreds of parameters would lead to overfitting, which needs to be controlled.

Two main approaches are used by the actuarial community:

Manage the number of parameters by carefully **selecting which transformations are used**:

- Polynomials
- Groupings
- ...

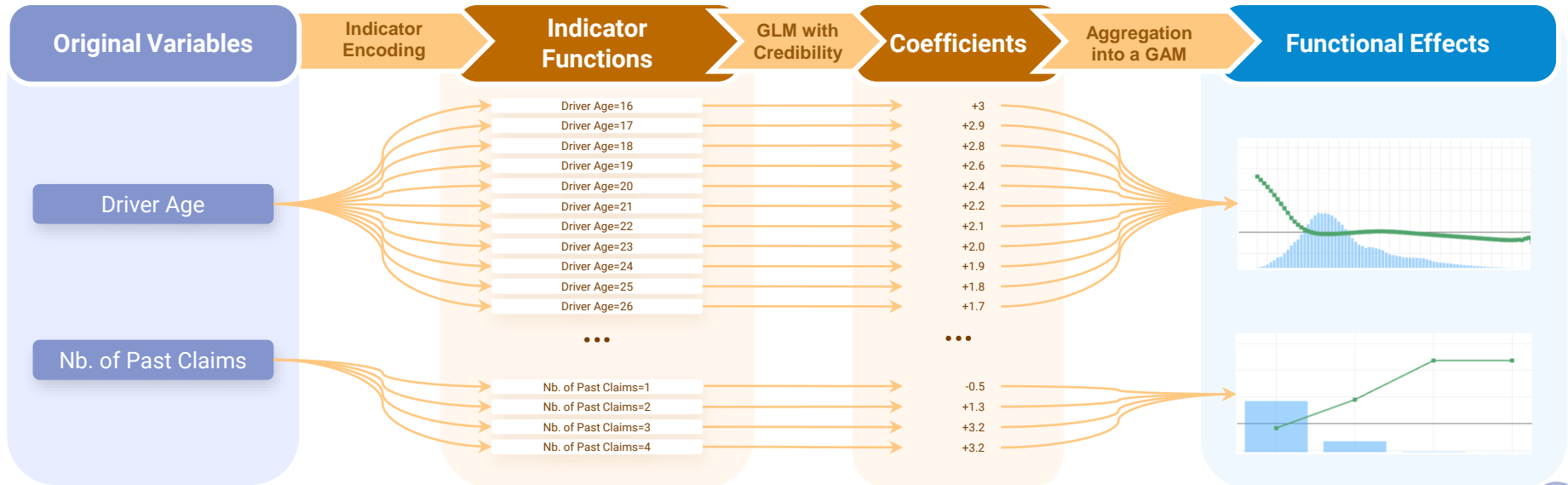
Integrate priors on the coefficients into the model creation:

- The priors will be directly included into the likelihood optimization.
- They will reduce the complexity of the models created.

# Automatic Modeling with Credibility

In order to remove the heavy and time-consuming data-preparation step, a **large number of indicator functions** are created - these functions equal one if a variable equals a given value, zero otherwise.

Then a model **fitted leveraging credibility** ensures the coherence between the different coefficients created.



# Automatic Modeling with Credibility

In order to remove the heavy and time-consuming data-preparation step, a **large number of indicator functions** are created - these functions equal one if a variable equals a given value, zero otherwise.

Then a model **fitted leveraging credibility** to ensure the coherence between the different coefficients created.



## Quick Reminder... What is credibility?



Credibility, simply put, is the weighting together of different estimates to come up with a combined estimate.

*Foundations of Casualty Actuarial Science*

Buhlmann credibility is the best-known approach. It is equivalent to a simple **Bayesian** framework, where a prior “knowledge” based on a model is updated based on observations.

Usually (after equations involving conditional probabilities), the output of a credibility approach is that the model predictions are a **weighted average** between the observations and the initial assumption.

The weight will depend on:

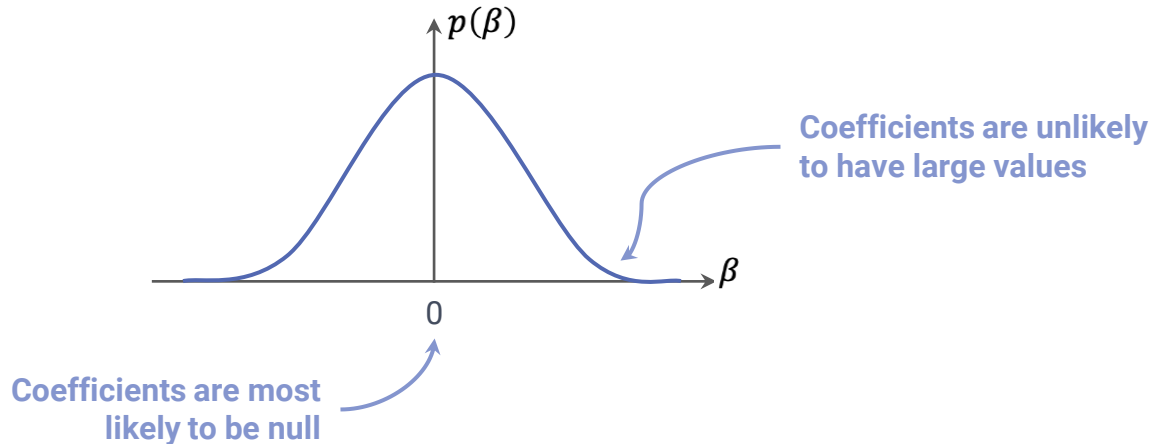
- the **quantity of data** (the larger the data, the higher the weight)
- the **strength of the prior** assumptions (a very reliable assumption with small variance will have a large weight).

# Prior and Credibility

A credibility framework is defined by the prior assumptions the modeller has on his model. These **assumptions represent a prior probability** distribution for the models coefficients.

For instance, **“simpler” models are usually assumed to be “more likely”**.

Classic prior assumptions can be: “The coefficients follow a Gaussian distribution, centered on 0”



# Prior and Penalized Regression

---

This prior is formalized as a distribution of probability for the coefficients:  $p_{prior}(\beta)$

The **Maximum of Likelihood approach** directly integrates the prior:

$$\beta^* = \mathit{Argmax}_{\beta} p(y|\hat{y}(X)) \times p_{prior}(\beta)$$

Taking the log, we get the Maximum of Likelihood problem:  $\beta^* = \mathit{Argmax}_{\beta} LL(x, y, \beta) + \log(p_{prior}(\beta))$

Or equivalently the Minimization of Error problem (or **Penalized Regression**):

$$\beta^* = \mathit{Argmin}_{\beta} \mathit{Errors}(y, \hat{y}) - \mathit{Penalty}(\beta)$$

# Prior $\Leftrightarrow$ Penalized Regressions

Some examples in the Linear Regression case

**Prior assumptions are at the center of penalized-regression methods** used to control high-dimensional or correlated data, such as Lasso or Ridge Regression. Controlling the distribution (through the  $\lambda$  parameter) allows for controlling the overfitting of the models.

Gaussian Hypothesis



**Prior:** Coefficients follow a Normal distribution  $N(0, 1/2\lambda)$ :



**Coefficients Distribution:**

$$p(\beta) \sim e^{-\lambda \beta^2}$$



**Log-Likelihood (incl. prior)**

$$LL(x, y, \beta) - \lambda \beta^2$$



Ridge Regression

Laplace Hypothesis



**Prior:** Coefficients follow a Laplace distribution  $L(0, 1/\lambda)$ :



**Coefficients Distribution:**

$$p(\beta) \sim e^{-\lambda |\beta|}$$



**Log-Likelihood (incl. prior)**

$$LL(x, y, \beta) - \lambda |\beta|$$



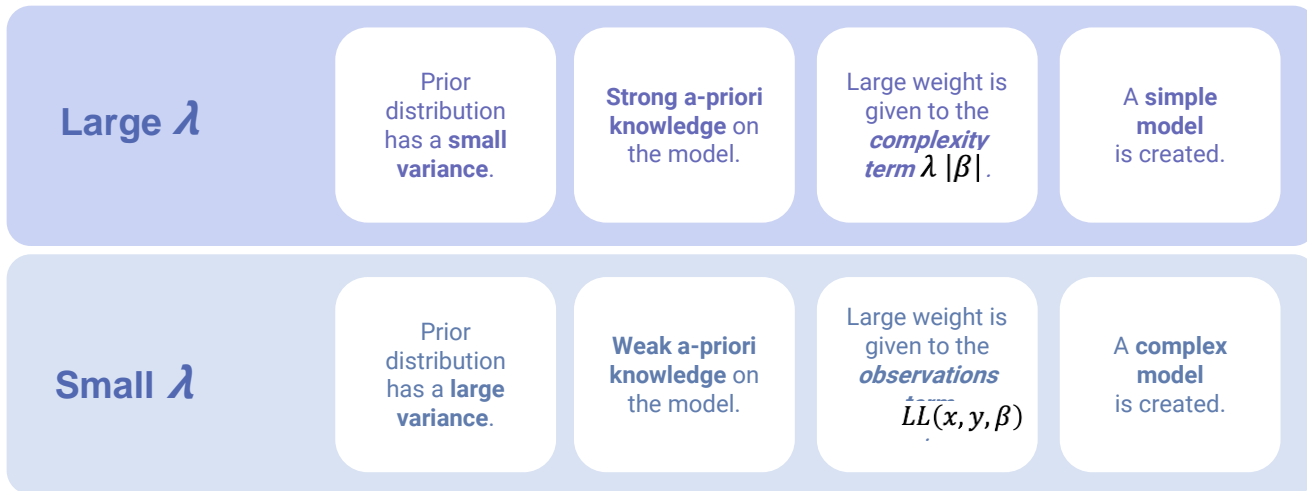
Lasso Regression

# Controlling the Prior distribution

With Lasso

In a Lasso, the prior follow a distribution  $P(\beta) \sim e^{-\lambda |\beta|}$   $2/\lambda^2$   
variance

The coefficients should **maximize**:  $LL(x, y, \beta) - \lambda |\beta|$

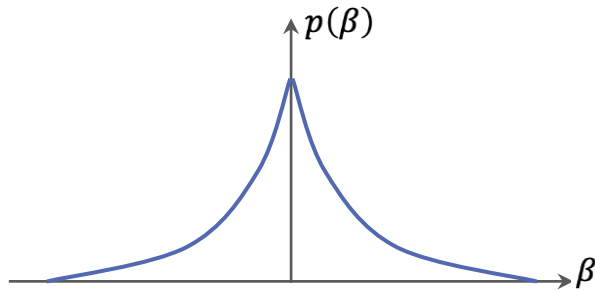


This approach has the **same roots as credibility** (it is based on a Bayesian framework) and the **same consequences** (it “shrinks” the coefficients towards our prior most likely estimate).

# Lasso and Hypothesis testing

**Lasso is especially popular as it is a good tool for variable selections:** models created with the Lasso framework are sparse - all the non-relevant coefficients equal zero.

The Laplace distribution that underlies the Lasso has a maximum at zero:



When used on binary explanatory variables, it is also equivalent to **hypothesis testing**:

**Null Hypothesis:**  $\beta = 0$  : “The coefficient is not significantly different from zero.”

- If the null hypothesis is **not rejected**, the coefficient value is zero.
- If the null hypothesis is **rejected**, the coefficient has a non-zero value.

# Lasso can be used in actuarial modelling...

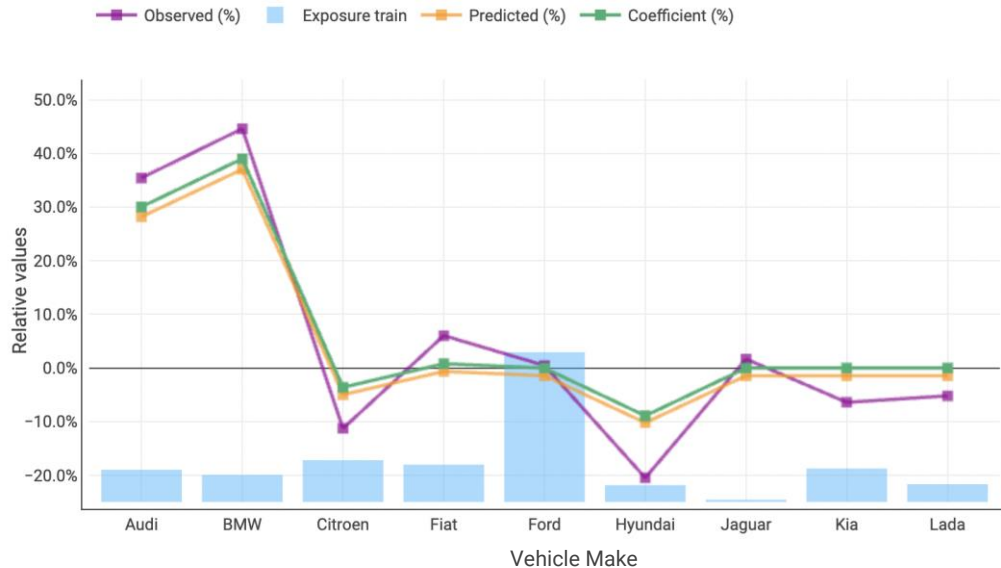
Lasso can be used to capture the signal on **categorical variables**.

Coefficients are created for each level of the data:

$$\hat{y}(X) = g^{-1} \left( \sum_{i,j} \beta_{i,j} \times I_{X_i=j} \right)$$

The result is coherent with a **credibility approach**: predictions are between their “pure GLM” values and the grand-mean of the observations.

Non-significant levels are grouped, with null coefficients.



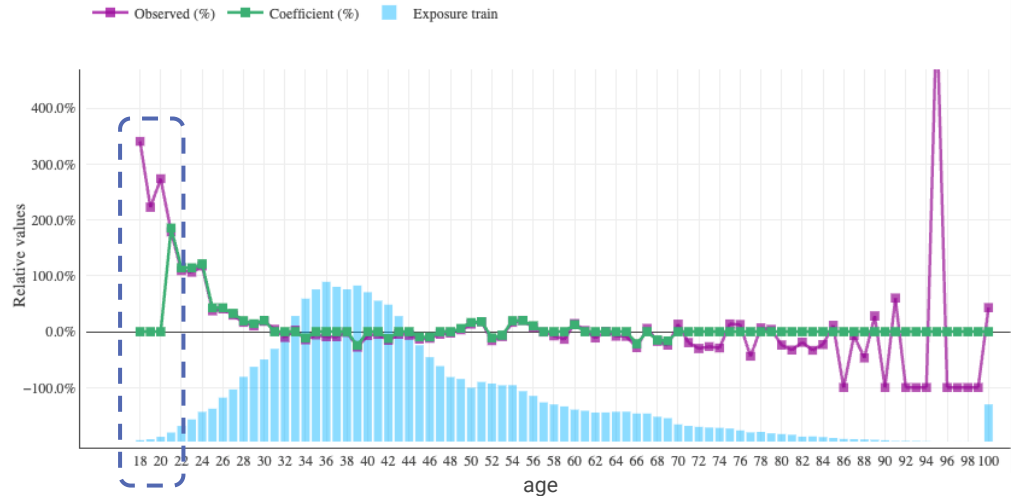
# ...but Lasso does not capture continuous effects!

While it is very powerful and well documented, the **Lasso can't be directly applied** to indicator-representation on the data to create a non-linear model:

$$\hat{y}(X) = g^{-1} \left( \sum_{i,j} \beta_{i,j} \times I_{X_i=j} \right)$$

All non-significant coefficients would be grouped at zero, which makes no sense.

A key piece of information: **the order of the levels would be lost in the process.**



No information in the data =  
The most likely coefficients  
are at zero.

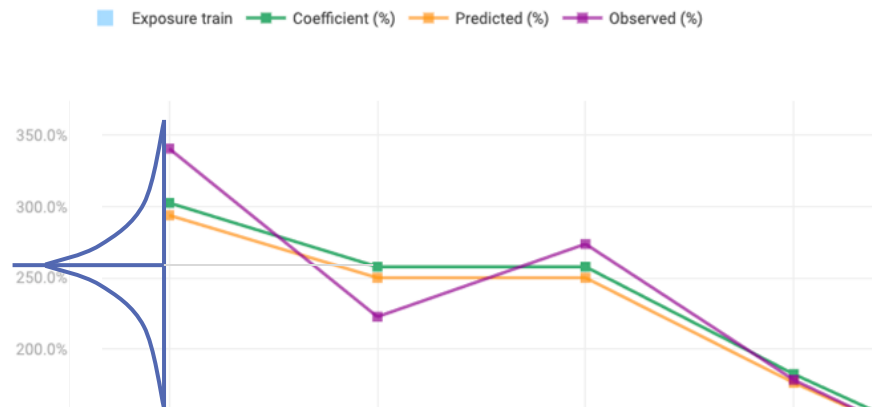
# Credibility on Ordered Variables

# Creating new Priors and Penalties

New priors have to be considered to take into account the structure of the models created.

In particular, for ordinal variables, two consecutive coefficients should:

- be more likely to be close than far apart if they are significantly different.
- or have the same coefficients if they are not significantly different...



This concept **generalizes the Lasso penalty to continuous function**, providing the high level of flexibility and stability necessary to create GAM models.

# Creating new Priors and Penalties

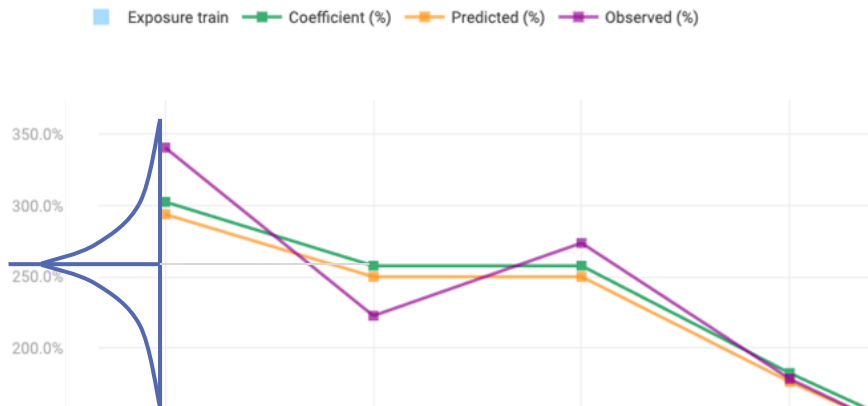
This means that the **derivative of the coefficient function**  $\beta'(X)$  follows a **Laplace distribution**:

As the values of the coefficients are discrete, the derivative can be written as:

$$p(\beta) \propto e^{-\lambda |\beta_i - \beta_{i+1}|}$$

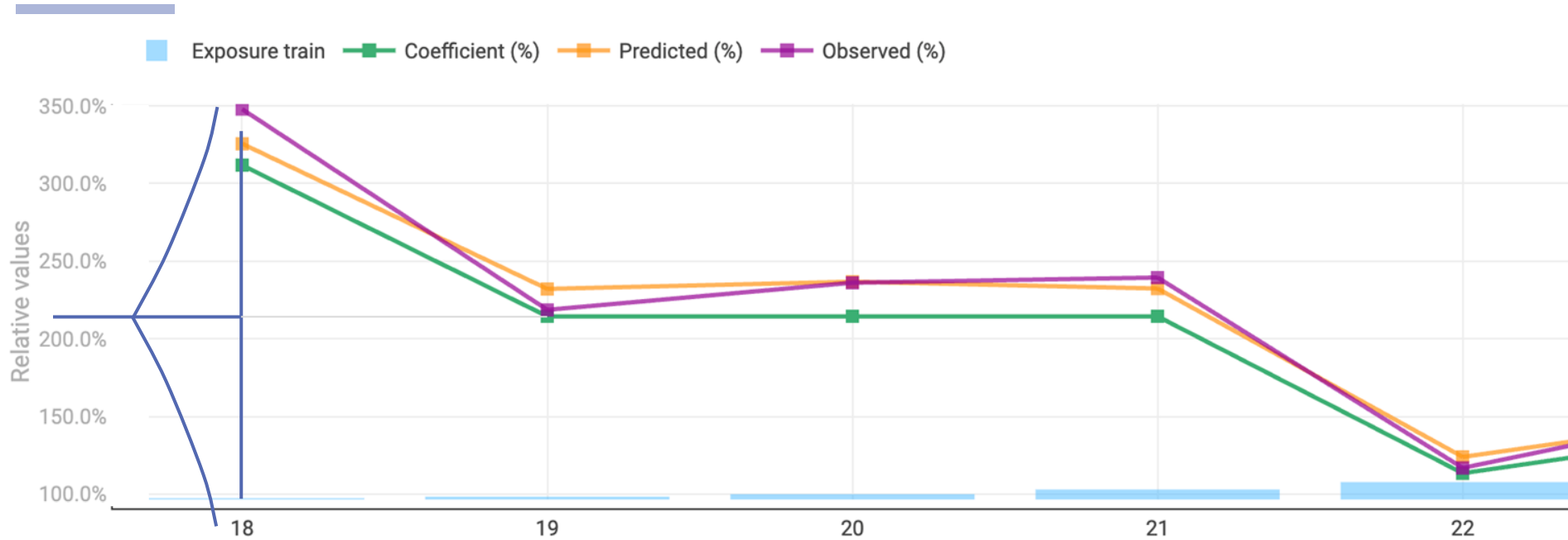
This distribution of probability is used as a **prior when maximizing the likelihood** to fit a model:

$$\beta^* = \text{Argmax}_{\beta} LL(x, y, \beta) - \lambda |\beta_i - \beta_{i+1}|$$



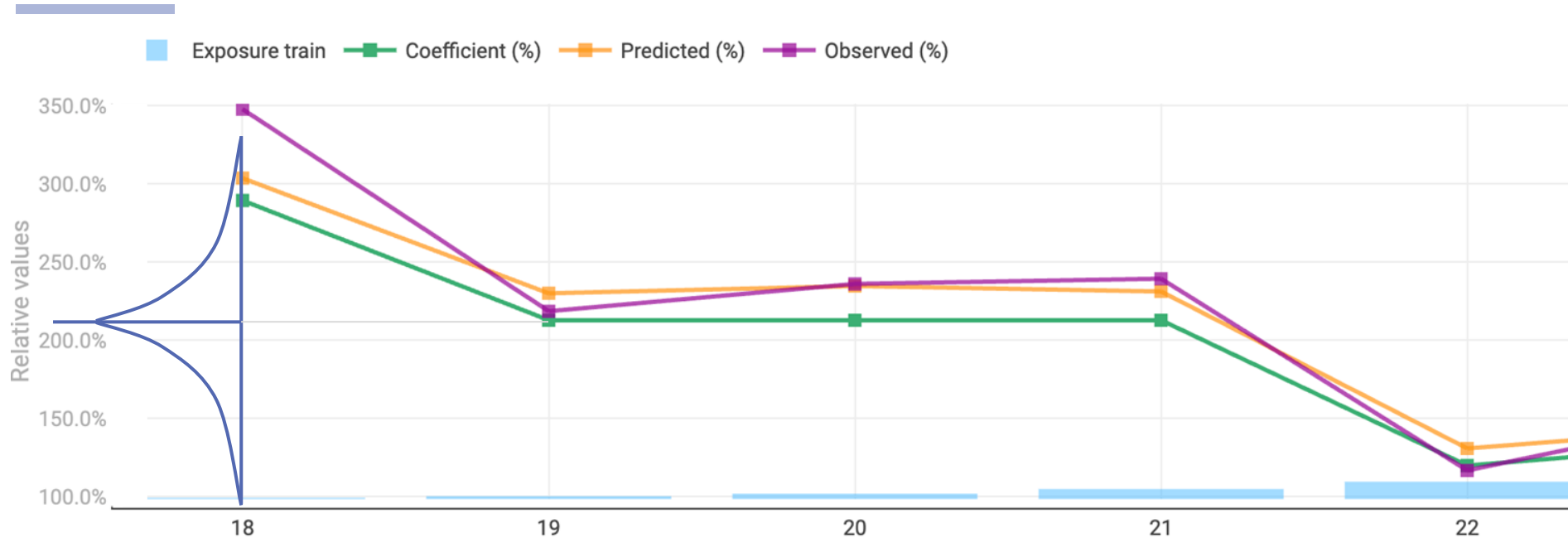
# Weak Prior $\Leftrightarrow$ Strong reliance on the observation

The prior has a very limited impact on the final model



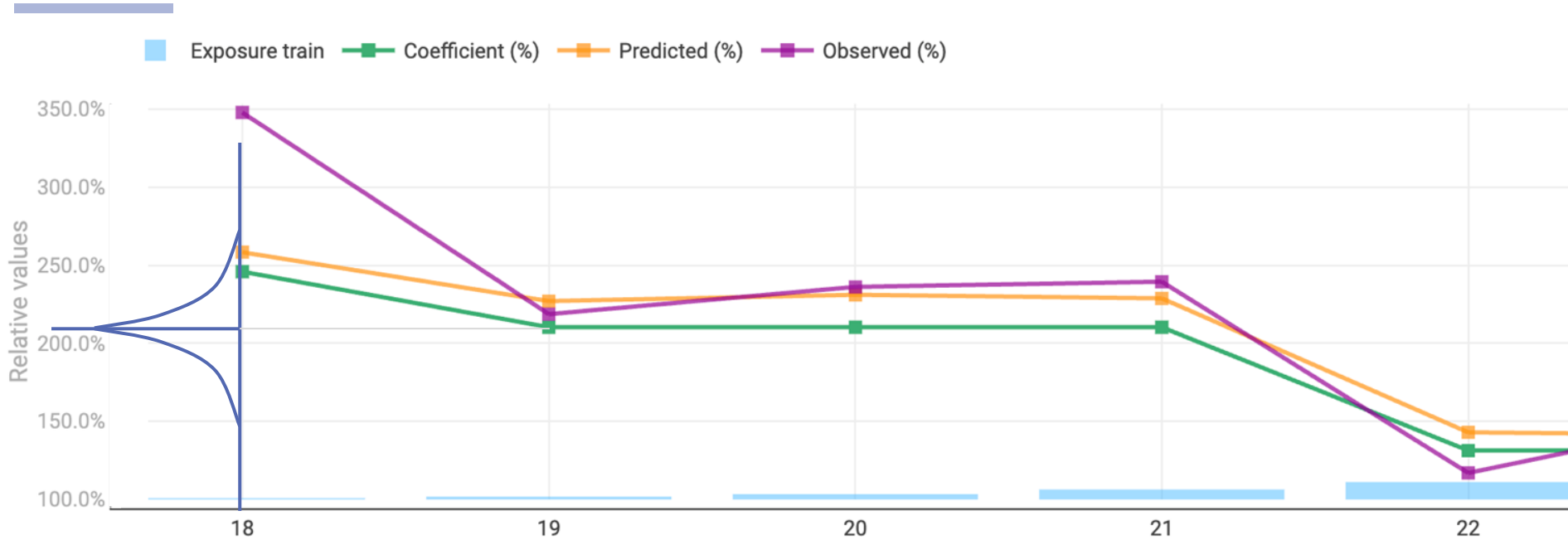
# Stronger Prior $\Leftrightarrow$ Weaker reliance on the observation

The final model is an average between the most likely coefficients according to the prior and the observations



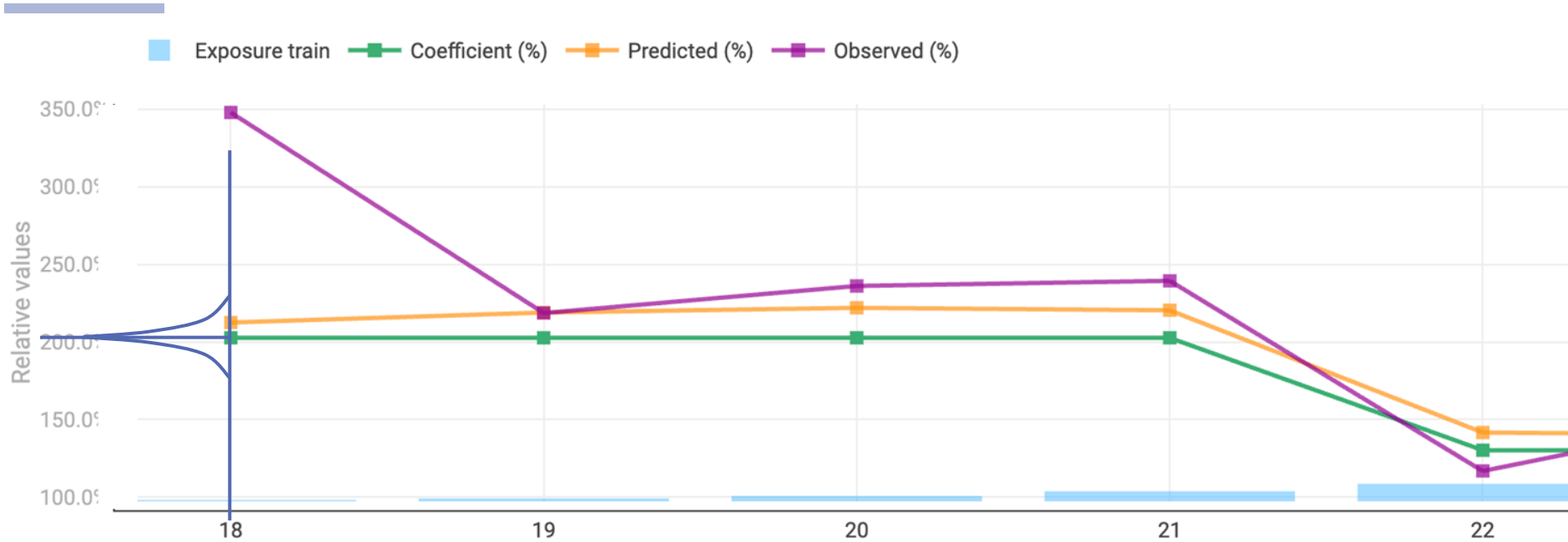
# Strong Prior $\Leftrightarrow$ Very weak reliance on the observation

The weight of the observation in the model is weaker than the priors



# Very Strong Prior $\Leftrightarrow$ Full reliance on the prior

The observations can't disprove such a strong prior - more data would be needed



This is equivalent to failing a significant test against the null hypothesis: “the first two coefficients are equal”.

A stronger effect - or more exposure - would be necessary to disprove it, and split the coefficients.

# Controlling the Prior distribution

To create smooth models

The prior follow a distribution  $p(\beta) \propto e^{-\lambda |\beta_i - \beta_{i+1}|}$  of variance  $2/\lambda^2$

The coefficients should **maximize**:  $LL(x, y, \beta) - \lambda |\beta_i - \beta_{i+1}|$

Large  $\lambda$

Prior distribution has a **small variance**.

**Strong a-priori knowledge** on the model.

Large weight is given to the ***smoothness term***.

A **smooth model** is created.

Small  $\lambda$

Prior distribution has a **large variance**.

**Weak a-priori knowledge** on the model.

Large weight is given to the ***observations term***.

A **noisy model** is created.

# Like for a Lasso, this is equivalent to a test!

The behavior is similar to a hypothesis-testing approach:

**A priori, we suppose the null-hypothesis:**  $\beta_{i+1} - \beta_i = 0$

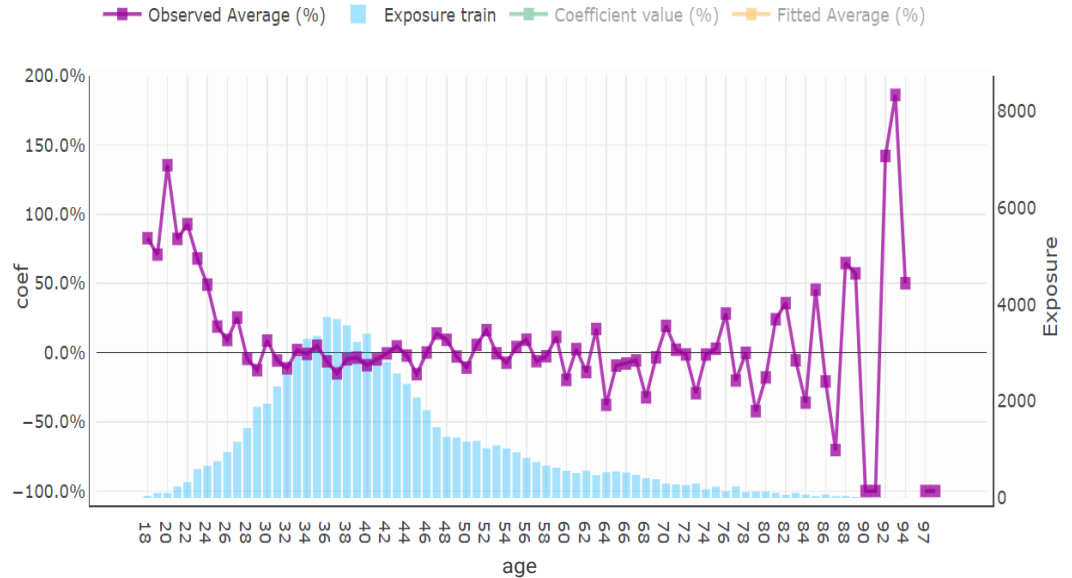
This null hypothesis is tested with the data, and potentially rejected.

This null hypothesis is equivalent to:  $\beta_{i,j} = \beta_{i+1,j}$

- If it is not rejected by the data, then the coefficients function is locally constant.
- If it is rejected by the data, then the coefficients function is not constant.

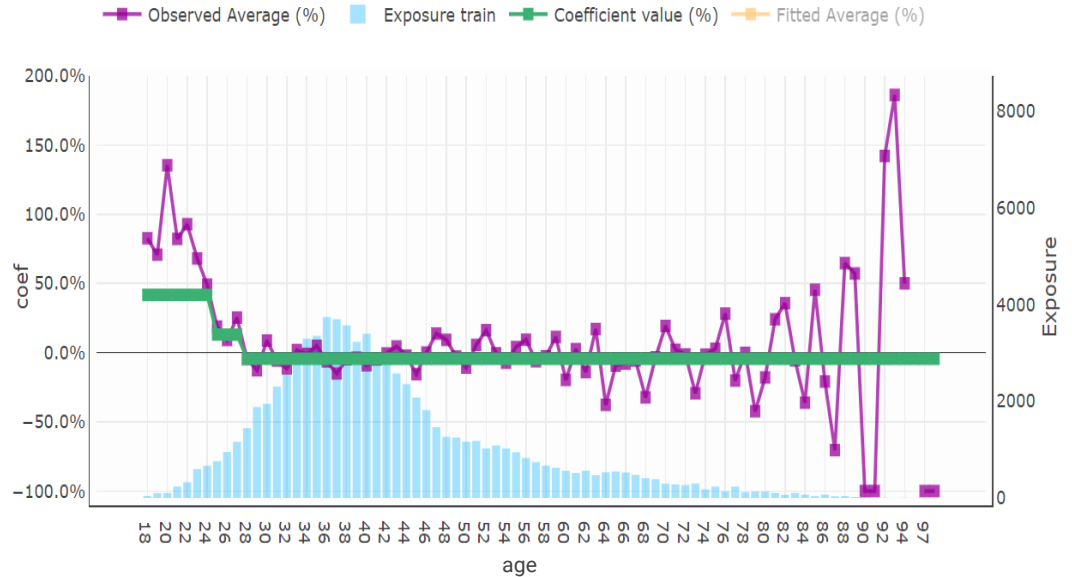
# Leveraging the prior on a full model scale

Data used to create the models are **naturally noisy.**



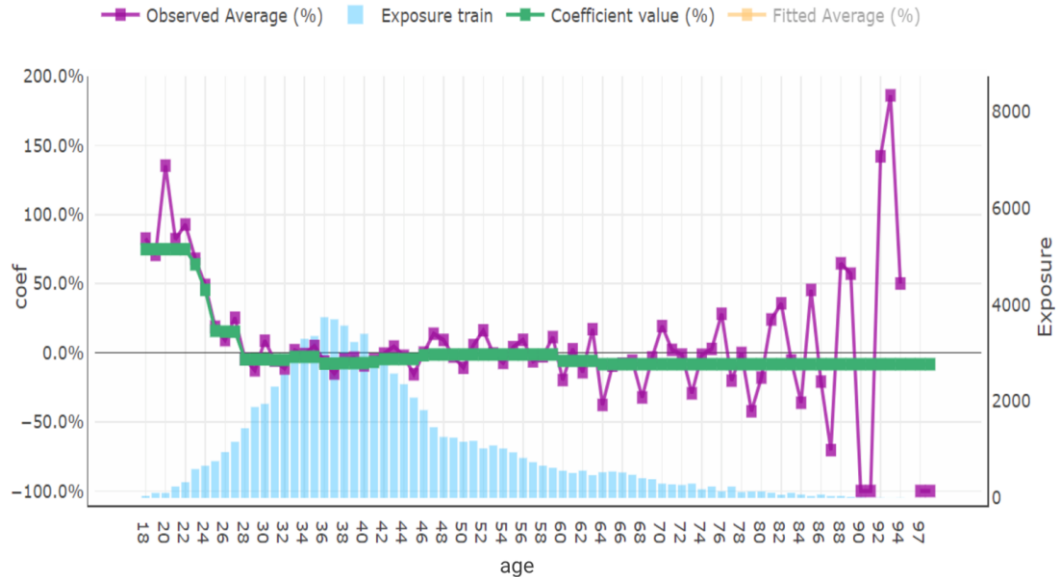
# Leveraging the prior on a full model scale

A very **strong prior**  
(with a small variance)  
leads to **robust models**.



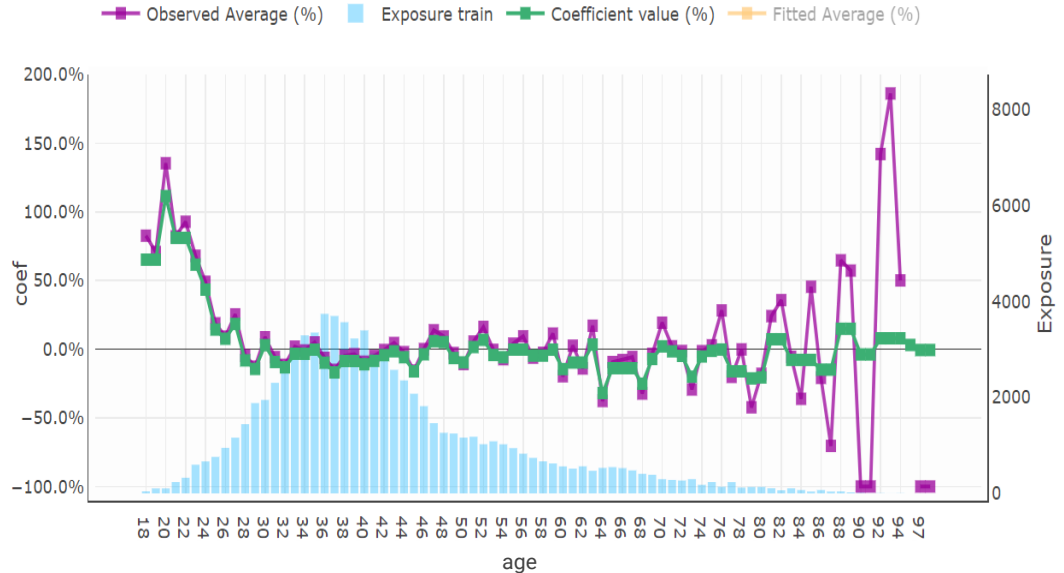
# Leveraging the prior on a full model scale

A more **balanced prior**  
(with a medium variance)  
leads to more **sensitive models**.



# Leveraging the prior on a full model scale

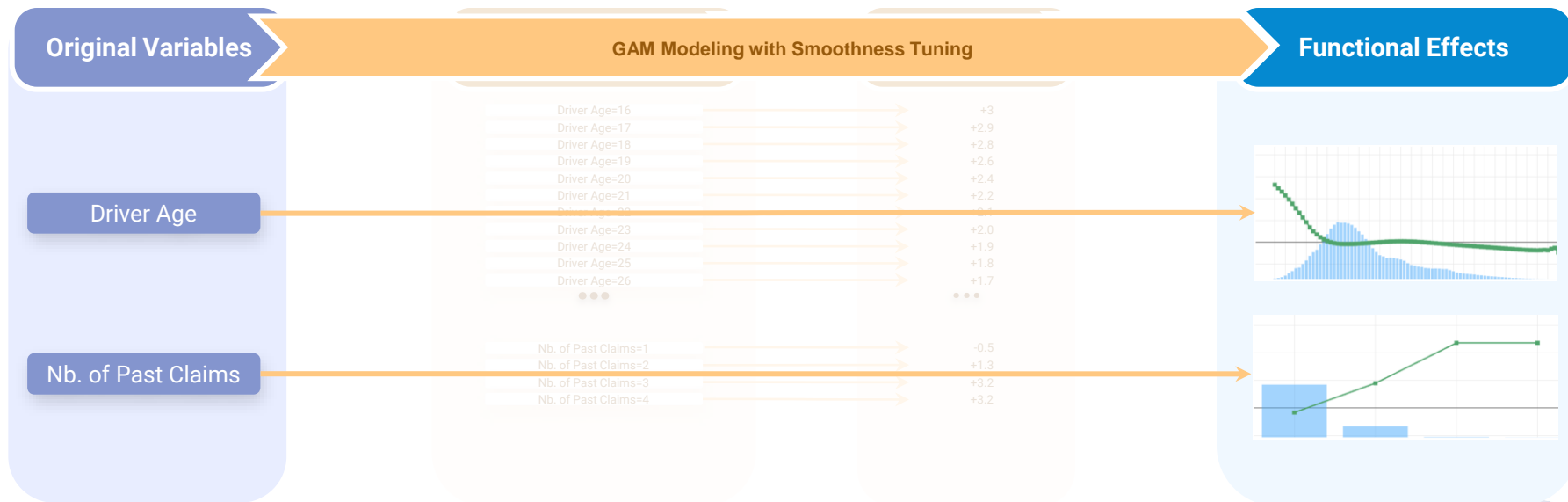
A very **weak prior**  
(with a large variance)  
leads to **noisy models**.



# Machine-Learning = GLM and Credibility

From a user's point of view, the creation of the models is **fully automated** and provides a unified machine-learning algorithm. As with all **machine-learning** techniques, the one presented today relies on a **solid statistical basis**.

A similar framework **can be leveraged to achieve variable selection**.



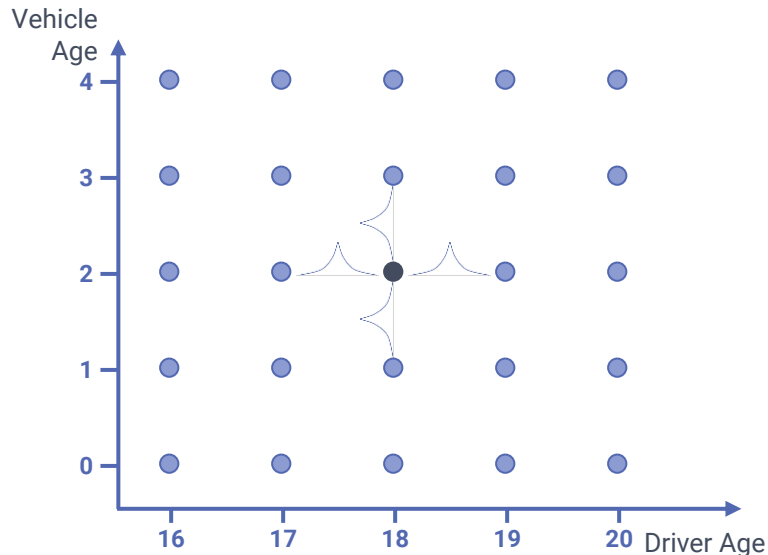
# Extending the framework

# Applying to Interactions

The same principle can be applied in **two dimensions, to fit interactions**. The prior there is slightly different to take into account the 2-D nature of the problem.

For instance, on an interaction between two ordered variables, we could suppose as prior that the differences between all the “connected” levels are supposed to follow a Laplace distribution.

The prior term would become:  $Penalty(\beta) = \dots + \lambda|\beta_{18,2} - \beta_{19,2}|$   
 $+ \lambda|\beta_{18,2} - \beta_{17,2}|$   
 $+ \lambda|\beta_{18,2} - \beta_{18,1}|$   
 $+ \lambda|\beta_{18,2} - \beta_{18,3}| + \dots$



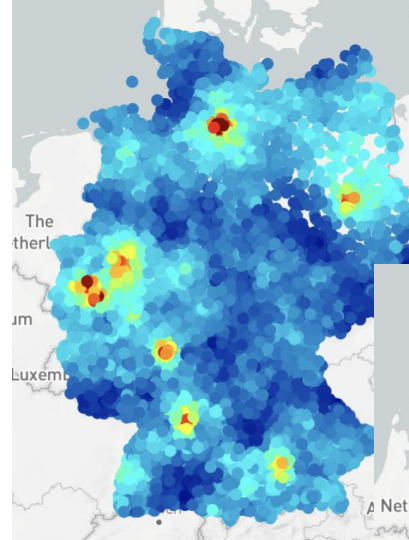


# Applying to Geography

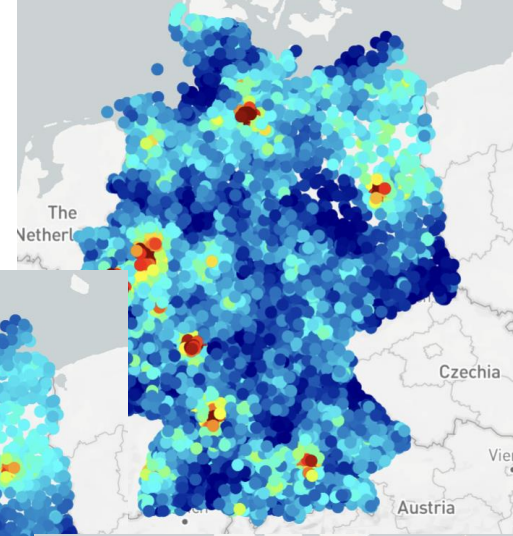
Geographic modeling can also be achieved with a similar method : the prior is that **nearby locations are expected to have similar risk levels.**

This has strong similarities to a **Gaussian Process** modeling.

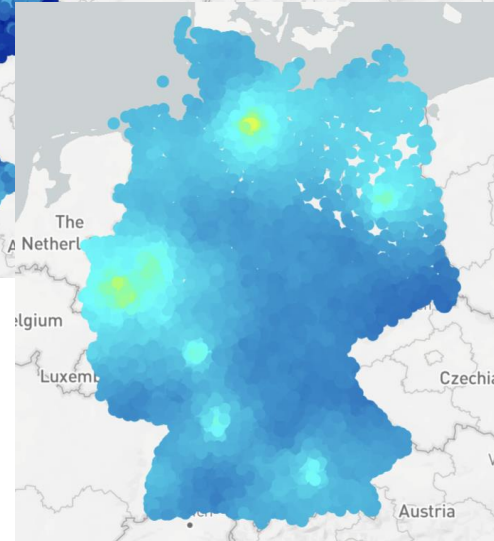
Intermediate Prior



Weak Prior



Strong Prior

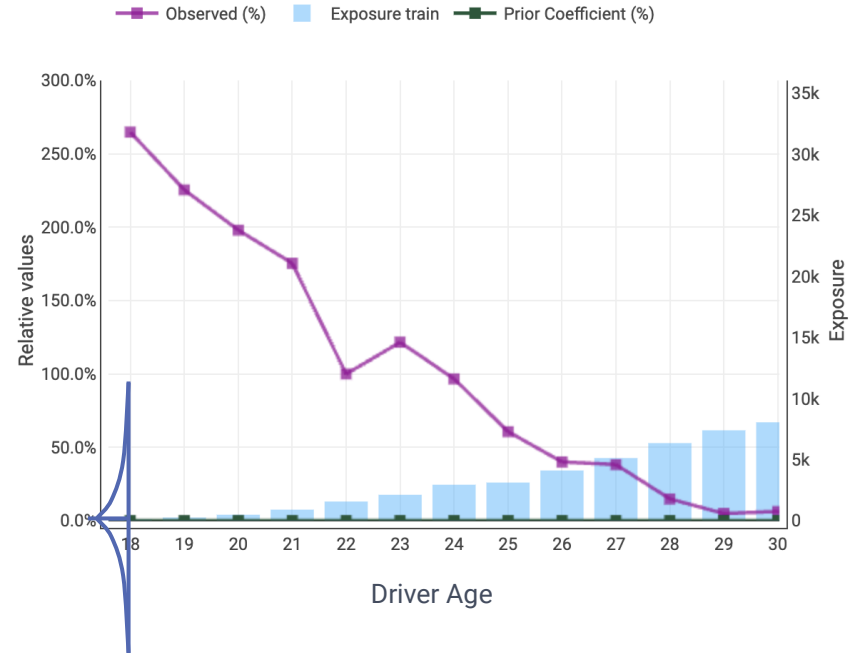


# Leveraging a different prior estimate

In all the examples above, the priors used were centered on zero:

- For categorical variables, we assumed the coefficients were **most likely to be null**.
- For ordered variables, we assumed the coefficients were **most likely to be constant**.

In both case, the most likely prior model is null.

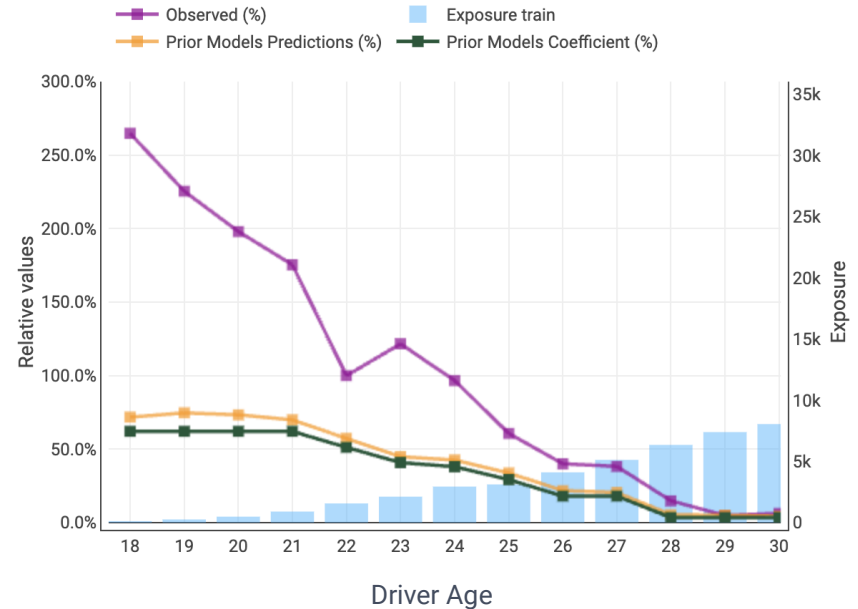


# Leveraging a different prior estimate

We could use a different set of assumptions. For instance we can assume that the coefficients are **most likely to follow a reference model**.

The assumed distribution would then be following that model:

- For categorical variables, the updated coefficients can be assumed to be **centered on the reference coefficients**.
- For ordinal variables, the derivative of the coefficients can be assumed to be **centered on the derivative of the reference coefficients**.

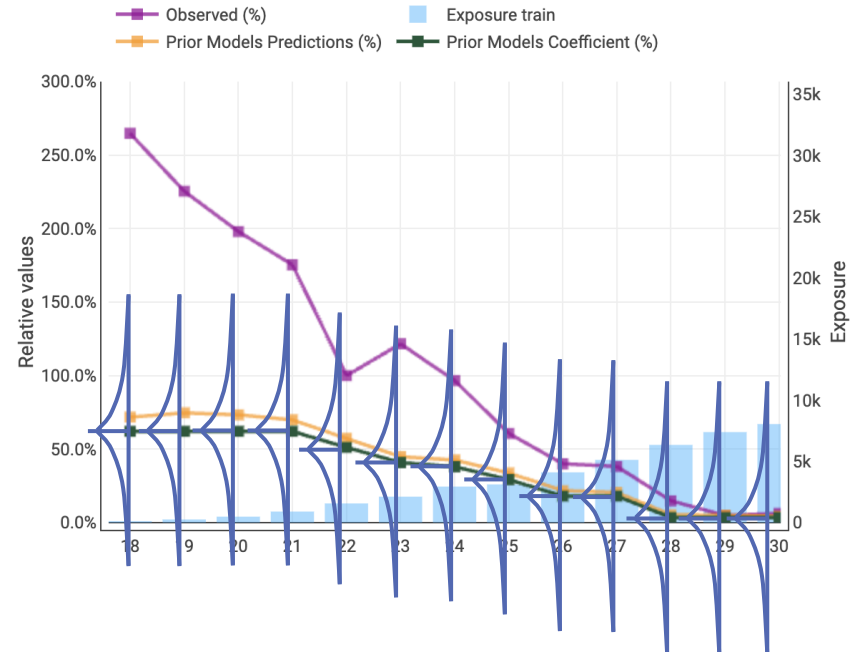


# Leveraging a different prior estimate

We could use a different set of assumptions. For instance we can assume that the coefficients are **most likely to follow a reference model**.

The assumed distribution would then be following that model:

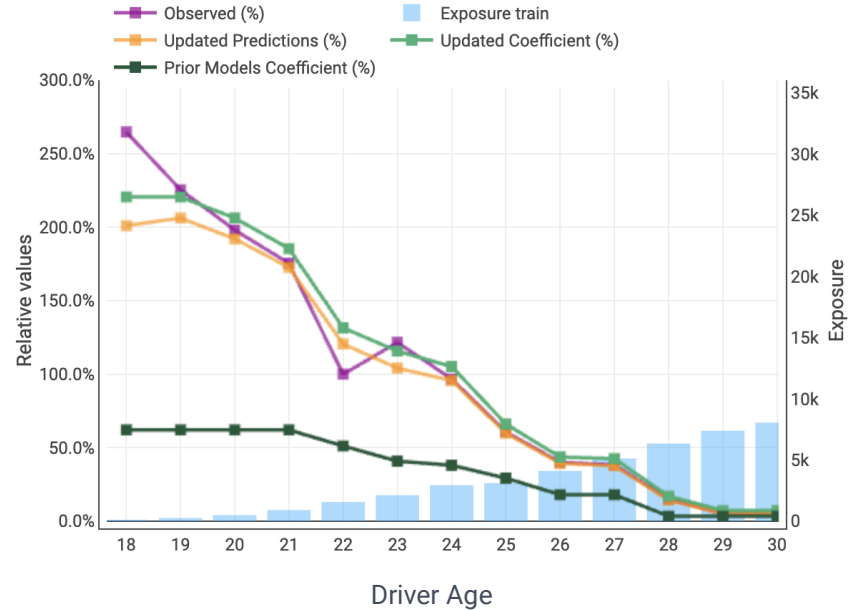
- For categorical variables, the updated coefficients can be assumed to be **centered on the reference coefficients**.
- For ordinal variables, the derivative of the coefficients can be assumed to be **centered on the derivative of the reference coefficients**.



# Leveraging a different prior estimate

We could use a different set of assumptions. For instance we can assume that the coefficients are **most likely to follow a reference model**.

The new model would then be a **balance between the reference model and the new observed data** (which is aligned with a credibility-based intuition).



# Data-driven Modelling

## A three-step process

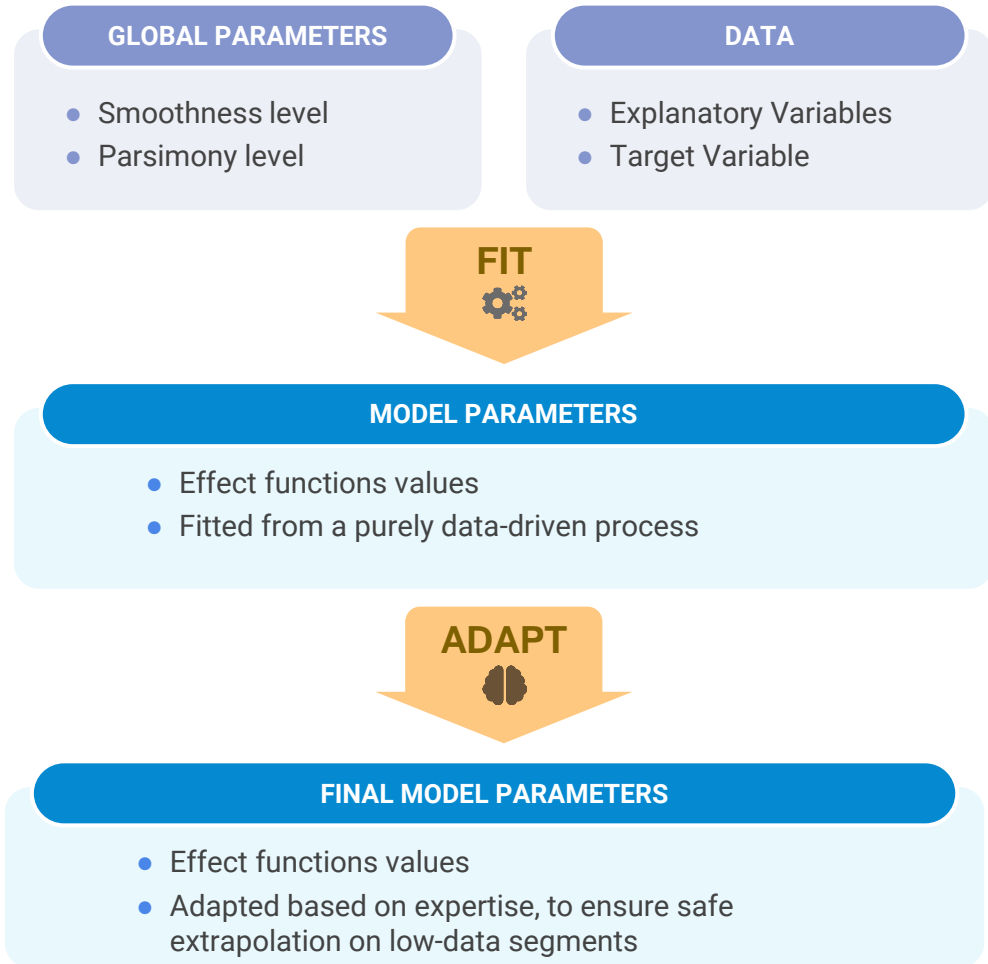
It is possible to directly leverage a model right out of the fit process.

This would be similar to a classic data-science approach.

However, handling transparent models opens the possibility of interacting with them, integrating expert knowledge in the modeling.

So the process is (on purpose) mixing elements of:

- Machine-Learning: **automated fit**, purely **data-driven** model creation, acting on **global parameters** to control overfitting.
- Direct interaction with the models: control of all the **effects** captured in the fitting model, analysis and potentially edition of the **effects** to ensure a good extrapolation of the model.



# Dankeschön!

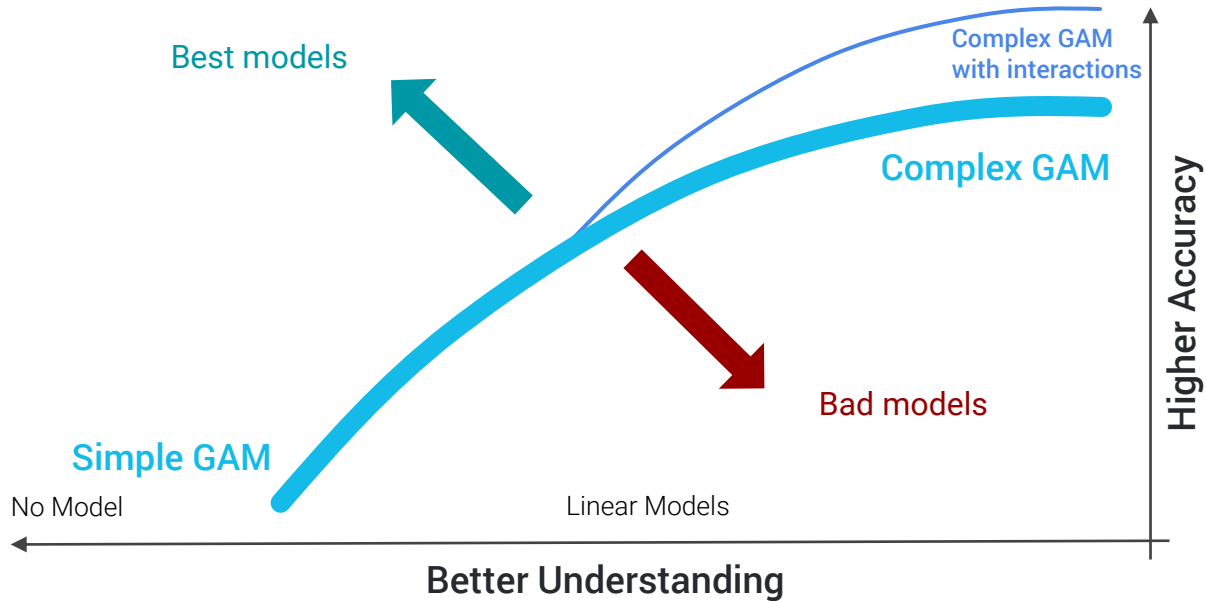
Ihre Fragen?

[jan.kuethe@akur8.com](mailto:jan.kuethe@akur8.com)

 **AKUR8**

# Extending to Variables Selection

# Complexity has a cost



The accuracy is measured on a back-test; actual results when moving to productions will not be

# Constraining the number of variables

---

While the control of a **models smoothness** provides a good management of overfitting, “**simpler**” **models, containing a limited number of variable, are easier to read.**

The creation of a simpler model is achieved by **adding a constraint in the computation of the coefficients**:

$$\left\{ \begin{array}{l} \beta^* = \text{Argmax}_{\beta} LL(x, y, \beta) - \lambda |\beta_i - \beta_{i+1}| \\ \#Variables < N \end{array} \right.$$

This constraint will ensure that the model is simple and keeps a good balance between transparency and predictiveness.

# Constraining the number of variables

$$\left\{ \begin{array}{l} \beta^* = \text{Argmax}_{\beta} LL(x, y, \beta) - \lambda |\beta_i - \beta_{i+1}| \\ \#Variables < N \end{array} \right.$$

This problem is equivalent to finding the optimal coefficients **placing a “cost” on the number of variables** in the model (using the optimization “lagrangian trick”).

$$\beta^* = \text{Argmax}_{\beta} \underbrace{LL(x, y, \beta)}_{\text{Minimize the Training Error}} - \underbrace{\lambda |\beta_i - \beta_{i+1}|}_{\text{Maximize the smoothness (control for overfitting)}} - \underbrace{\lambda_v \cdot \#Variables}_{\text{Minimize the number of variables (control for quality)}}$$

The Lagrangian trick is a common optimization technique that allows to transform a problem with multiple constraints into a simpler optimization problem\*.

# Solving the Constrained Optimization Problem

So we want to fit a model by finding the optimal coefficients  $\beta^*$ :

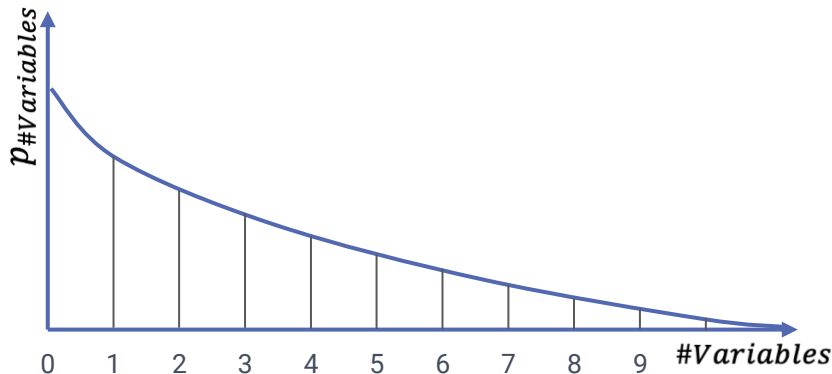
$$\beta^* = \text{Argmax}_{\beta} LL(x, y, \beta) - \lambda |\beta_i - \beta_{i+1}| - \lambda_v \cdot \#Variables$$

Interestingly, it is possible to relate back the terms to different priors:

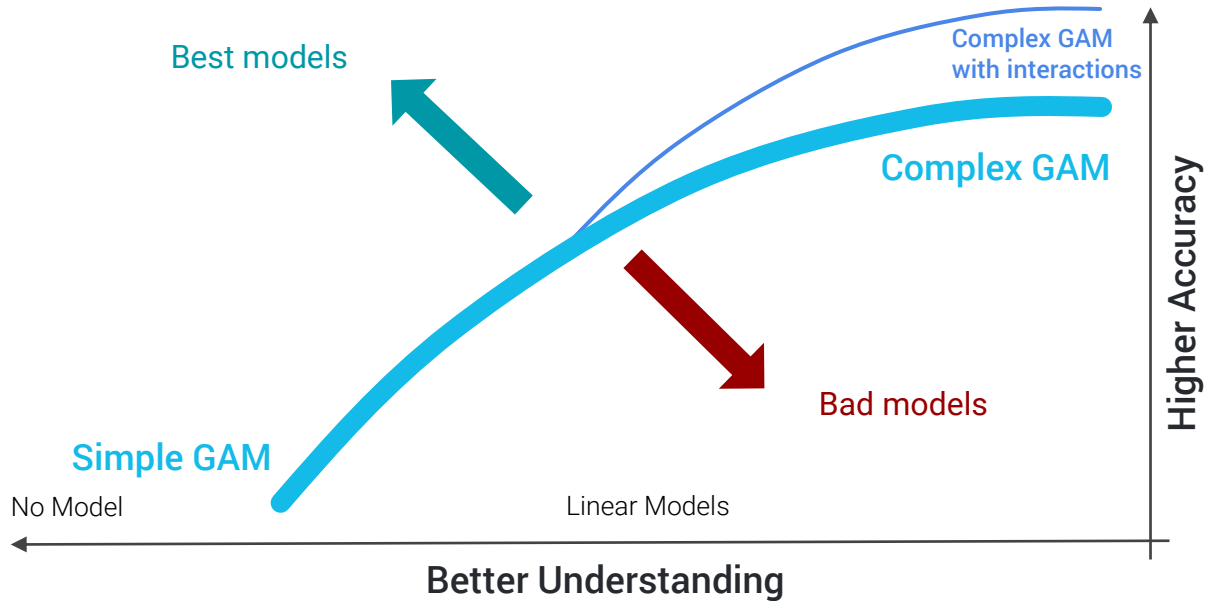
$$\beta^* = \text{Argmax}_{\beta} p(y|\hat{y}(X)) \times p_{\text{prior}}(\beta) \times p_{\#Variables}(\beta)$$

with  $p_{\#Variables}(\beta) \sim e^{-\lambda_v \cdot \#Variables}$

**Models with many variables are less likely** than the simpler ones.

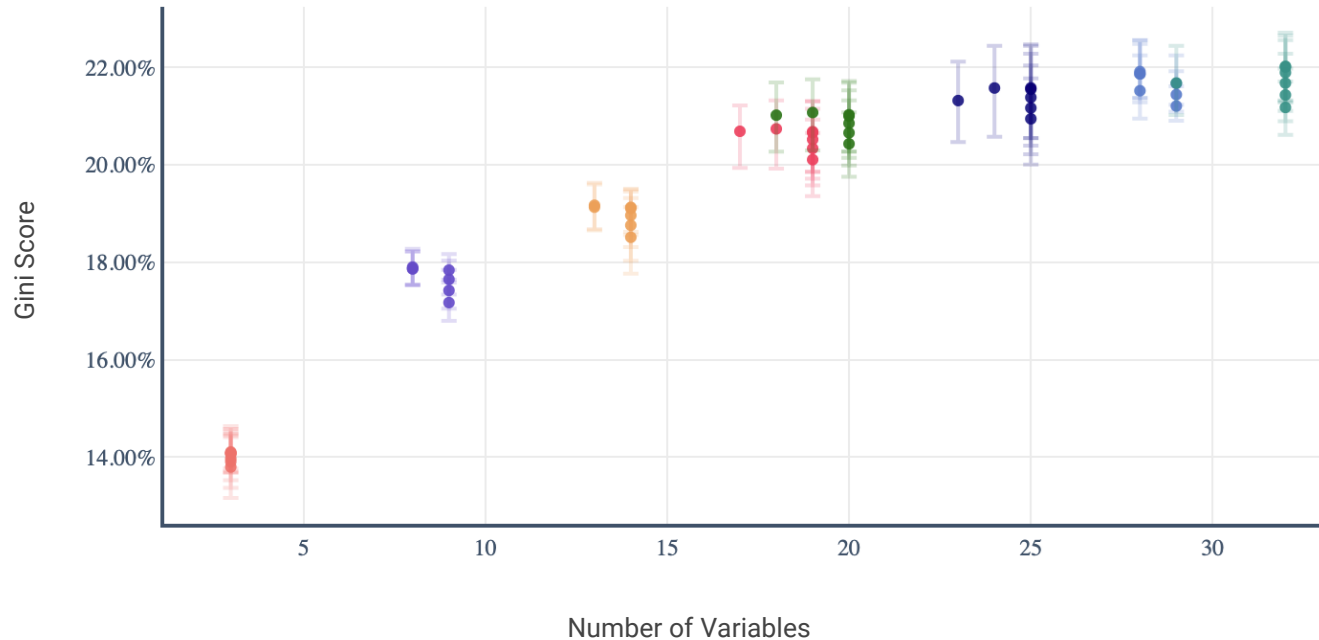


# Complexity has a cost



The accuracy is measured on a back-test; actual results when moving to productions will not be

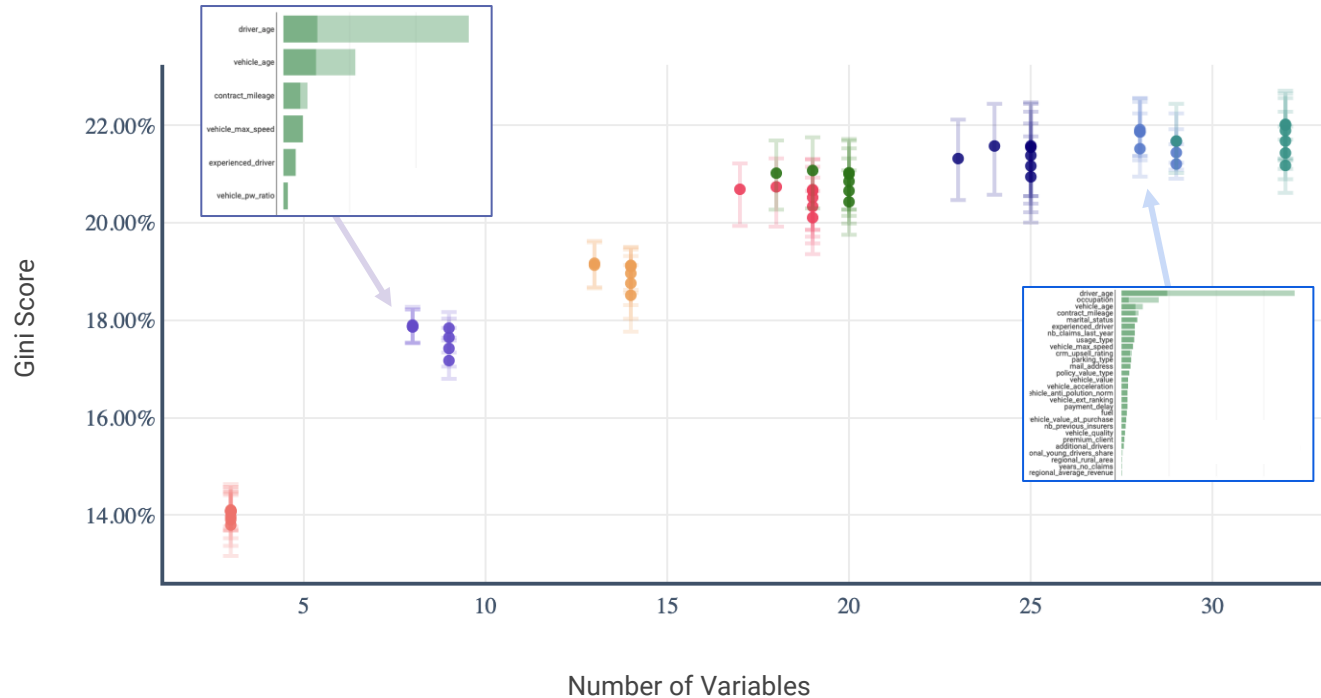
# Complexity has a cost



**Grid-search results:**  
each **point** represents  
one **model**.

The gain in models quality and  
the fading marginal  
improvement are clearly  
visible.

# Complexity has a cost



**Grid-search results:**  
each **point** represents  
one **model**.

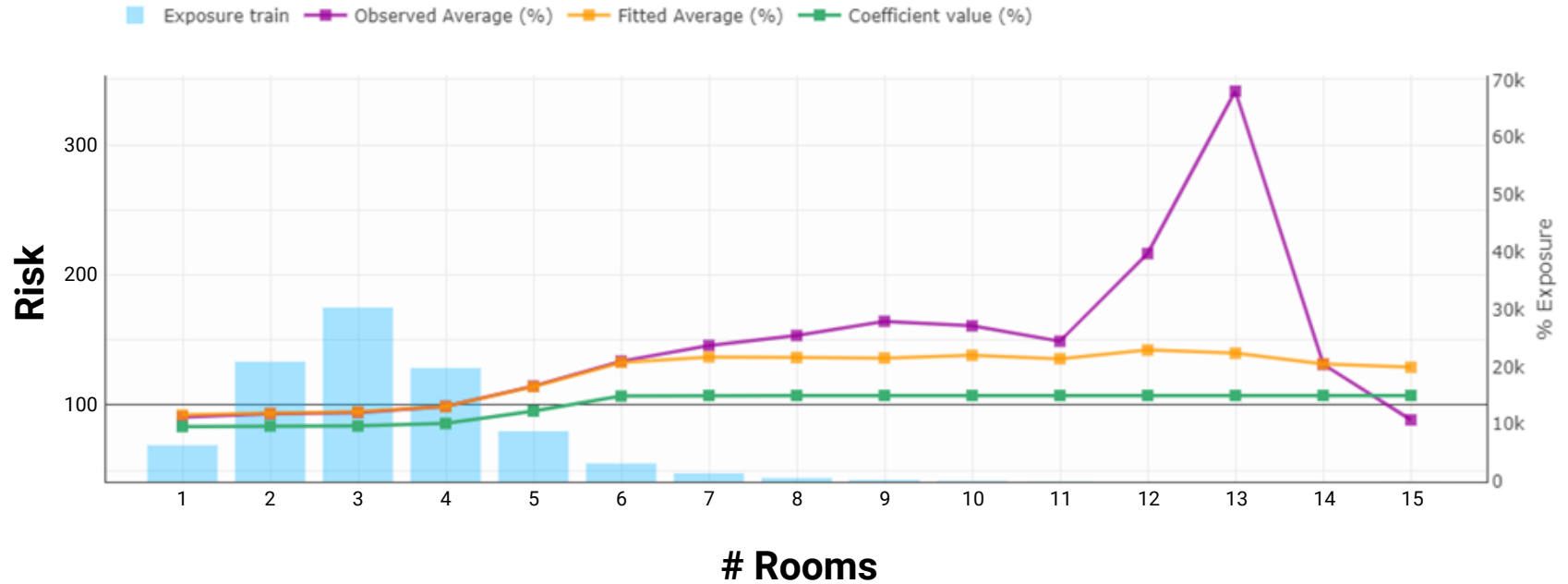
The gain in models quality and  
the fading marginal  
improvement are clearly  
visible.

# APPENDIX

# Practical Impact

# Interact with the models

Spotting the issues is nice..



# Interact with the models

... solving the issues is better !

