



IACA WEBINAR

## Explainable machine learning for actuaries

SEPTEMBER 2023



THIS PAGE IS LEFT BLANK INTENTIONALLY

# SPEAKER INTRODUCTION



## Xavier Maréchal

- CEO Reacfin
- MSc. in Engineering Science (Applied Mathematics), MSc. in Actuarial Sciences and MSc. in Business Management
- Co-author of “Actuarial Modelling of Claim Counts: Risk classification, Credibility and Bonus-Malus Systems”
- Consultant for 20 years in Non-Life and Health insurance (Pricing, DFA models, Solvency 2)



Reacfin s.a. is a **consulting firm**, spin-off of the University of Louvain (Louvain-la-Neuve – Belgium).

We develop, in partnership with our clients, **actuarial** & quantitative **financial** solutions, building on strong **data analytics**, while securing full transparency and integral knowledge transfer.



We offer consulting services in actuarial science & quantitative finance, including a.o. capital, portfolio, product, risk and liquidity management. We build our expertise on broad data analytics capacities.



We develop solutions in partnership with our clients, i.e. we integrate our solutions in our client's systems and processes and we secure full knowledge transfer (e.g. open source code).



We share our knowledge with our clients. We offer a comprehensive learning platform, including on-site trainings, e-learning modules, e-classrooms and webinars.

# GOALS OF THIS PRESENTATION

## The problem

---

- Whereas advanced Machine learning (ML) techniques (e.g. random forest or neural networks) usually have a better predictive power than statistical techniques (e.g. GLM), their main drawback is that they are black-box and their results are difficult to understand/interpret.

## Two different strategies to use ML for practical applications

---

- There are basically 2 strategies to use ML techniques in predictive modelling
  1. **Replacing** traditional models (e.g. GLM) by ML models
  2. **Combining** the pros of traditional and ML models to improve predictive modelling
- The goals of this presentation are therefore to
  - Briefly **remind some useful machine learning techniques** and explain why it is difficult to interpret their results
  - Present several techniques that have been developed in order to **better understand the results** of machine learning techniques
  - Explain how these **interpretation techniques** can be used to implement the 2 strategies presented above and improve predictive modelling

# AGENDA

A non-exhaustive reminder to some useful ML techniques

Adding complexity means increasing need for interpretability

An introduction to ML interpretation tools

Conclusions: how to make the most of ML techniques

# WHAT IS MACHINE LEARNING?

## Objectives of Machine Learning (“ML”)

**ML algorithms aim at finding by themselves the method that best predicts the outcome of the studied phenomenon.**

## Supervised vs. Unsupervised learning

### ■ Supervised learning:

- Inputs and examples of their desired outputs are provided
- The goal is to learn a **general rule that maps inputs to outputs**.

➔ *Given a set of training examples  $(x_1, x_2, \dots, x_n, y)$ , where  $y$  is the variable to be predicted, what is the most efficient algorithm to best approximate the realizations of  $y$*

- 2 main techniques
  - **Classification** : inputs are divided into two or more classes, and the learner must produce a model that assigns unseen inputs to one (or multi-label classification) or more of these classes.
  - **Regression**: the outputs are continuous rather than discrete.

### ■ Unsupervised learning:

- No labels are given to the learning algorithm
- The goal is to **find structure in its input** (discovering hidden patterns in data).
- Main technique
  - **Clustering**: a set of inputs is to be divided into groups. Unlike in classification, the groups may not be known beforehand.

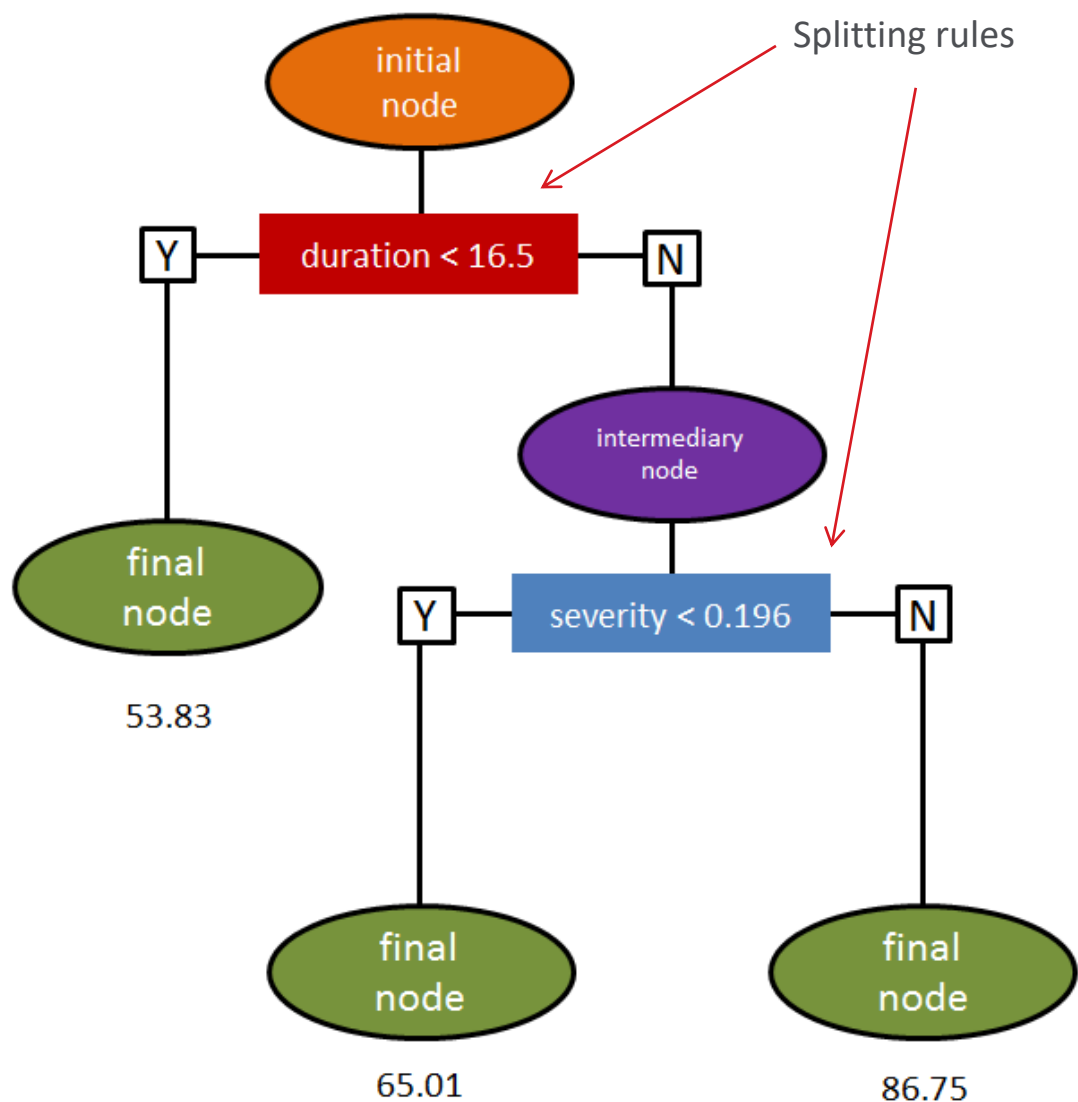
## Examples of use

- Typically used to model **pricing or underwriting related target variables** in function of available **features**
  - Regression: frequency (#claims) or severity (claims cost)
  - Classification: lapse rates, conversion rates
- Typically used for **features engineering** (i.e. creating new variables)
  - E.g. vehicle classification, zoning,...

For a more complete presentation of some supervised models, check Reacfin webinar on “Machine learning applications to non-life pricing” <https://www.reacfin.com/index.php/reactfinacademy-2-2/guided-training/>

**Focus on supervised models**

# A FIRST SIMPLE ML MODEL: CLASSIFICATION AND REGRESSION TREES (CART)



## Purpose

- Tree enables to **segment the predictor space** into a number of simple homogenous regions defined according to the covariates
- **Splitting rules** can be summarized in a tree view
- For each region the prediction is set as the **region average**

## Definitions

- The *root node* in orange:
  - at the top of the tree
  - contains the whole population
- The splitting rules set aim at segmenting the predictor space into a number of **simple regions** that are as **homogeneous** as possible with respect to the response variable
- The *leaves nodes* in green at the bottom of the tree: that is a node that is not further split.

# AN EXAMPLE OF A MORE COMPLEX ML MODEL: BOOTSTRAP AGGREGATION (BAGGING)

## Main idea

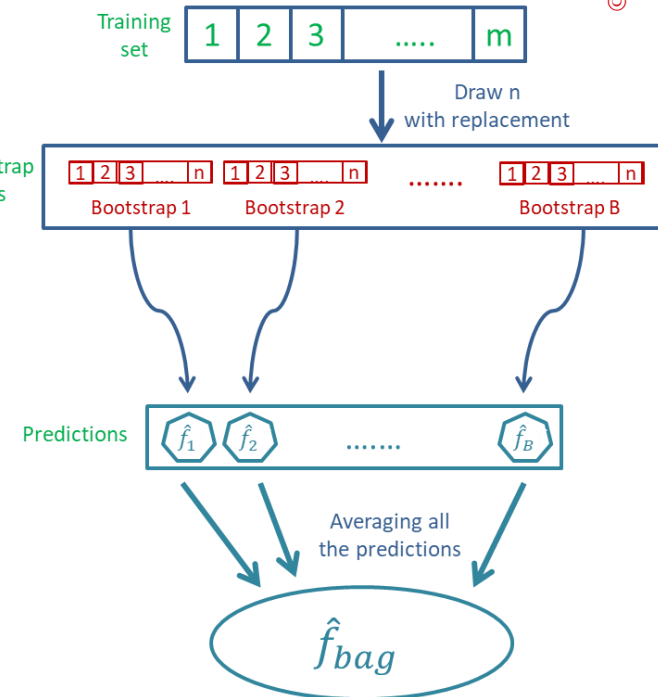
- Bootstrap **aggregation**, or **Bagging**, is a general-purpose procedure for reducing the variance of a statistical learning method
- Recall that given a set of  $n$  independent observations  $Z_1, Z_2, \dots, Z_n$  each with variance  $\sigma^2$ , the variance of the mean  $\bar{Z}$  of the observations is given by  $\frac{\sigma^2}{n}$ .
- **Averaging a set of observations reduces variance.**

## Algorithm

1. Bootstrap, by taking **repeated samples** from the training data set.
2. Generate  $B$  different training data sets.
3. **Train our method** (e.g. regression tree) on the  $b$ th bootstrapped training set to get  $\hat{f}_b(x)$  the prediction at point  $x$ .
4. We then **average all the predictions** to obtain:

$$\hat{f}_{bag}(x) = \frac{1}{B} \sum_{b=1}^B \hat{f}_b(x)$$

- ➡ The final prediction is **difficult to understand** as it is an **average** of many “intermediate” predictions
- ➡ Similar difficulty in interpreting results is also an issue for other widely used ML models (Random Forests, Gradient Boosting Models, Artificial Neural Networks,...)





# AGENDA

A non-exhaustive reminder to some useful ML techniques

**Adding complexity means increasing need for interpretability**

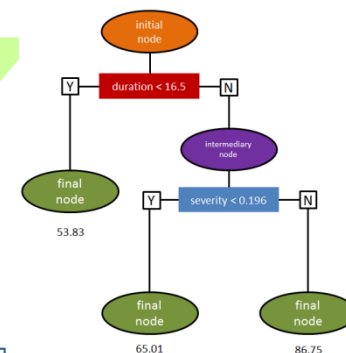
An introduction to ML interpretation tools

Conclusions: how to make the most of ML techniques

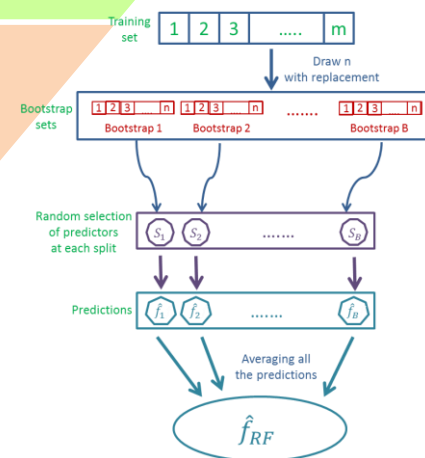
# SOME MACHINE LEARNING TECHNIQUES ARE BLACK BOXES AND INTERPRETATION OF THE RESULTS CAN BE QUITE DIFFICULT

Understanding the results of ML techniques is not easy

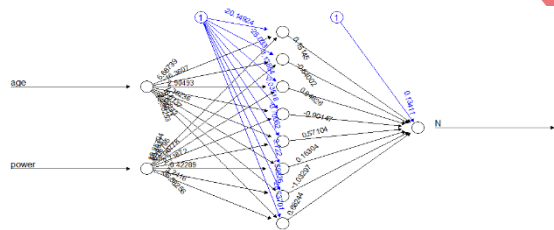
In the case of regression trees, understanding how the model predicts claims' cost or frequency values for new data points is not a problem, as it is very intuitive.



In the case of more complex methods such as Bagging and Random forests, even understanding how the model predicts values for new data points is rather difficult.



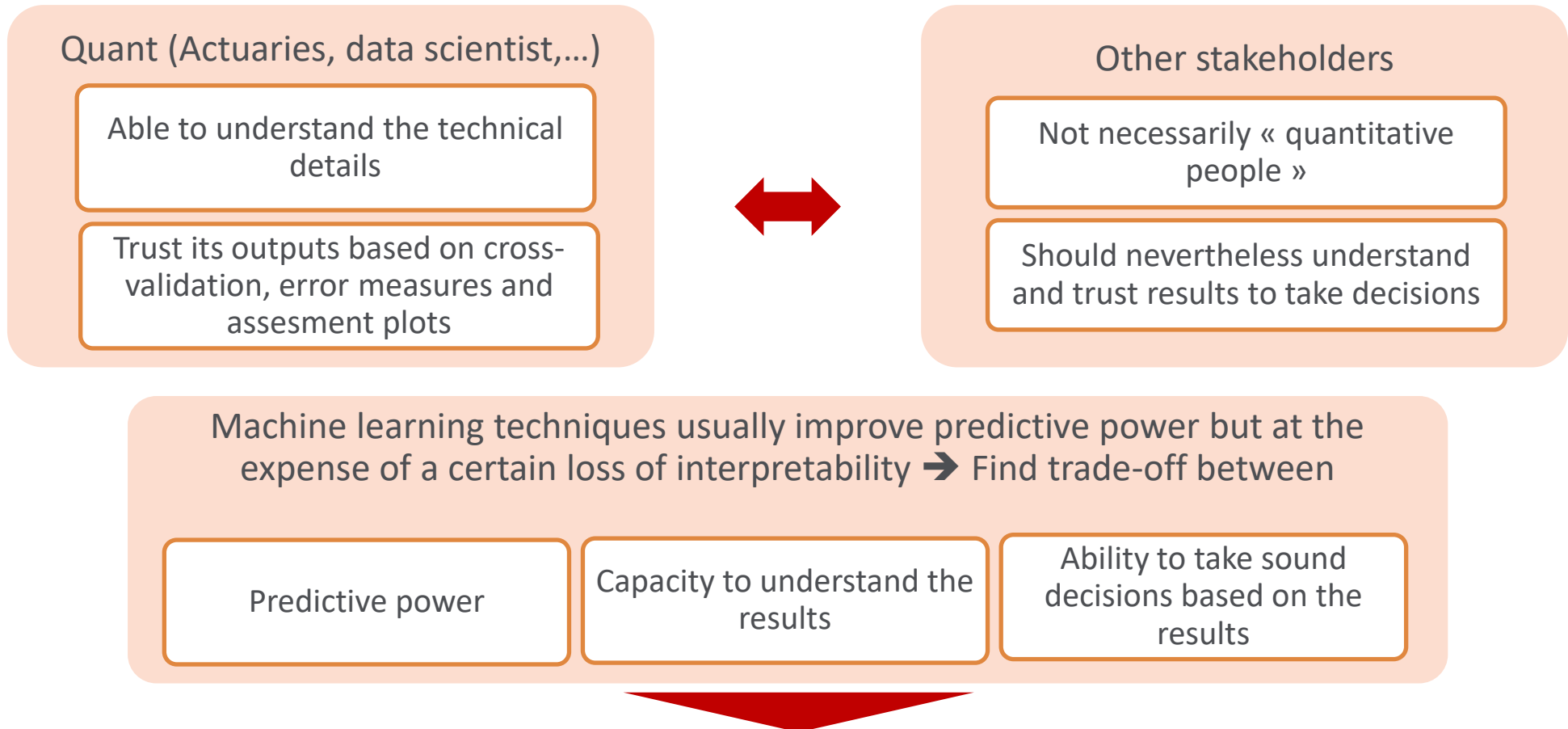
Things may be even worse for GBM and NN.



Complexity

Interpretability

# UNDERSTANDING THE RESULTS OF ML MODELS IS NEVERTHELESS KEY FOR SOUND BUSINESS DECISION-MAKING AS MANY STAKEHOLDERS USE THE RESULTS OF THE MODELS



## High-end questions

**Who will use the results? For what purpose? With which impact?**

# AGENDA

A non-exhaustive reminder to some useful ML techniques

Adding complexity means increasing need for interpretability

**An introduction to ML interpretation tools**

Conclusions: how to make the most of ML techniques

# GLOBAL VS LOCAL INTERPRETABILITY OF ML TECHNIQUES

## ■ Global Model Interpretability

### ○ How does the trained model make predictions?

- Which features are important and what kind of interactions between them take place?
- Global model interpretability helps to understand the distribution of your target outcome based on the features.
- Global model interpretability is very difficult to achieve in practice → Any model that exceeds a handful of parameters or weights is difficult to understand
- Some models are interpretable at a parameter level :
  - For linear models, the interpretable parts are the weights,
  - For trees interpretable parts are the splits (selected features plus cut-off points) and leaf node predictions.

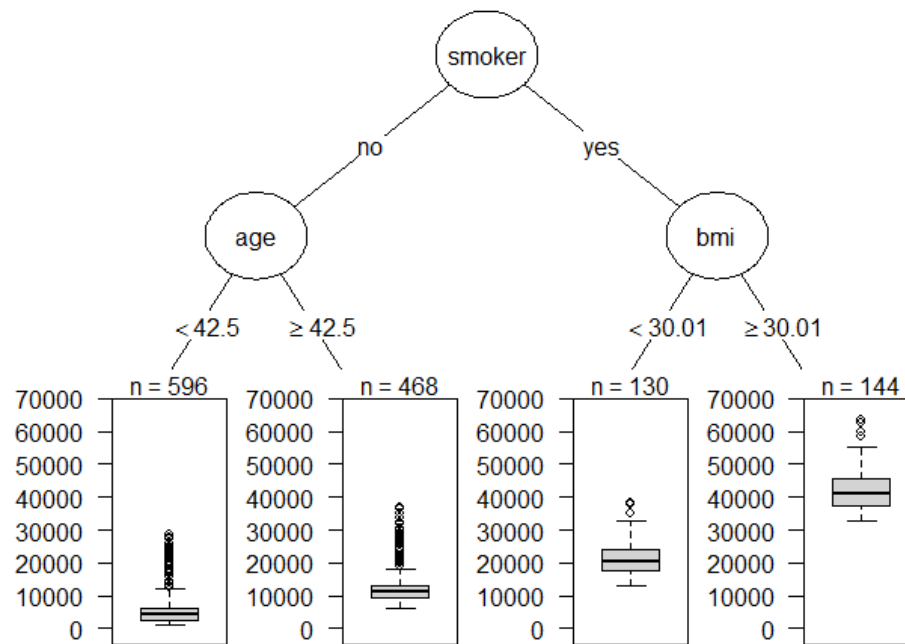
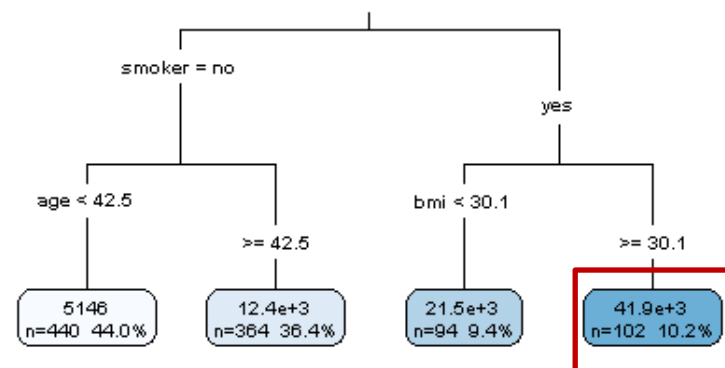
### ○ Global Interpretable tools

- Interpretable Models by nature (eg. Linear models, Regression Tree)
- Feature Importance
- Partial Dependence Plot (PDP), Individual Conditional Expectation (ICE) and Accumulated Local Effects (ALE)
- Interaction Measures (H-statistic)

# GLOBAL MODEL INTERPRETATION

## Interpretable Models by nature

- By Nature regression tree are **easy to interpret** :
  - Starting from the root node,
  - Go to the next nodes and the split rules tell you which subsets you are looking at.
  - Once you reach the leaf node, the node tells you the predicted outcome.
- In this dummy example, the smokers with a large Body Mass Index (>30) have the highest average claims amount (41,9k€ - see first graph)
- We can even obtain the distribution of the observed claims in each segment (second graph)

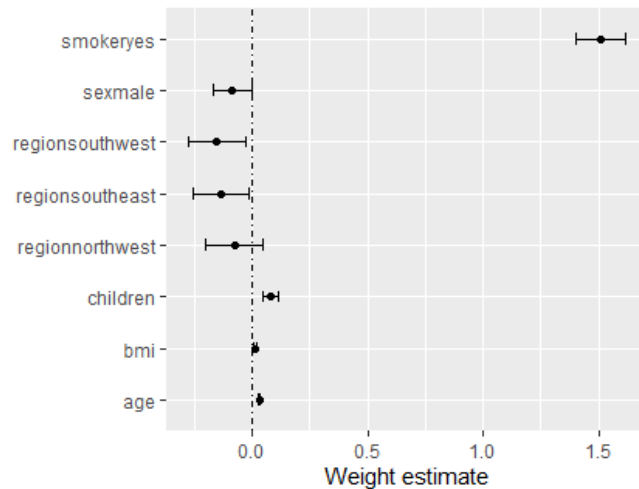


# GLOBAL MODEL INTERPRETATION

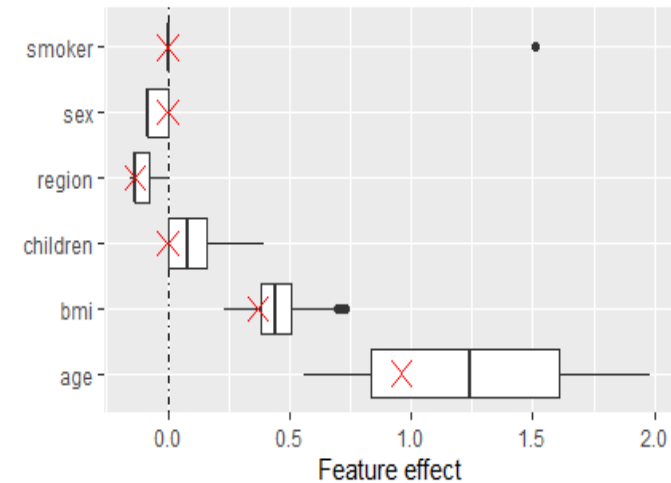
## Interpretable Models by nature

- One should also pay attention when interpreting linear model
- The problem with the coefficients is that the features are measured on different scales
- Features effects of the linear regression model can be more meaningfully analyzed as they are coefficients multiplied by the actual feature values

### Model Coefficients



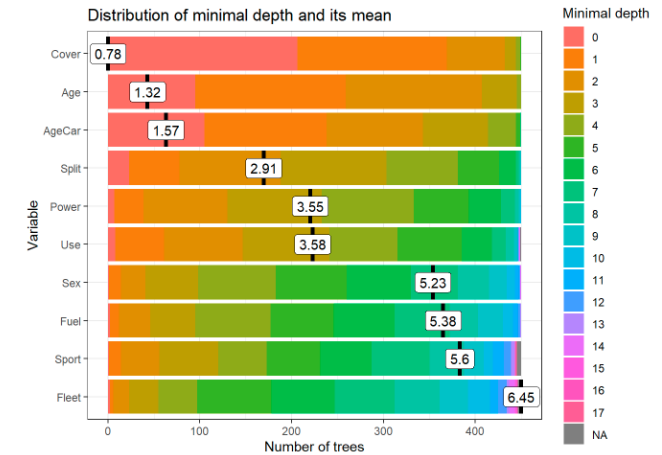
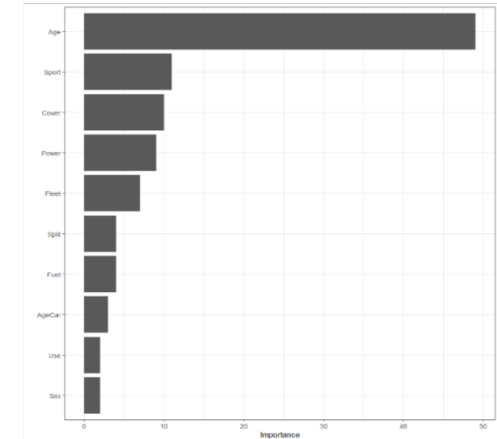
### Feature Effects



# GLOBAL MODEL INTERPRETATION

## Features Importance

- Features Importance:
  - In a tree-based method : Go through all the splits for which the feature was used and measure **how much it has reduced the Loss Function** (eg. Gini, MSE, Poisson Deviance,...) compared to the parent node.
  - The sum of all importance measures is scaled to 100.
  - This means that each variable importance can be interpreted as share of the overall model importance
- One can get additional measures such as:
  - Minimal depth and its mean :
    - ✓ Which variables were the most often on the top of the tree
    - ✓ Mean depth of first split
- Features Importance can be used as a **features' selection tool**
  - Goal: Identify the **most relevant variables**
  - Pay attention: when some variables are correlated, their **global impact can be spread** between them, therefore reducing individual importance of each variable





# GLOBAL MODEL INTERPRETATION

## Partial dependence plot

### ■ Partial Dependence Function/Plot

- Partial dependence plot (short PDP or PD plot) shows the **marginal effect one or two features** have on the predicted outcome of a machine learning model
- Partial dependence plot can show whether the **relationship between the target and a feature** is linear, monotonic or more complex. It can be estimated as

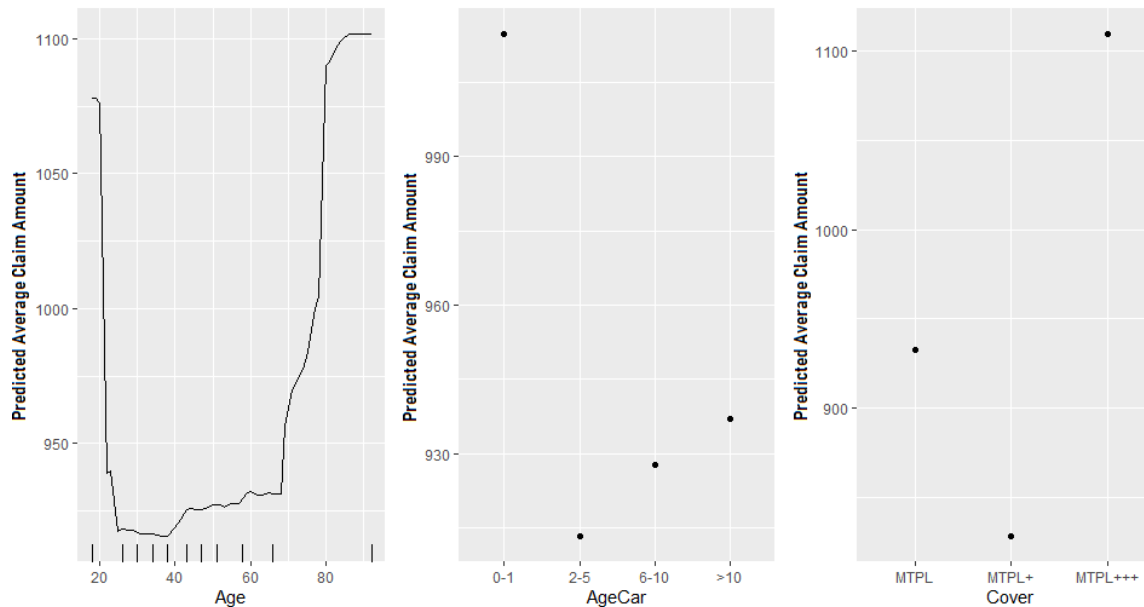
$$PD_{age}(age) = \frac{1}{n} \sum_{i=1}^n \hat{f}(age, agecar^i, cover^i, \dots)$$

- In this formula,  $agecar^i, \dots$  are actual features' values from the dataset for the features in which we are not interested,  $\hat{f}$  is the trained model and  $n$  is the number of instances in the dataset.
- So we marginalize model outputs over the distribution of the features we are not interested in (e.g.  $agecar$ ,  $cover$ , ...)
  - the function shows the relationship between the feature  $age$  we are interested in and the predicted outcome.
  - By marginalizing over the other features, we get a function that depends only on feature  $age$ , interactions with other features included.

# GLOBAL MODEL INTERPRETATION

## Partial dependence plot

- Example of Partial Dependence Plot (1D) on Average Claim Amount :

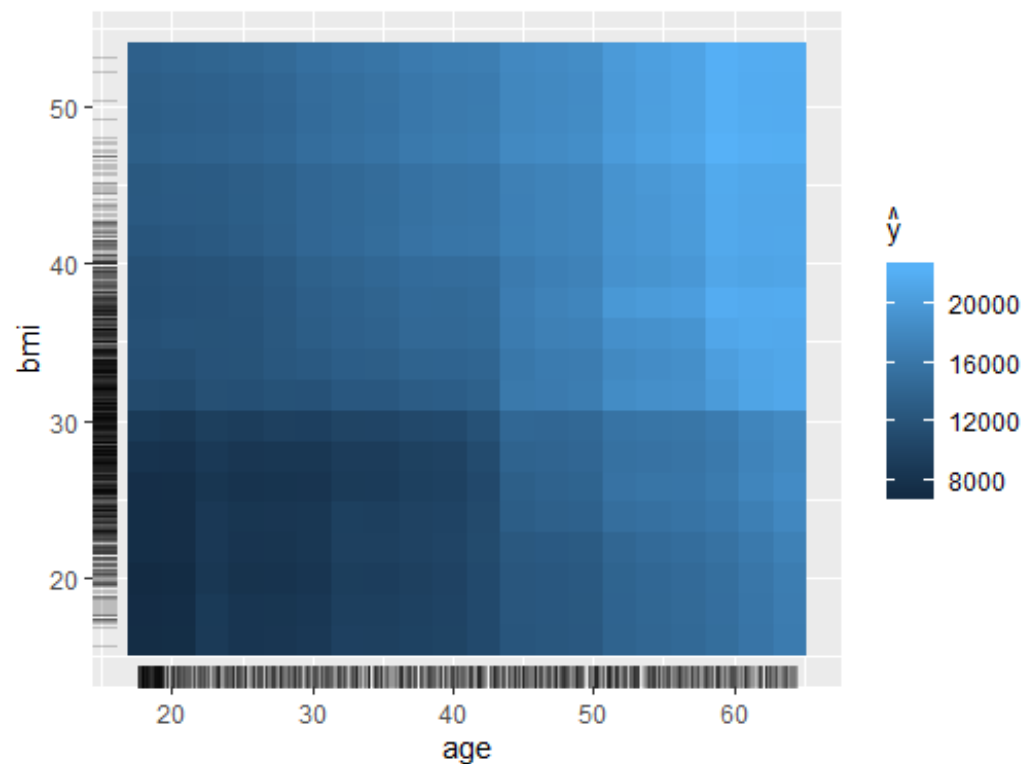


- Partial dependence plot can be used as a **features' impact explanation tool**
  - It allows to better understand the marginal impact of a variable on the prediction
  - It is very similar to the interpretation of the multiplicative factors we obtain in a GLM or GAM model

# GLOBAL MODEL INTERPRETATION

## Partial dependence plot

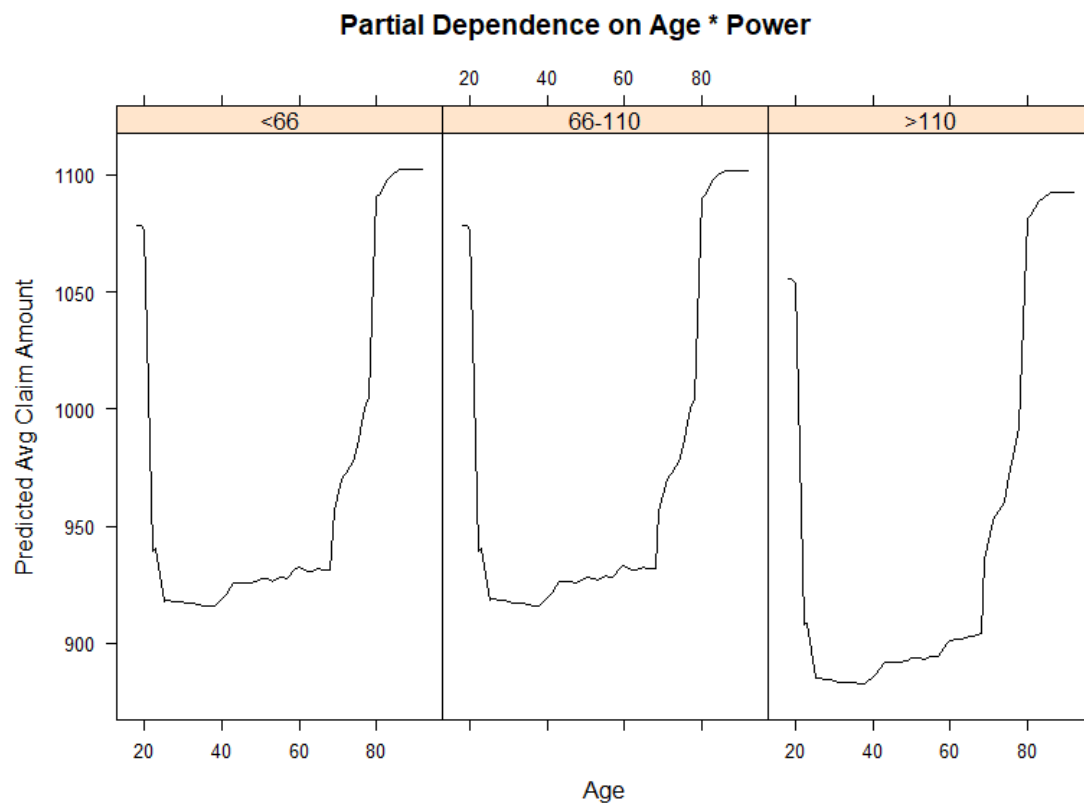
- Example of Partial Dependence Plot (2D) :
  - PD can be generalized to more than one feature
  - PDP -2D can be very useful to highlight interactions



# GLOBAL MODEL INTERPRETATION

## Partial dependence plot

- Example of Partial Dependence Plot (2D) on Average Claim Amount :



# GLOBAL MODEL INTERPRETATION

## Basic example with only two features

Original data features		
Age of the driver	License Age	Predicted Frequency
20	1	6,2%
35	10	5,1%
20	3	5,5%
55	32	4,2%
60	40	4,3%

PDP Age of the driver		
Age of the driver	License Age	Predicted Frequency
20	1	6,2%
<b>20</b>	10	<b>3,0%</b>
20	3	5,5%
<b>20</b>	32	<b>1,0%</b>
<b>20</b>	40	<b>0,5%</b>
PDP Age(20)=		<b>3,2%</b>

### Partial Depend Plot

- PDP computes what the model predicts on average when each data instance has the value 20 for driver age.
- Weird instances are created during the calculation process (see yellow rows)
- Marginal distribution is used so all instances in the data set enter in the calculation for each driver age computation.
- Computation time can be huge with large dataset.

# GLOBAL MODEL INTERPRETATION

## Partial dependence plot

### ■ Attention point with Partial Dependence Plot

- Correlated features :
  - **With correlated features**, computation of a PDP involves averaging predictions of artificial data instances that can be **unlikely in reality**.
  - E.g. “Age of the driver” and “License Age” in motor insurance : we don’t expect a 20 years old policyholder with 10 years of license whereas PDP computation process will consider this type of instance...
- 1D Flat PDP does not imply that the feature has no influence!
  - Interaction effect might still be there
  - E.g. half of the instance have a positive impact on the prediction and the other half has a negative impact. Both effects could cancel each other in the PDP.
  - These interactions effects can be observed in Individual Conditional Expectations (ICE – see further)

*Nice alternative to PDP are Accumulated local effect plot (ALE)*

# GLOBAL MODEL INTERPRETATION

## Example

Original data features		
Age of the driver	License Age	Predicted Frequency
20	1	6,2%
35	10	5,1%
20	3	5,5%
55	32	4,2%
60	40	4,3%

M-Plot Age of the driver		
Age of the driver	License Age	Predicted Frequency
20	1	6,2%
35	10	5,1%
20	3	5,5%
55	32	4,2%
60	40	4,3%
M-Plot Age(20)=		5,9%

### M-Plot

- M-plot (marginal plot) computes what the model predicts on average for policyholders that are **close to 20** years old.
- **Conditional distribution** is used so only instance where driver age is close to 20 are used in the calculation
- **Attention point:** The effect observed in M-Plot could be due to that feature, but also due to another correlated features (like License Age in our example)

### ALE-Plot

- Based on **Conditional distribution** (like M-Plot) but use **the sum of incremental effects** of the feature of interest in order to avoid effects of correlated features.
- Calculation out-of-the scope of this presentation (see *Daniel W. Apley and Jingyu Zhu 2019*)

# GLOBAL MODEL INTERPRETATION

## Individual Conditional Expectation

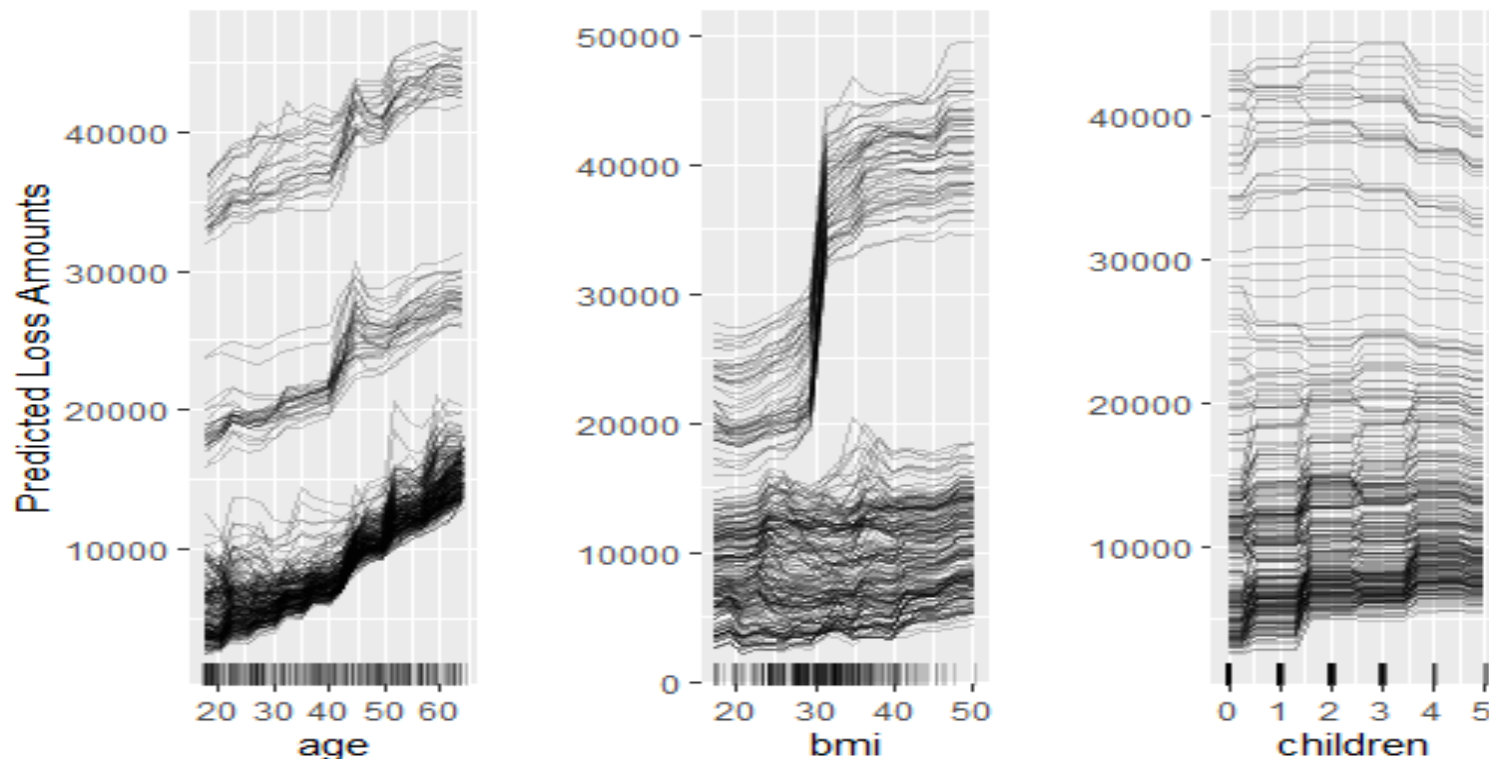
- One line per instance that shows how the instance's prediction changes when a feature changes
- An ICE plot visualizes the dependence of the prediction on a feature for each instance separately → one line per instance compared to one line overall in PDP.
- A PDP is the average of the lines of an ICE plot.
  
- **Advantage over PDP :**
  - In case of interactions, the ICE plot will provide much more insight.
  
- **How to compute ICE ?**
  - Creating variants of an observation by replacing the feature of interest value with values from a grid
  - Keeping all other features the same
  - Make predictions with the black box model for these newly created observations.
  - The result is a set of points for an original observation with the feature value from the grid and the respective predictions.



# GLOBAL MODEL INTERPRETATION

## Individual Conditional Expectation

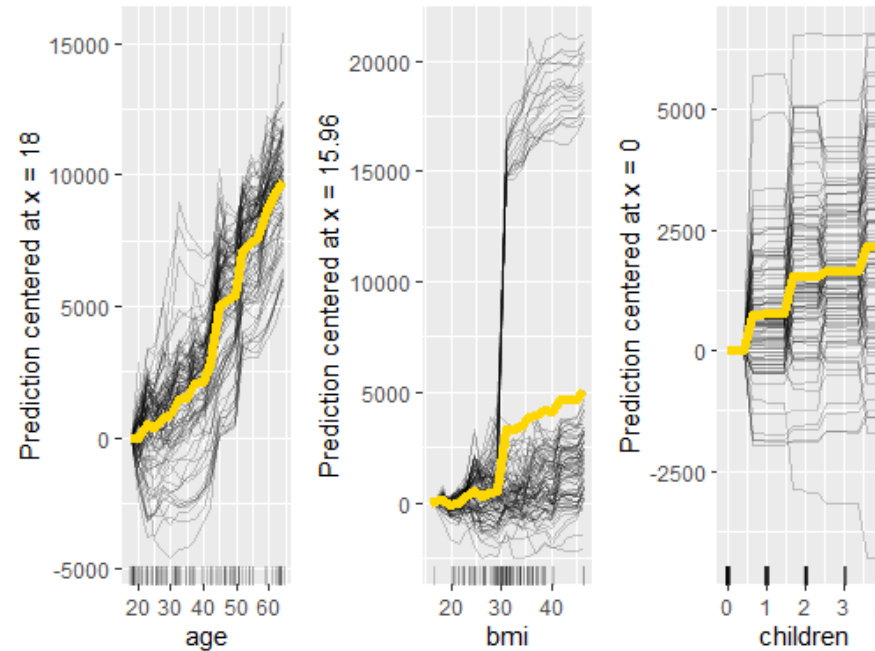
- Do you notice the interaction?



# GLOBAL MODEL INTERPRETATION

## Centered ICE

- Individual Conditional Expectation (ICE) :
  - It can be hard to tell whether the ICE curves differ between individuals because they start at different predictions.
  - A simple solution is to centre the curves at a certain point in the feature and display only the difference in the prediction to this point.



# DETECTION OF INTERACTION BETWEEN VARIABLES

## H-Statistics

### ■ Interaction Measures (H-Statistics)

- In case of interaction prediction cannot be expressed as the sum of the feature effects, because the effect of one feature depends on the value of the other feature
- How to measure the level of interaction between two features?

→ Have a look at **H-Statistic**, below the main idea :

- If two features do not interact, we can decompose the partial dependence function

$$PD_{age,power}(age,power) = PD_{age}(age) + PD_{power}(power)$$

- Measure the difference between the observed partial dependence function and the decomposed one without interactions.

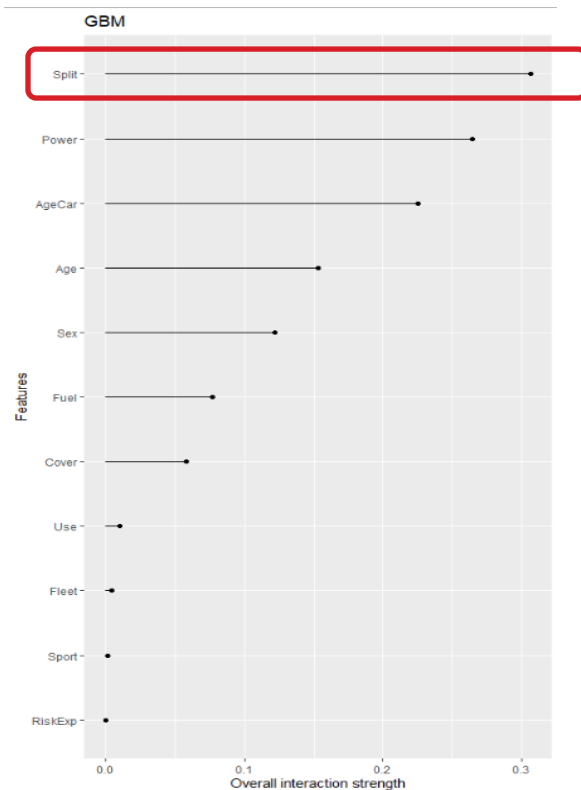
$$H^2 = \frac{\sum_{i=1}^n [PD_{age,power}(age^i, power^i) - PD_{age}(age^i) - PD_{power}(power^i)]^2}{\sum_{i=1}^n PD_{age,power}^2(age^i, power^i)}$$

- *H is 0 if there is no interaction at all*
- *A value H of 1 between two features means that each single PD function is constant and the effect on the prediction only comes through the interaction.*

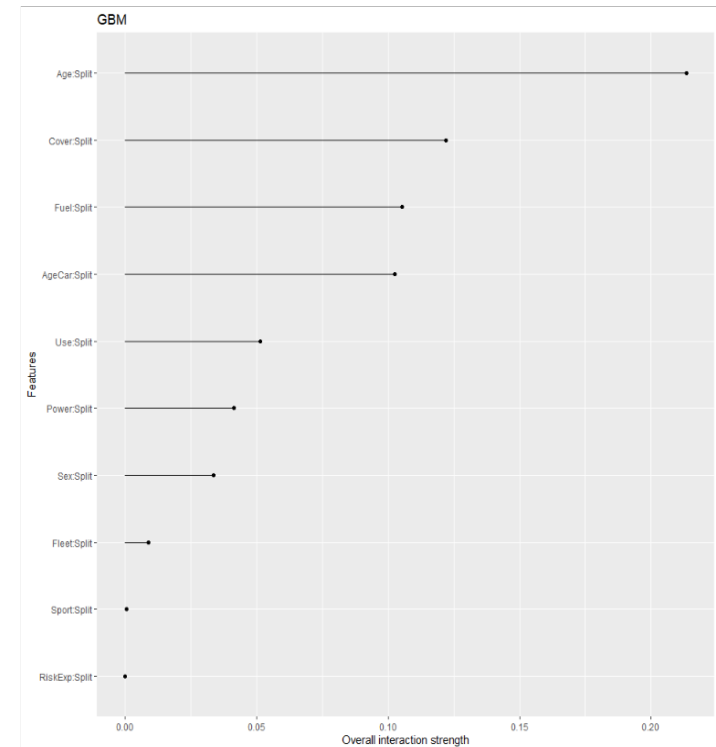
# DETECTION OF INTERACTION BETWEEN VARIABLES

## Interaction measure

(H-statistic) for each feature with all other features



2-way interactions between the cover and the other features



As the computation of the interactions is time consuming, it is better to:

- First, highlight the feature(s) interacting the most with the all other features
- For this (these) specific feature(s), to draw the 2-way interactions

### ■ H-Statistics can be used as a **features' interaction identification tool**

- It allows to identify features strongly interacting with other features
- It can then be used for **features engineering** (e.g. creating a new feature as an interaction between 2 features)

# GLOBAL VS LOCAL INTERPRETABILITY OF ML TECHNIQUES

## ■ Local Interpretability for a Single Prediction

### ○ Why did the model make a certain prediction for an instance?

- If you look at an individual prediction, the behavior of the otherwise complex model might behave more pleasantly.
- You can **zoom in on a single instance** and examine what the model predicts for this input and explain why.
  - Shapley Value
  - Breakdown

## ■ Local Interpretability for a Group of Predictions

### ○ Why did the model make specific predictions for a group of instances?

- Model predictions for multiple instances can be explained either with global model interpretation methods or with explanations of individual instances.
- The global methods can be applied by taking the group of instances, treating them as if the group was the complete dataset, and **using the global methods with this subset**.
  - LIME (Local Interpretable Model-agnostic explanations)
  - LIVE
- The individual explanation methods can be used on each instance and then listed or aggregated for the entire group.

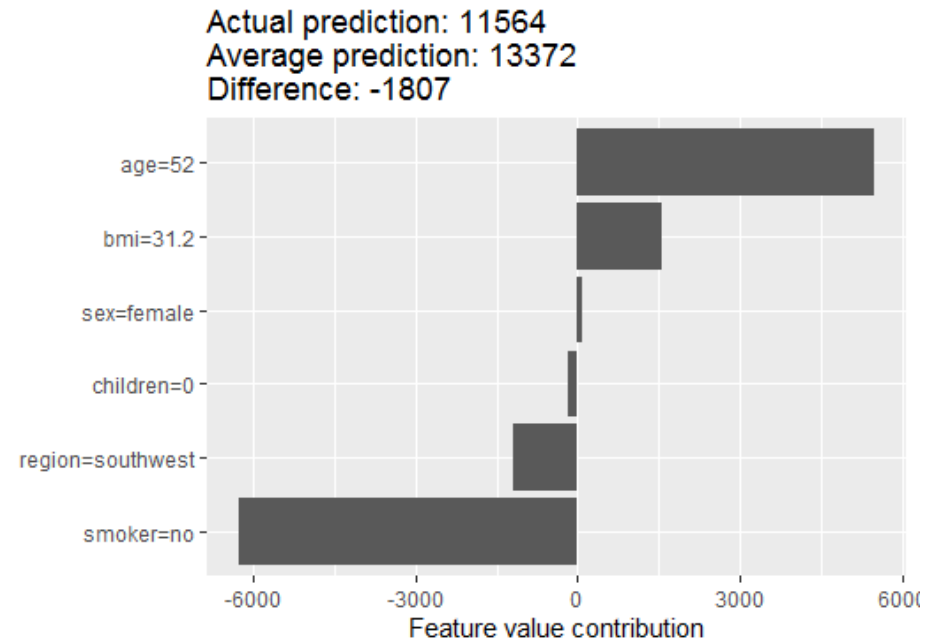
# LOCAL INTERPRETABILITY FOR A SINGLE PREDICTION

## ■ Shapley Value :

- The shapley value measures for a single prediction **how much each specific feature value will contribute to make the instance prediction different from the overall prediction**
- The computation time increases exponentially with the number of features.

### *From Game Theory*

- *The Shapley value is the average marginal contribution of a feature value across all possible coalitions (= sets composed of different number of features).*
- *For each of these coalitions we compute the prediction with and without the feature value of interest and take the difference to get the marginal contribution.*
- *The Shapley value is the (weighted) average of marginal contributions across all the coalitions.*



# AGENDA

A non-exhaustive reminder to some useful ML techniques

Adding complexity means increasing need for interpretability

An introduction to ML interpretation tools

**Conclusions: how to make the most of ML techniques**

# HOW TO MAKE THE MOST OF ML TECHNIQUES IN INSURANCE ANALYTICS?

## Two different strategies

1. Replacing traditional models (e.g. GLM) by ML models
2. Combining the pros of traditional and ML models to improve existing ones

## Replacing traditional models by ML models

- The main drawback of this approach is the black-box effect of the ML results
- There is therefore a strong need in using interpretation tools
  - **Features importance** to select the most relevant variable (e.g. if we have too many features available and/or we want to limit the number of features)
  - **PDP, ICE, ALE and/or H-Statistics** to understand the impact of the selected variables on the prediction and identify the potential interactions
  - **Shapley value** to better understand the prediction on specific data points



# HOW TO MAKE THE MOST OF ML TECHNIQUES IN INSURANCE ANALYTICS?

## Combining traditional and ML models

- ML methods would then be used to perform features extraction, features selection and/or features engineering
  - **Feature extraction** = reducing the dimensionality of too voluminous datasets (in terms of # features)
  - **Feature selection** = selecting the most relevant variables to our problem
  - **Feature engineering** = identifying the best representation of the sample data to learn a solution to your problem (e.g. interactions)
- The selected/engineered variables could then be introduced in our usual model (e.g. GLM) to obtain easily interpretable results combined with more insights

## REFERENCES

- Christoph Molnar, “Interpretable Machine Learning: a guide to making black box models explainable”. Available: <https://christophm.github.io/interpretable-ml-book/>
- Michaël Lecuivre and Samuel Mahy, “Machine Learning Interpretability: A toolbox to better understand your ML results – With application in insurance pricing”, Reacfin White Paper Vol.1 2022 – February 2022. Available: <https://www.reacfin.com/wp-content/uploads/2016/12/20220224-Reacfin-White-Paper-Machine-Learning-Interpretability-1-1.pdf>
- Jean Dessain, Fabien Vinas, Nora Bentaleb: “Cost of Explainability in AI: an Example with Credit Scoring Models” (to be published in November). Presentation available on <https://www.reacfin.com/wp-content/uploads/2016/12/2023-07-26-XAI-Conference-Cost-of-explainability-with-credit-scoring-JDE-presentation-v2.0.pdf>

## CONTACT DETAILS



**Xavier Maréchal**

CEO – Managing Partner

M +32 497 48 98 48

[xavier.marechal@reactfin.com](mailto:xavier.marechal@reactfin.com)



# Reactfin

Place de l'Université 25

B-1348 Louvain-la-Neuve (Belgium)

T +32 (0) 10 68 86 07

[www.reactfin.com](http://www.reactfin.com)

**DISCLAIMER**

*The recipient of this document should treat all information as strictly confidential and only in the context stated below. Information may not be disclosed to any third party without the prior joint-consent of Reacfin.*

*Estimates given in this presentation are based on our current knowledge, they can be based upon our previous experience within the Undertaking, as well as taking into account similar projects in the same context as the Undertaking, either locally, within majority of the EU countries as well as overseas.*

*This presentation is only the supporting document of a verbal presentation. Hence, it is not intended to be exhaustive. Quoting or using this document on its own might be misleading. As a result, these materials may not be used by anybody except their authors nor should they be relied upon in any way for any purpose other than as contemplated by joint written agreement with Reacfin.*

**Reacfin**

Place de l'Université 25  
B-1348 Louvain-la-Neuve

[www.reacfin.com](http://www.reacfin.com)