

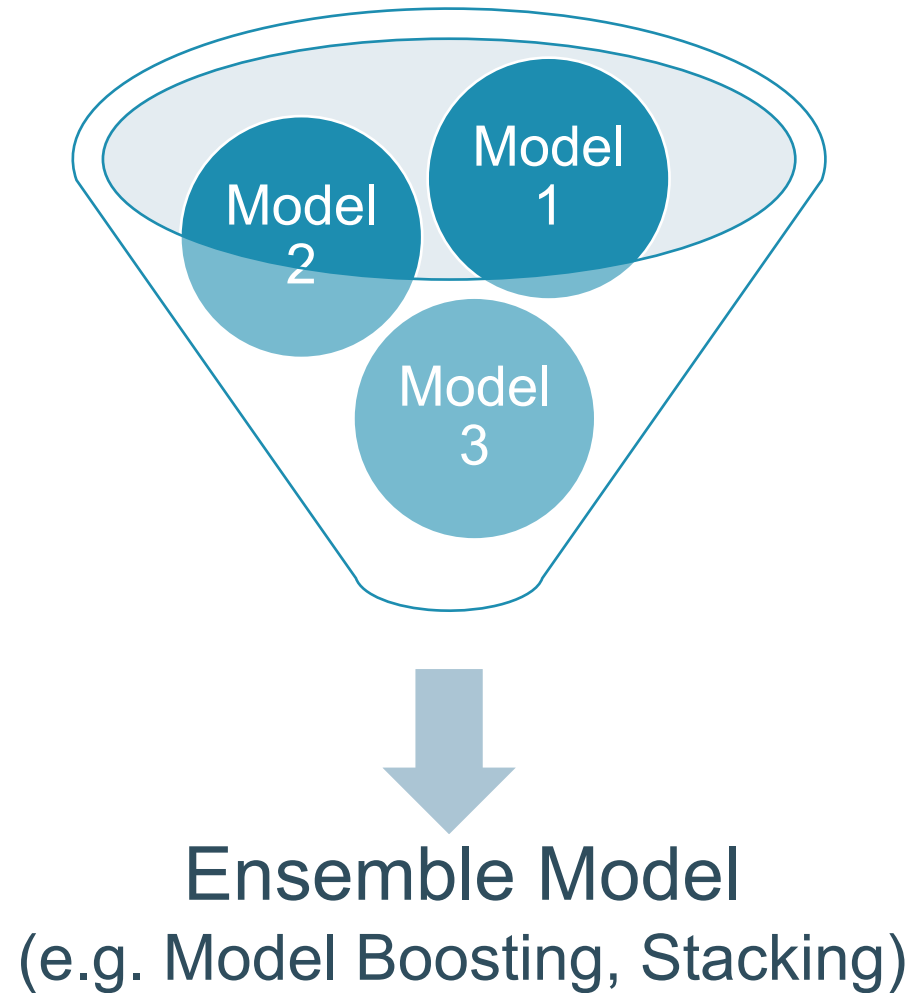
Model boosting and stacking for insurance pricing

Master thesis of Ine Fransen

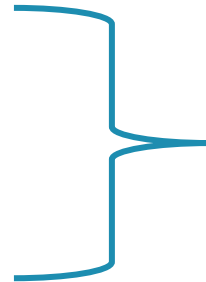
Master of Actuarial and Financial Engineering

Prof. Dr. K. Antonio, Ir. R. Henckaerts

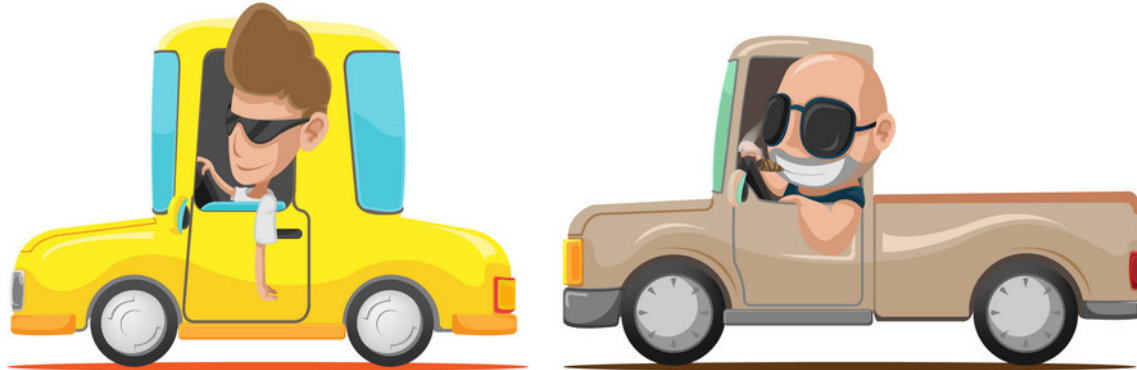
Ensemble approaches combine different models



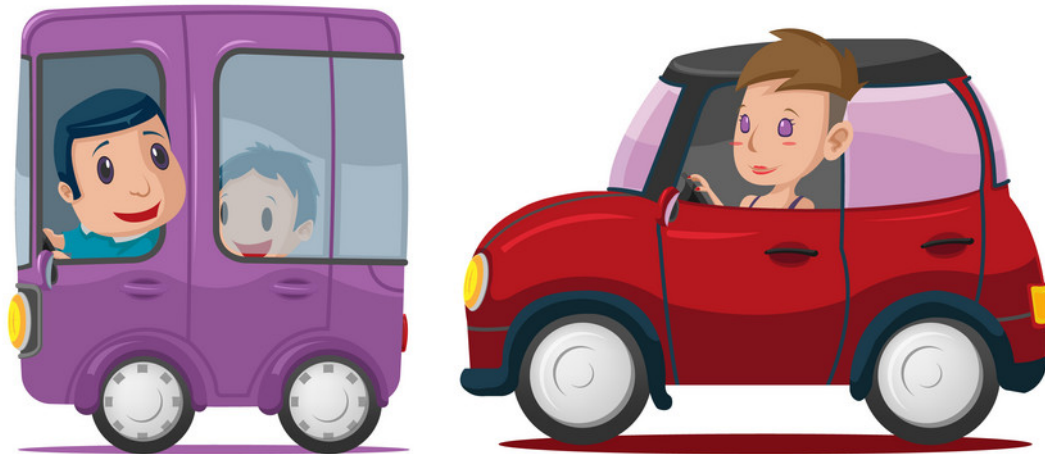
Goals of the thesis

- Present the concept of **model boosting**
 - Present the concept of **stacking**
 - Compare model performance and interpretability
- 
- Apply to car insurance pricing

Portfolio in non-life insurance

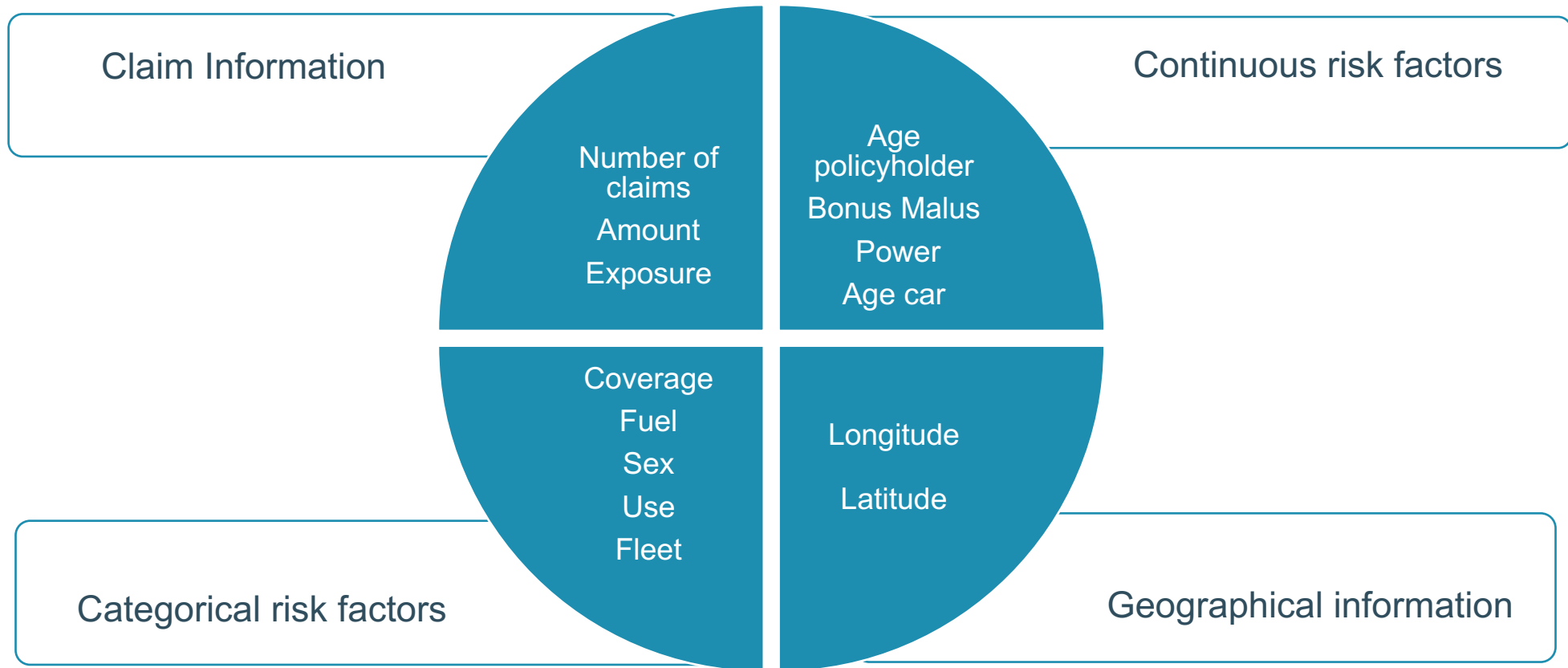


Construction of a pricing model



<https://www.vectorstock.com/royalty-free-vector/cars-driver-cartoon-collection-set-vector-14066666>

Motor Third Party Liability (MTPL) data 1997



Calculation of the premium

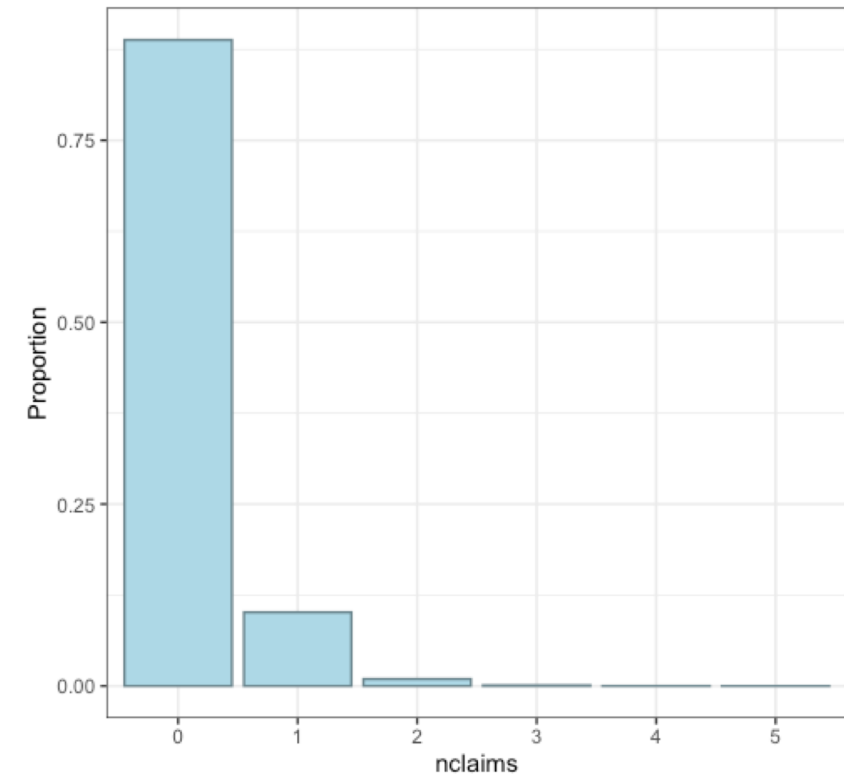
$$\pi_i = \mathbb{E}[F_i] \times \mathbb{E}[S_i]$$

Frequency

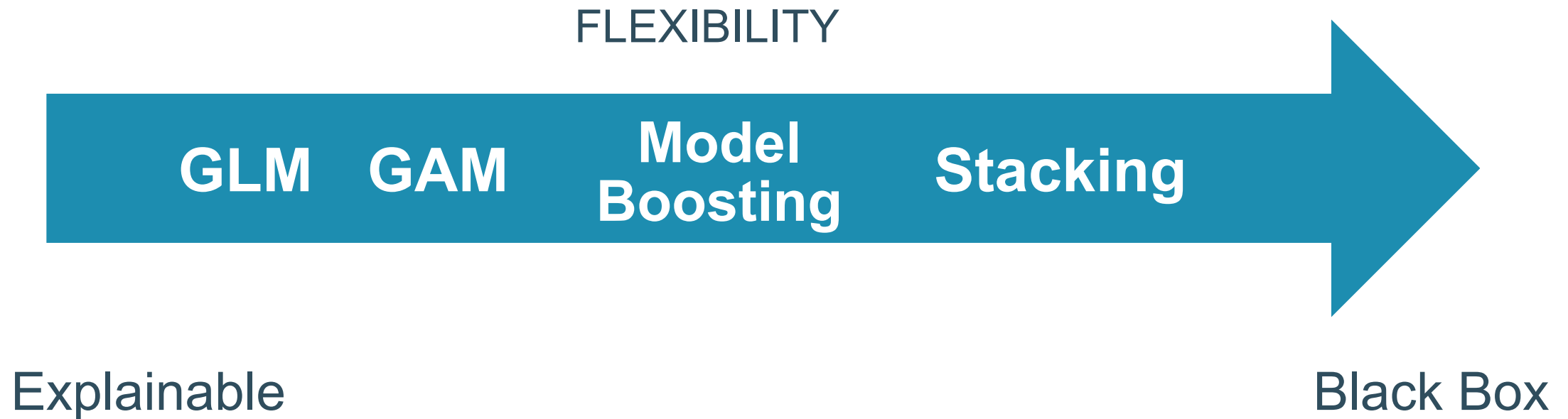
Severity

Challenges

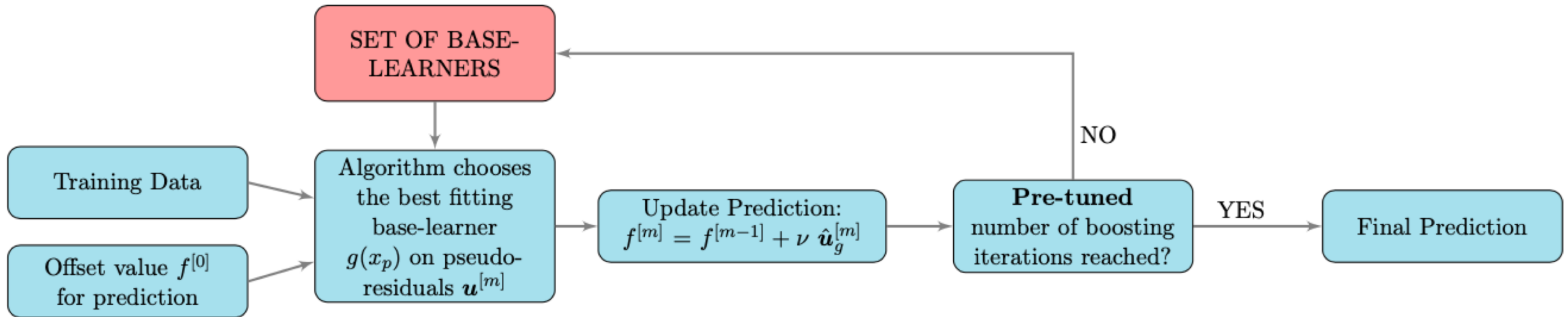
- Selection of relevant variables
- Different types of risk factors
- Interpretability of the pricing model
- Distribution of the target variable



Demand for flexible yet explainable models



Model Boosting

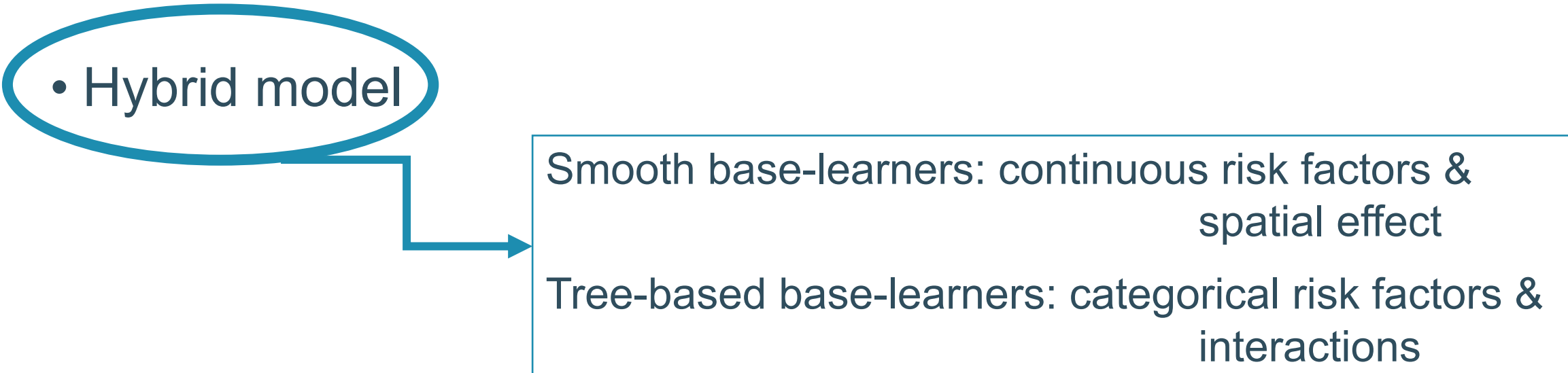


Model Boosting

- Base-learners:
 - **Linear** effects: e.g. βsex
 - **Smooth** effects: e.g. $\sum_{j=1}^t B_j(\text{age}, q)$
 - **Tree-based** effects: e.g. $\sum_{j=1}^{J_1} \hat{y}_{R_j} \mathbb{I}(\text{age} \in R_j)$
- Automatic variable selection!

Three boosting models for claim frequency

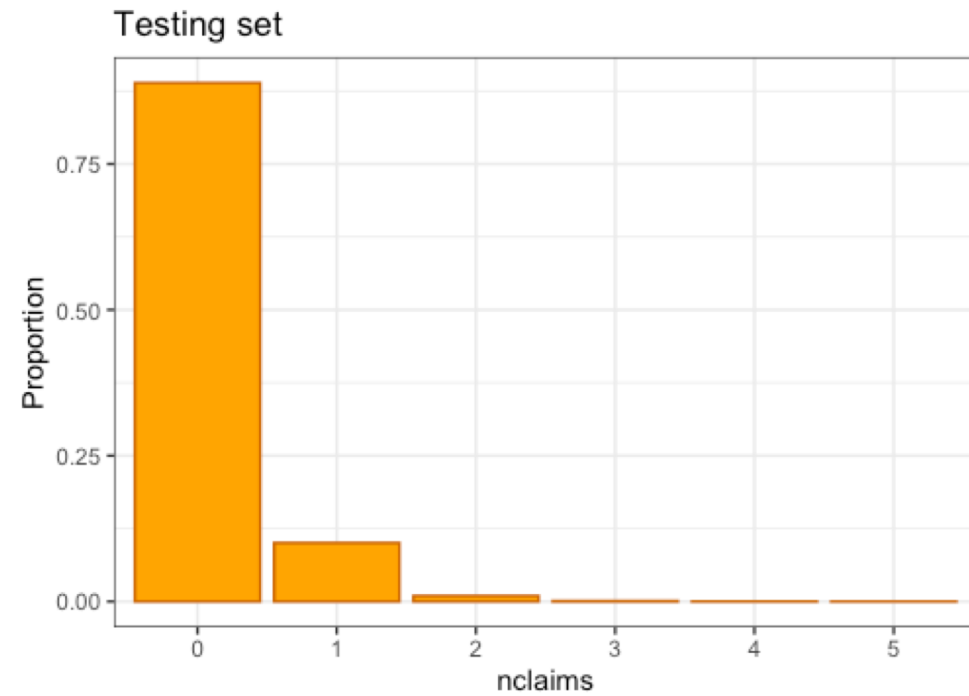
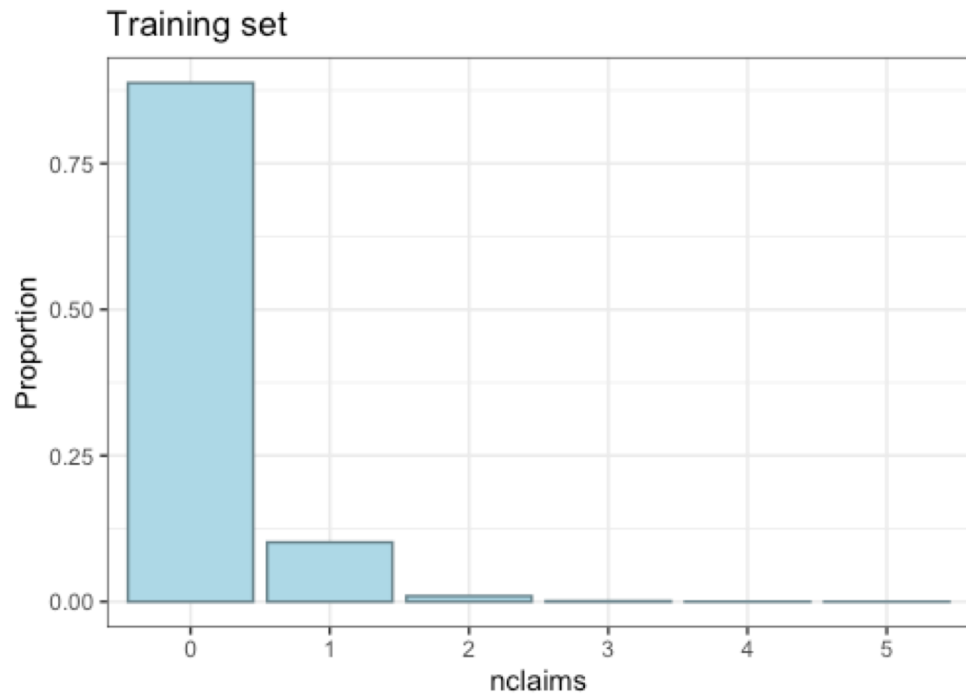
- Smooth model
- Tree model
- Hybrid model



Smooth base-learners: continuous risk factors & spatial effect

Tree-based base-learners: categorical risk factors & interactions

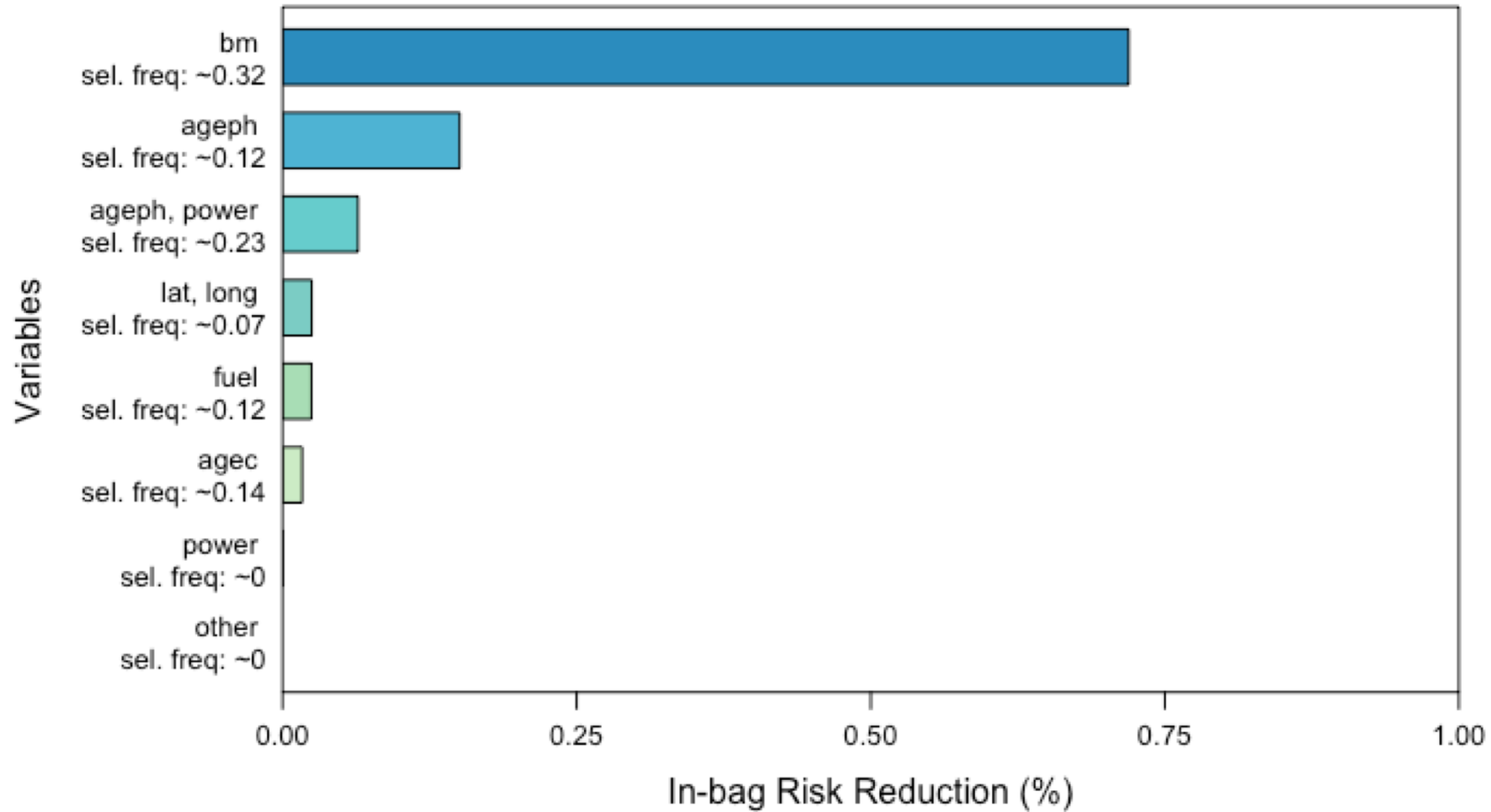
Training and test set



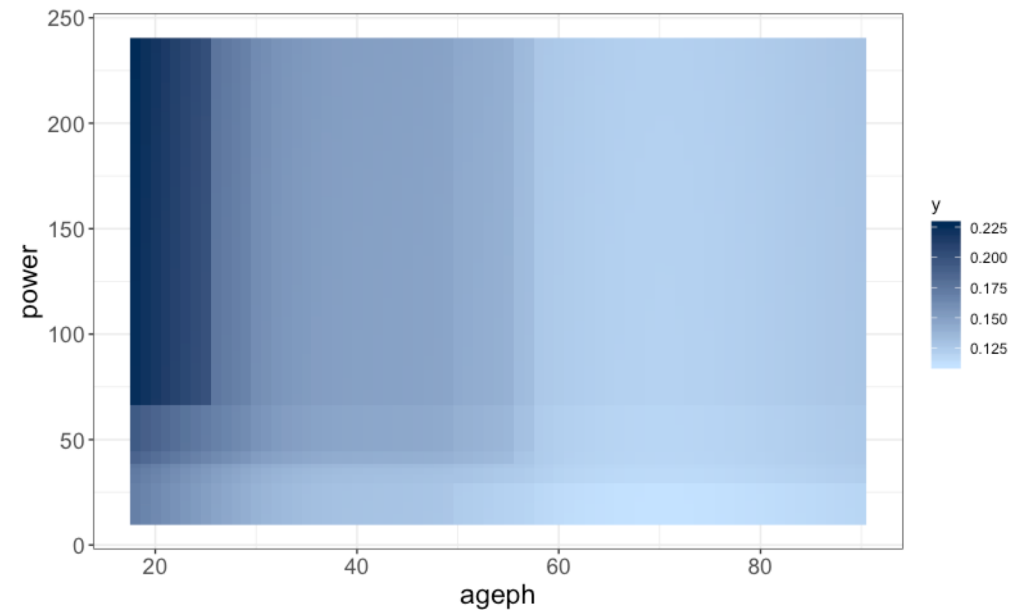
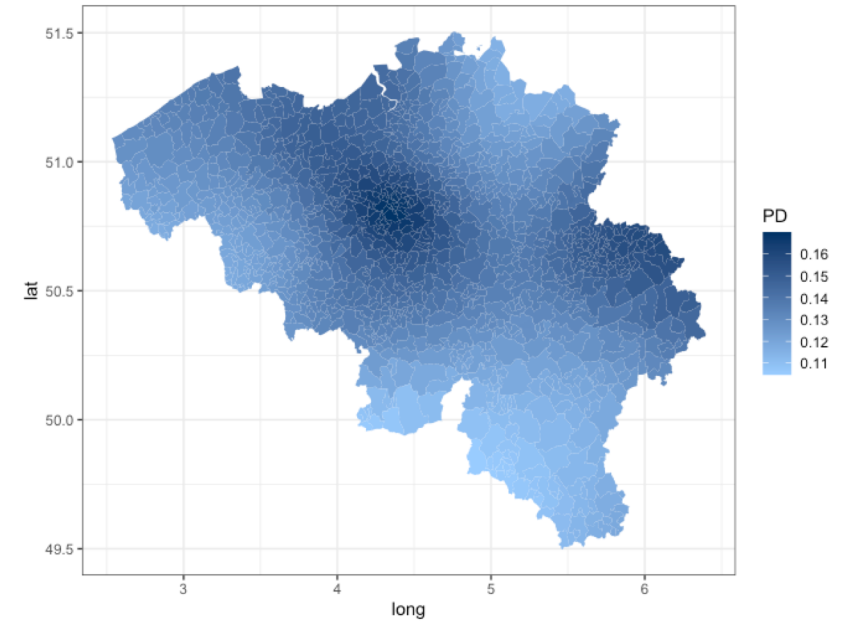
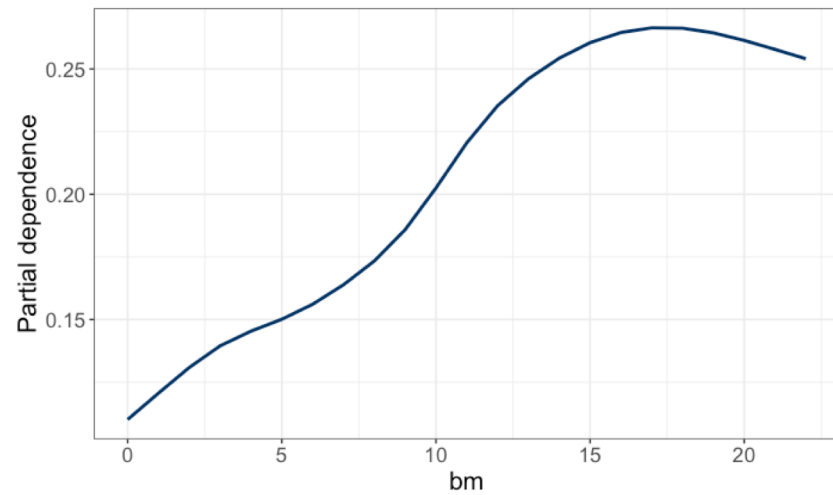
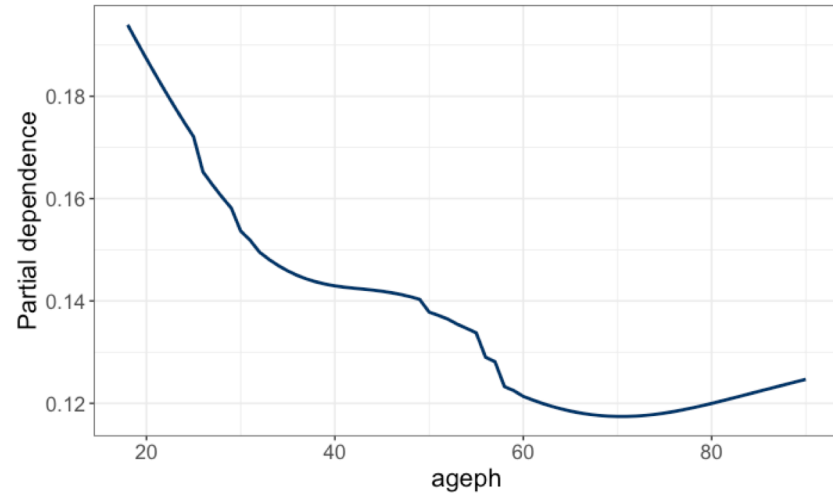
Comparison with GLM

- Boosting can be more accurate
- Boosting is computationally costly
- Interpretability of boosting can be increased using
 - variable importance
 - partial dependence plots

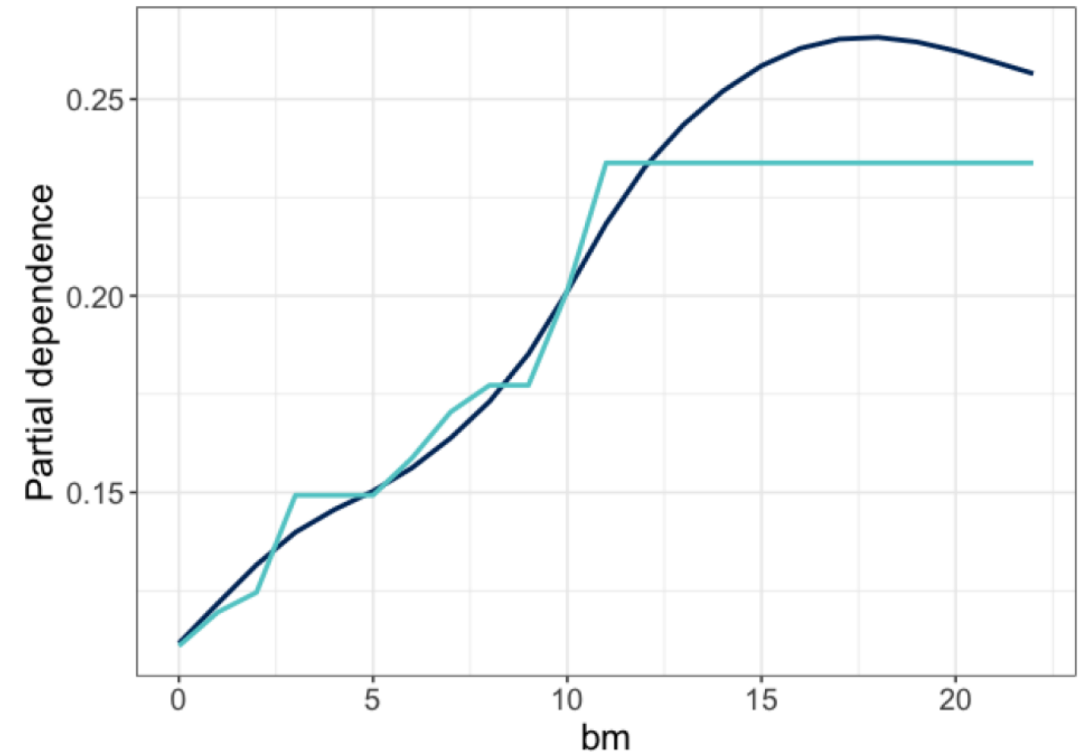
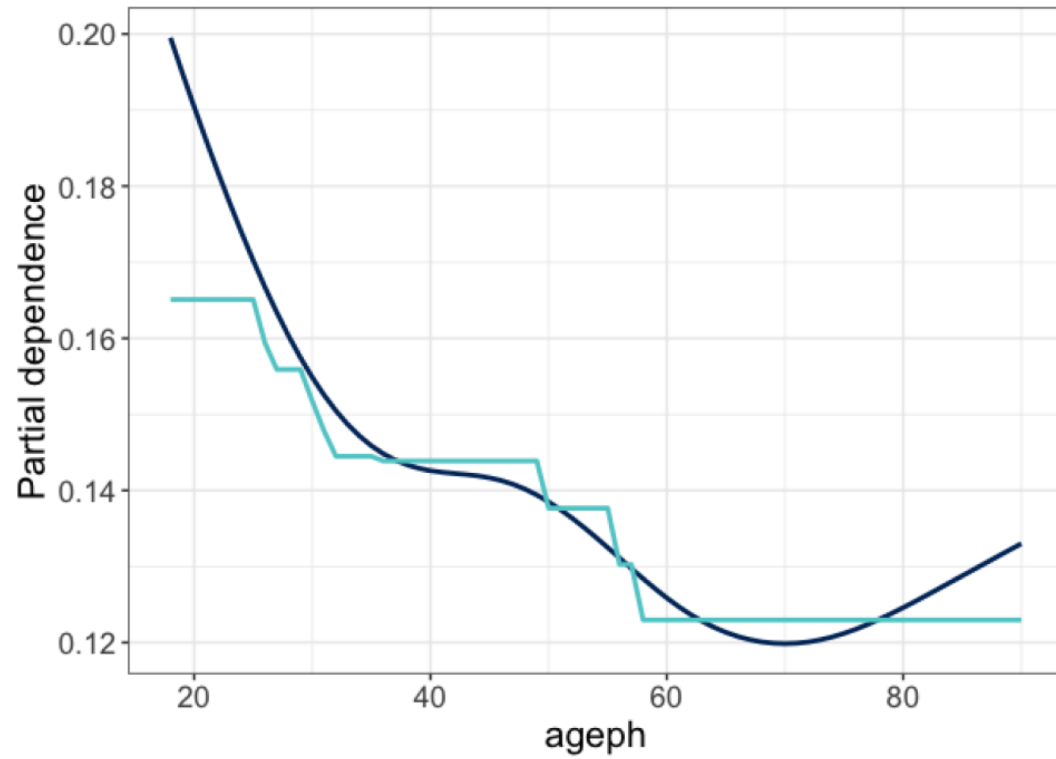
Variable importance



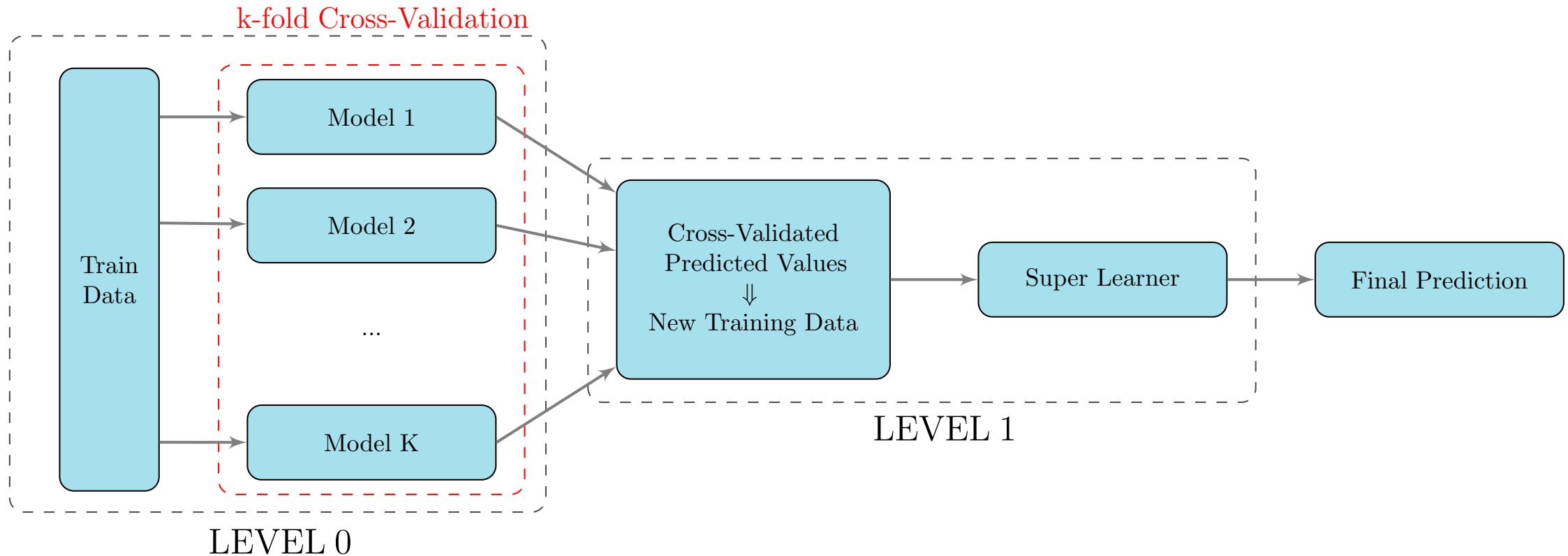
Partial dependence plots



Approach for automatic binning



Stacking



Four different models are stacked

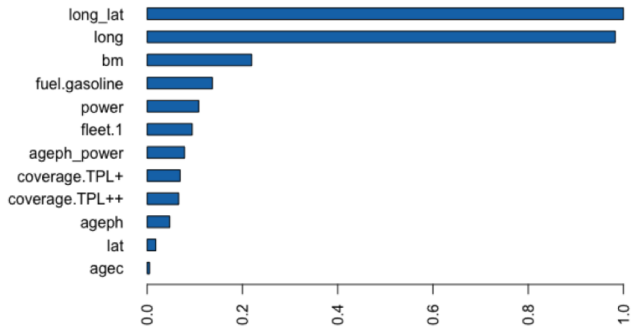
- GLM
- GBM
- XGBoost
- Random Forest

Comparison

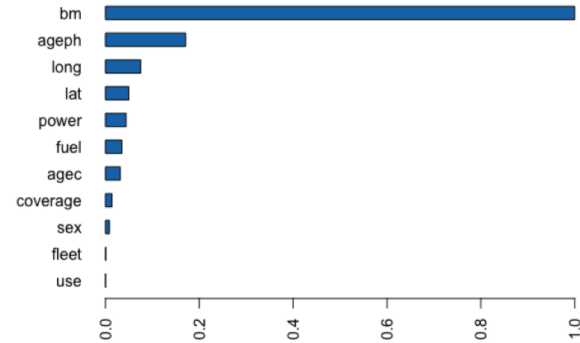
- Better accuracy than single models (GLM, Boosting models,...)
- Reduced interpretability
 - Variable ~~importance~~
 - Partial dependence plots
- Increased computation time
- Overkill?

Variable importance

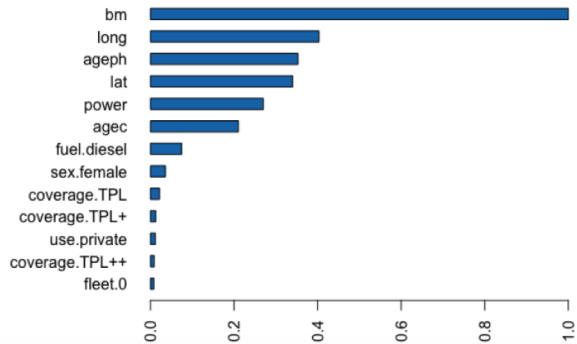
Variable Importance: GLM



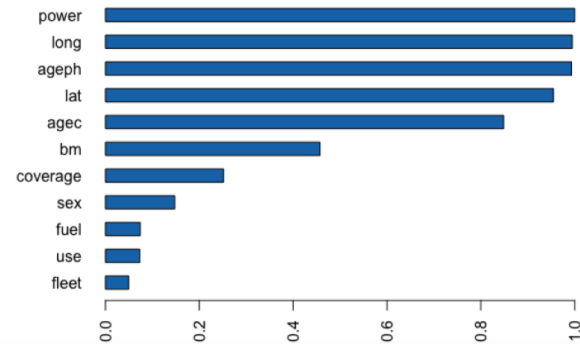
Variable Importance: GBM



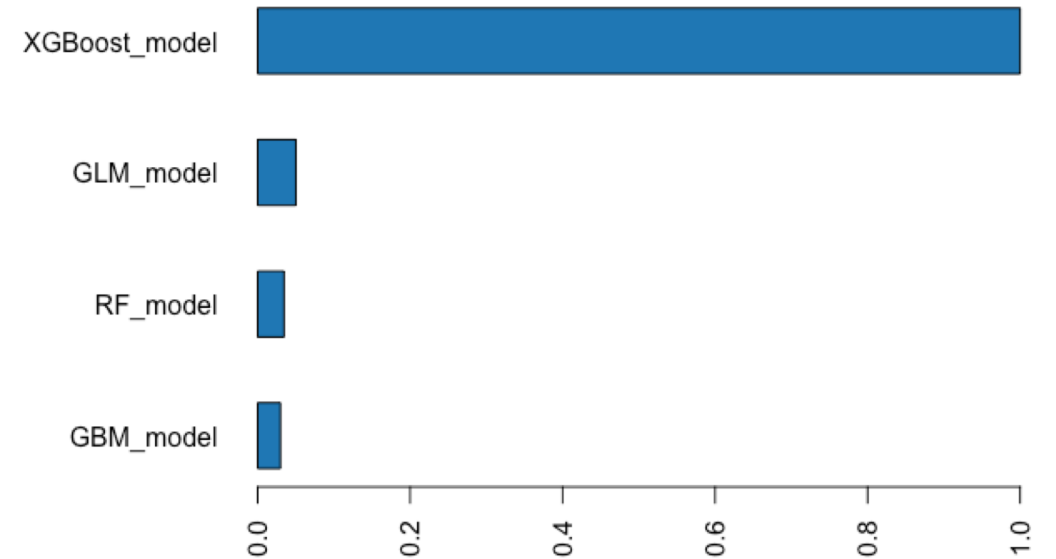
Variable Importance: XGBOOST



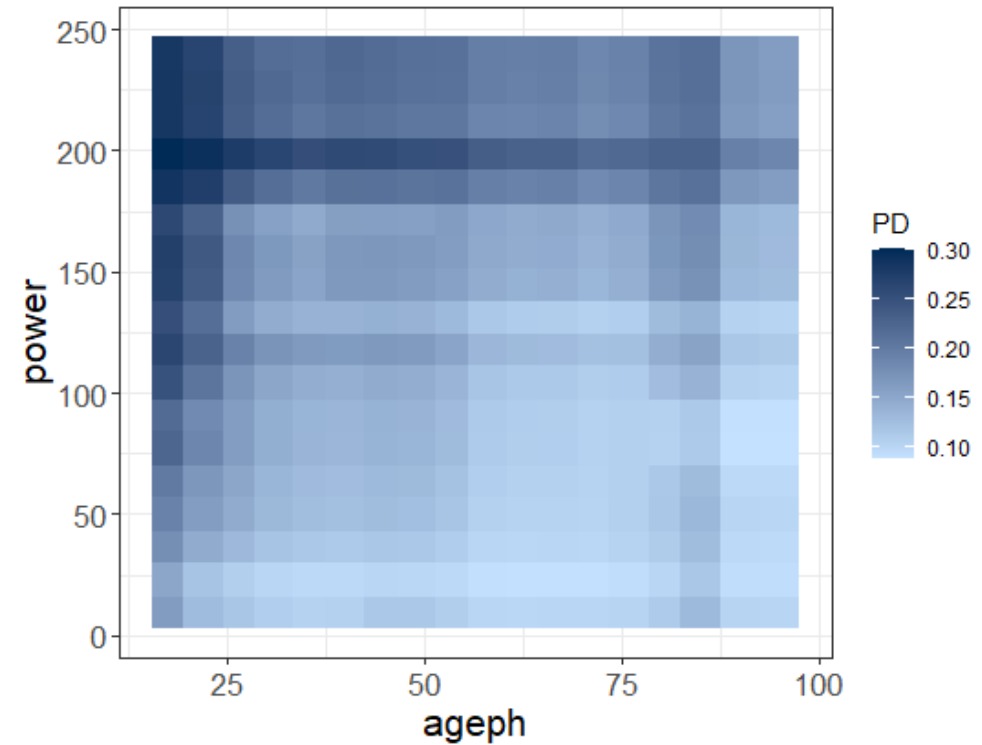
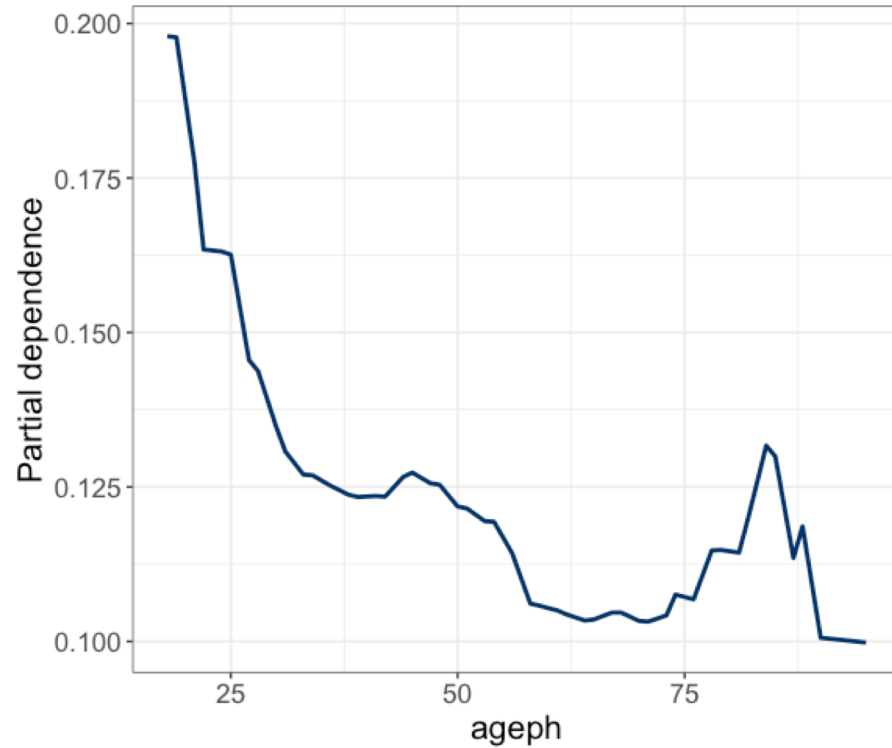
Variable Importance: DRF



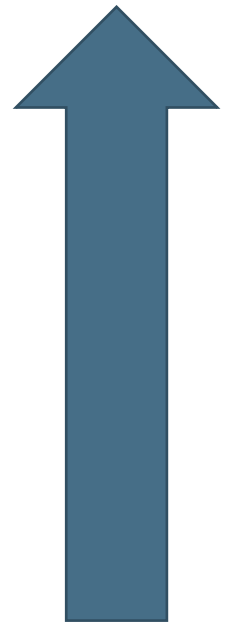
Variable Importance



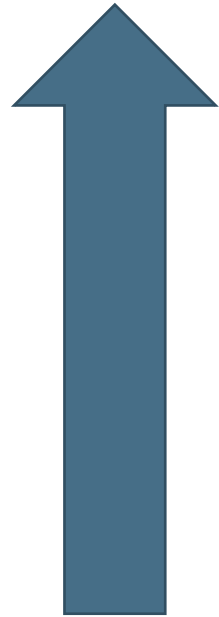
Partial dependence plots



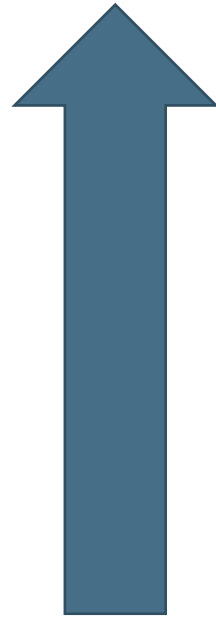
Conclusion



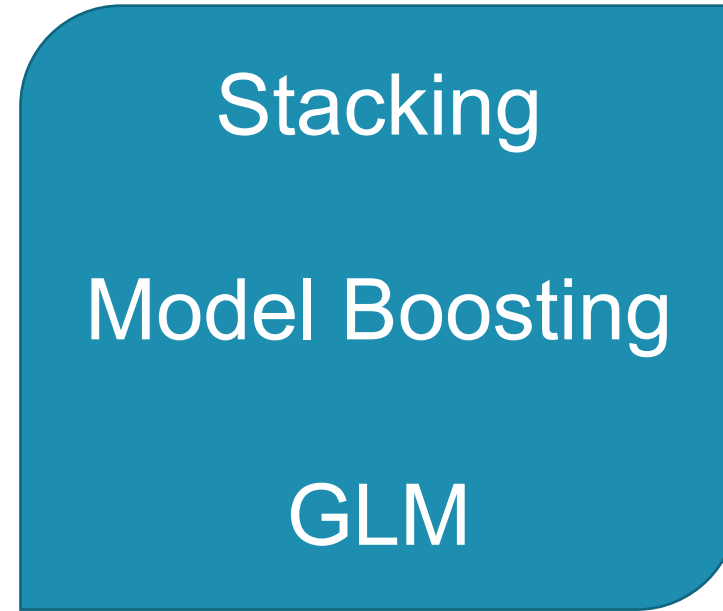
NEED FOR
COMPUTATIONAL
POWER



ACCURACY



FLEXIBILITY



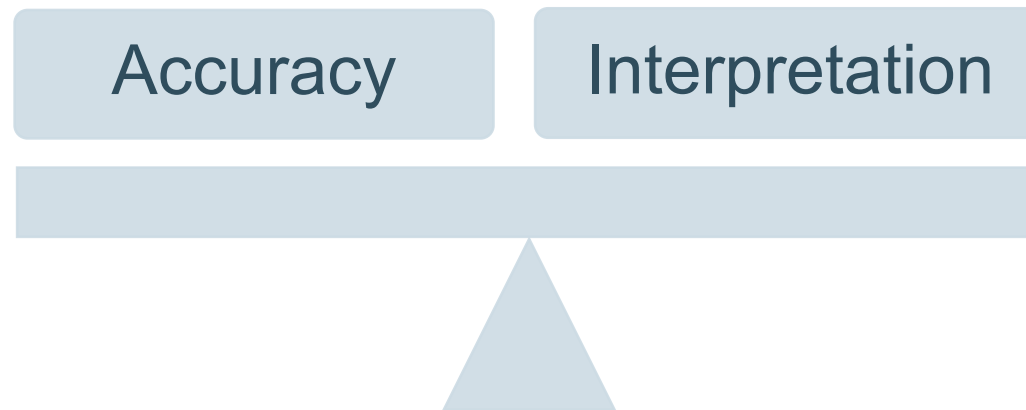
INTERPRETABILITY

Future prospects

- Increased data size and complexity
- More powerful machines



More complex models



Thank you for your attention!

