

Interpretable Machine Learning with R

Dean Allsopp

SatRday Newcastle

April 6, 2019

Why
“Interpretable”
Machine Learning?

In the beginning...

Linear Models

Let there be light...

- Coefficients
- Residual plots
- QQ plots
- Leverage plots

Kaggle Era DS...

It's all about the error score

Darkness falls...

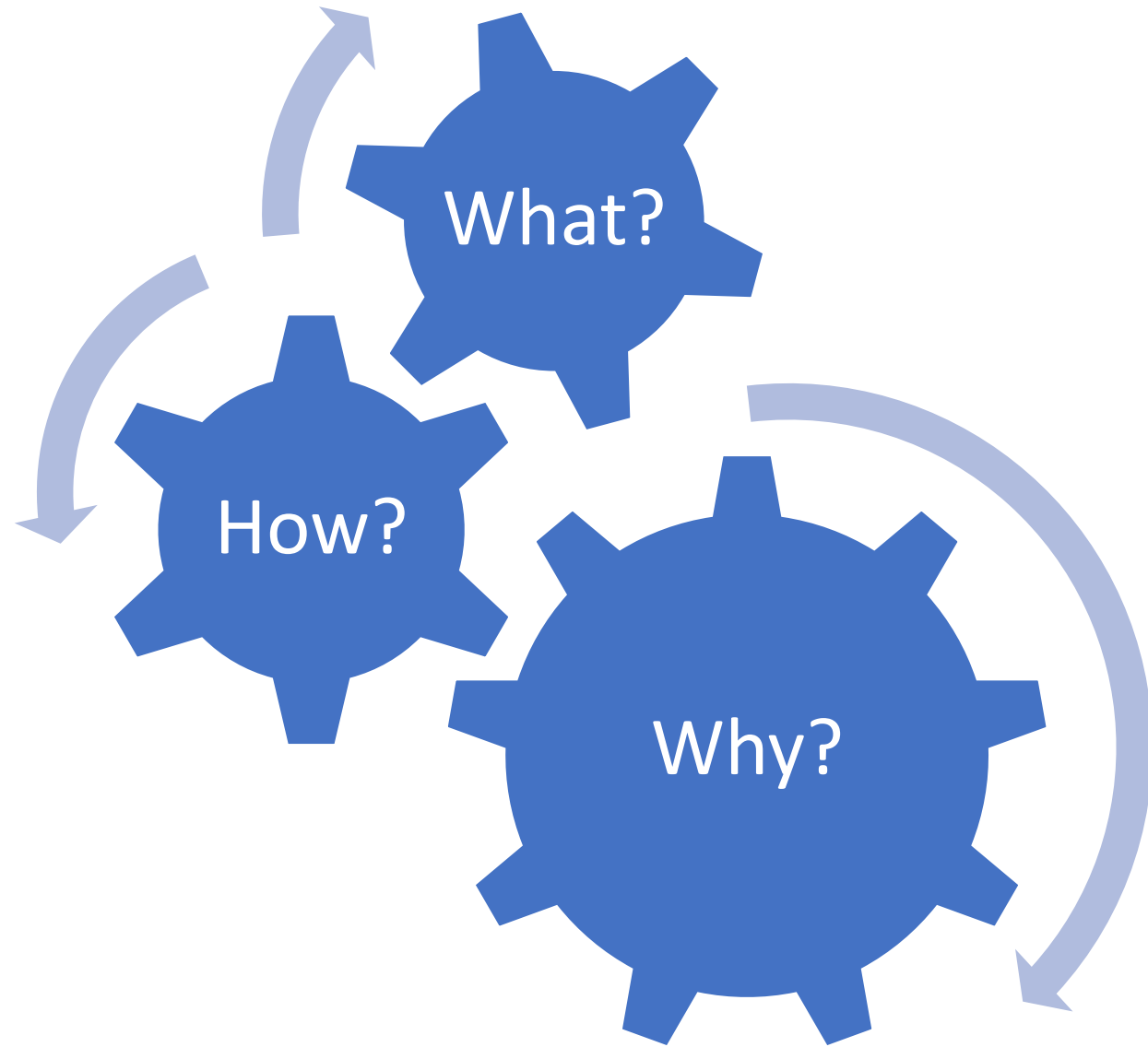
- Complex Features
- Complex Models
- Hyperparameter Search

But what is it actually doing to
predict this result?

Just hypothetically important?

What constitutes the best
medical treatment for recovery?

See Rich Caruana



What? – Data (Features/Instances)

How? – does the Model work?

Why? – do we get this Prediction?

How well do you understand
the data?

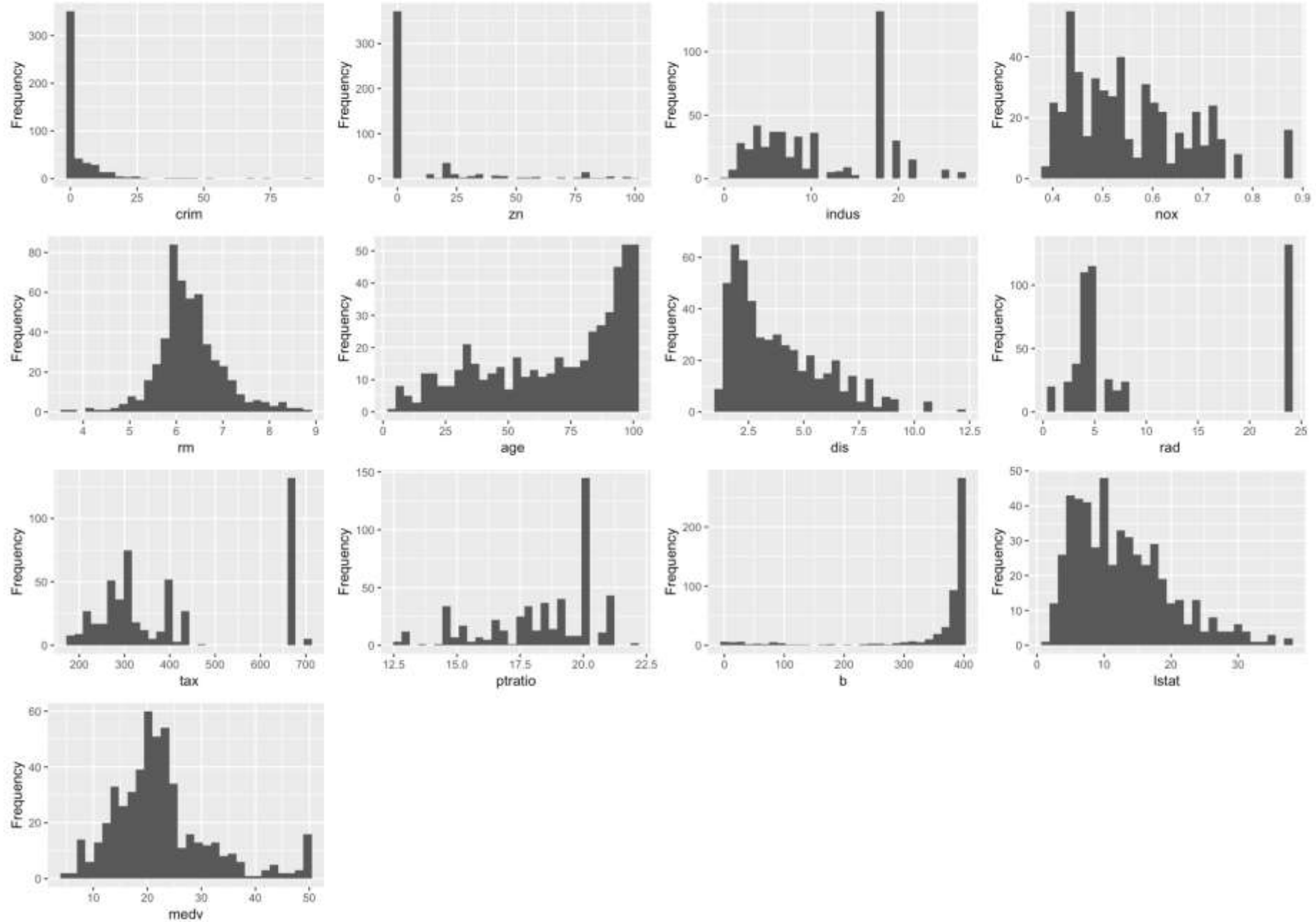


EDA

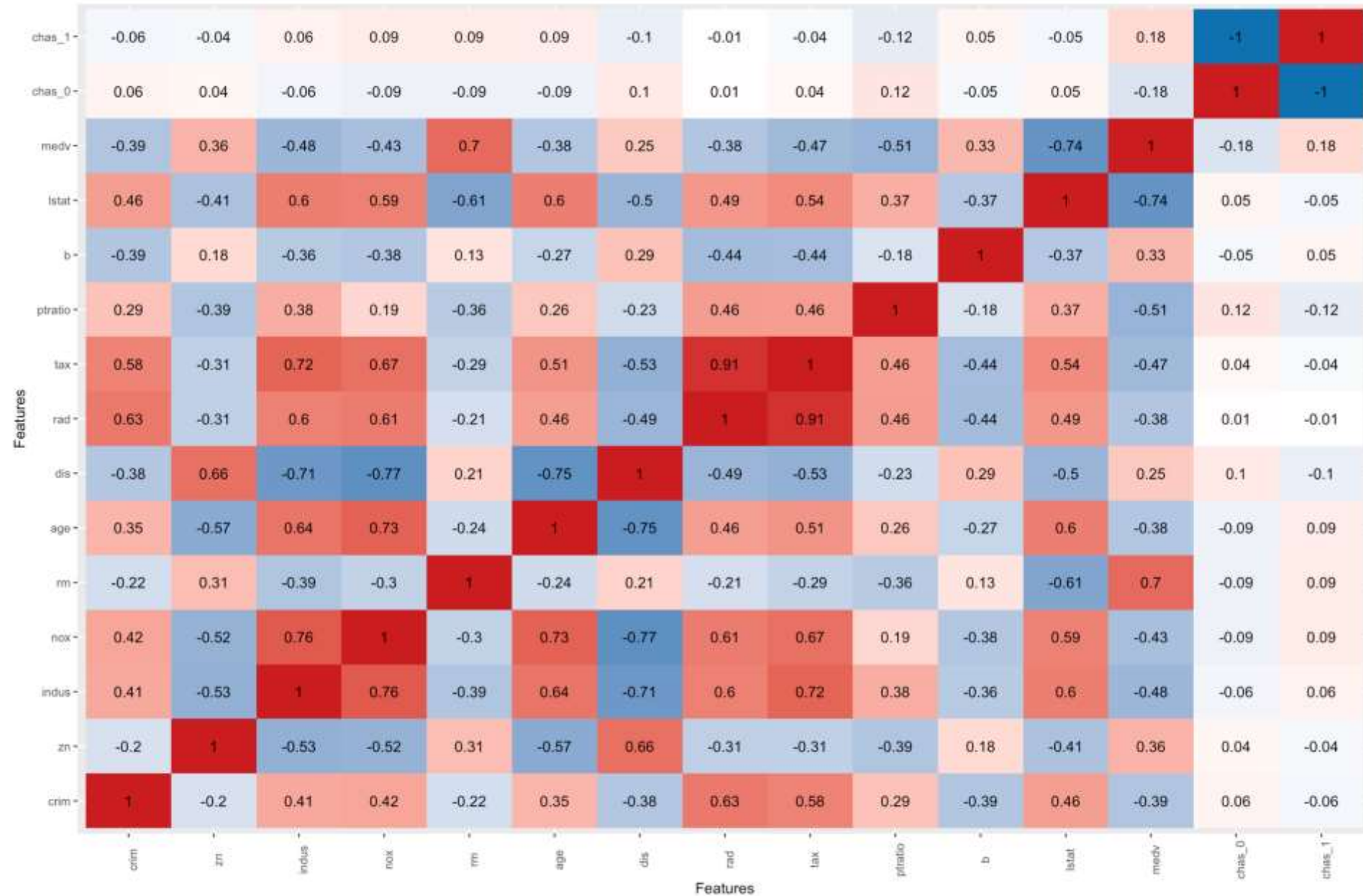
Data Explorer

Univariate Distribution

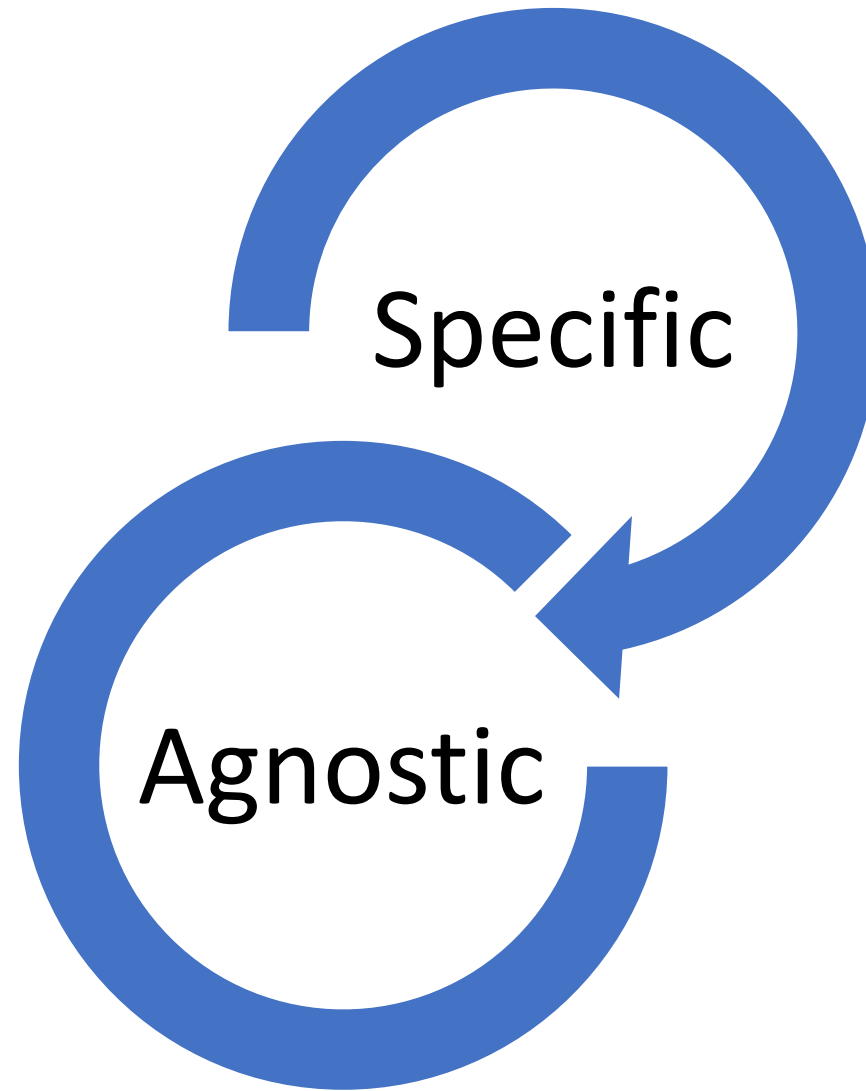
Histogram



Correlation Analysis



Kind of Model Interpretation?

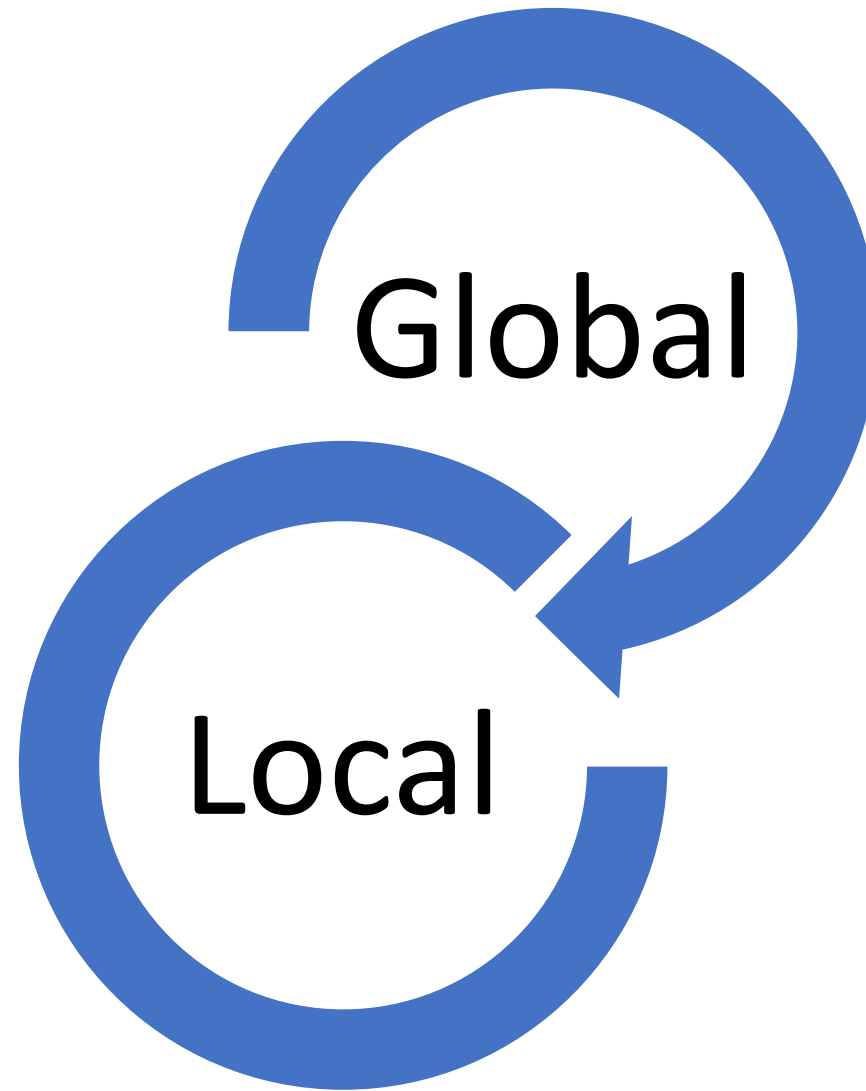


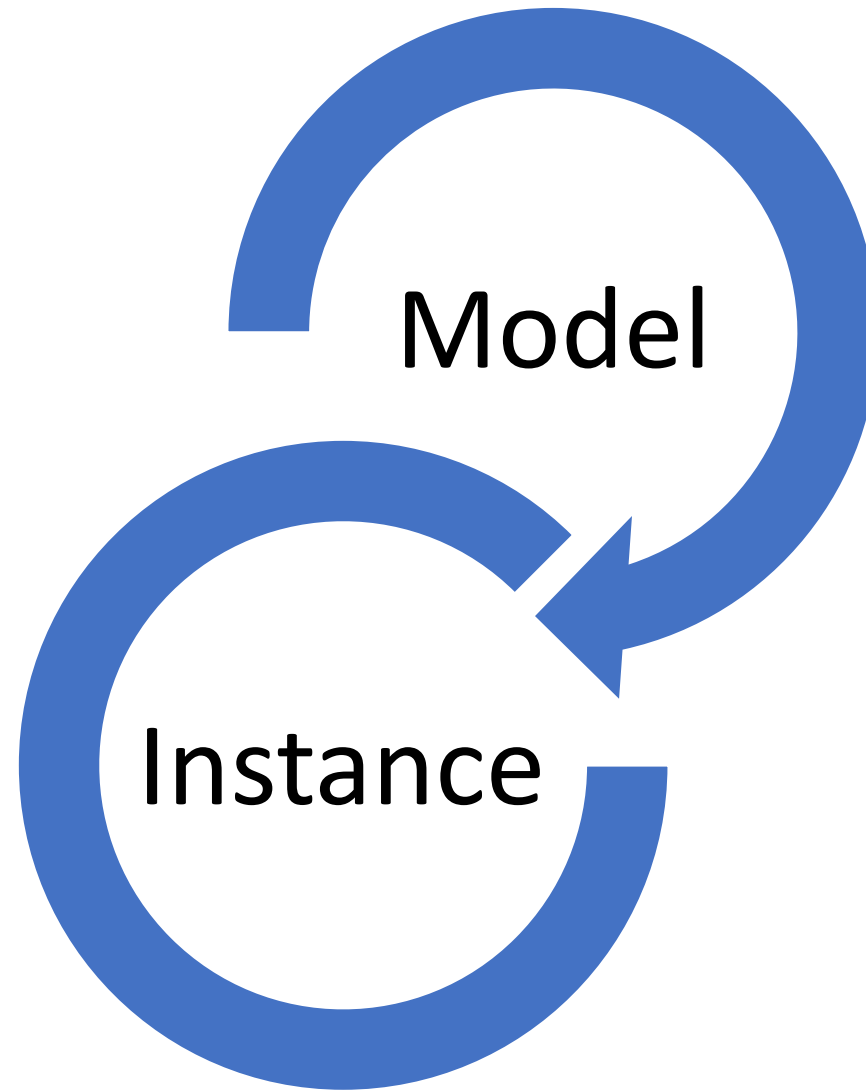
Exploit Internal Logic

Versus

System Identification Approach

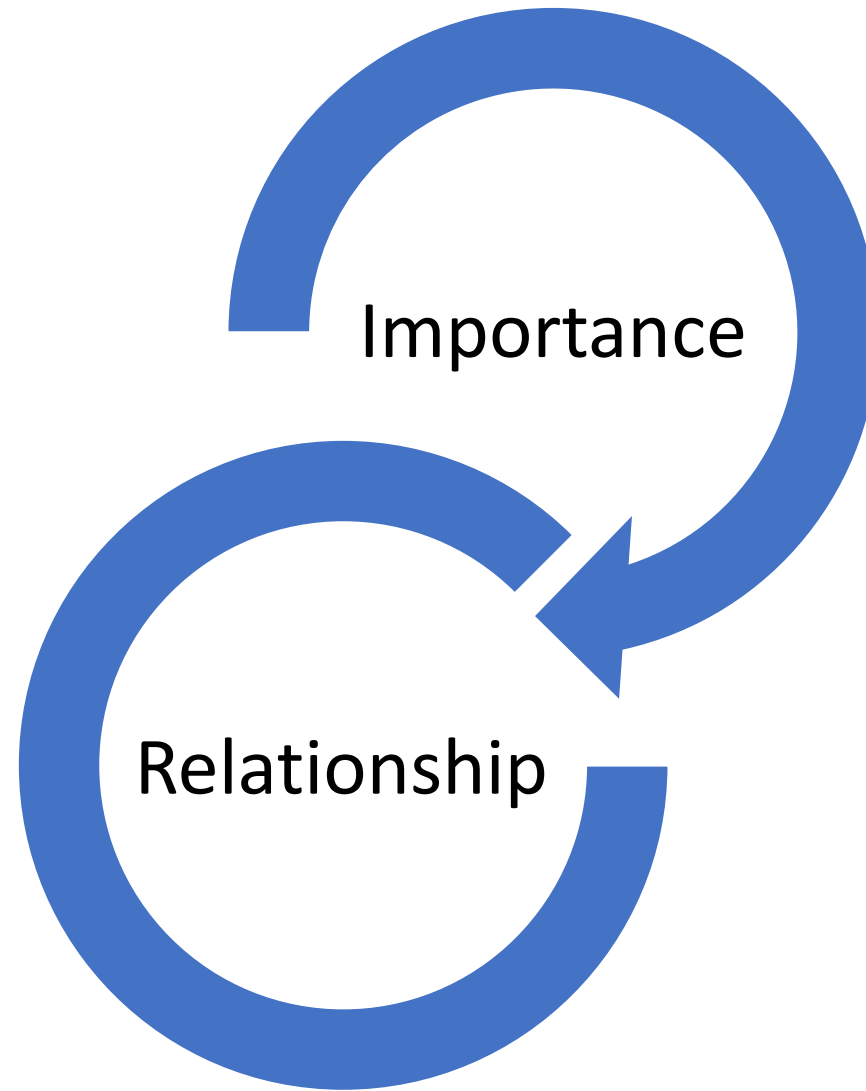
Scope of Model Interpretation?





#1: Global

What do you want to understand
about the features and the
model?



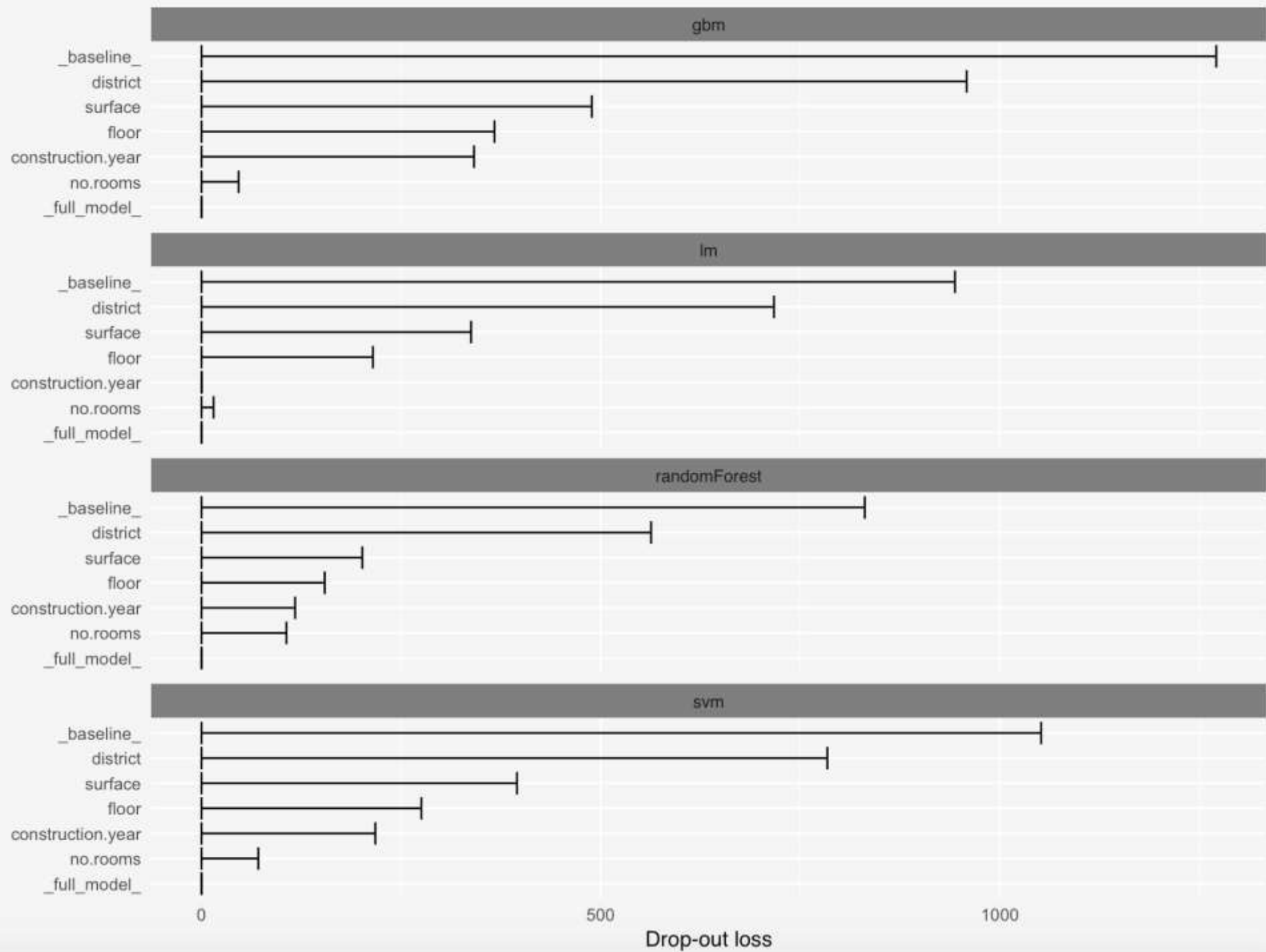
Variable Importance Plot

DALEX/VIP (agnostic)

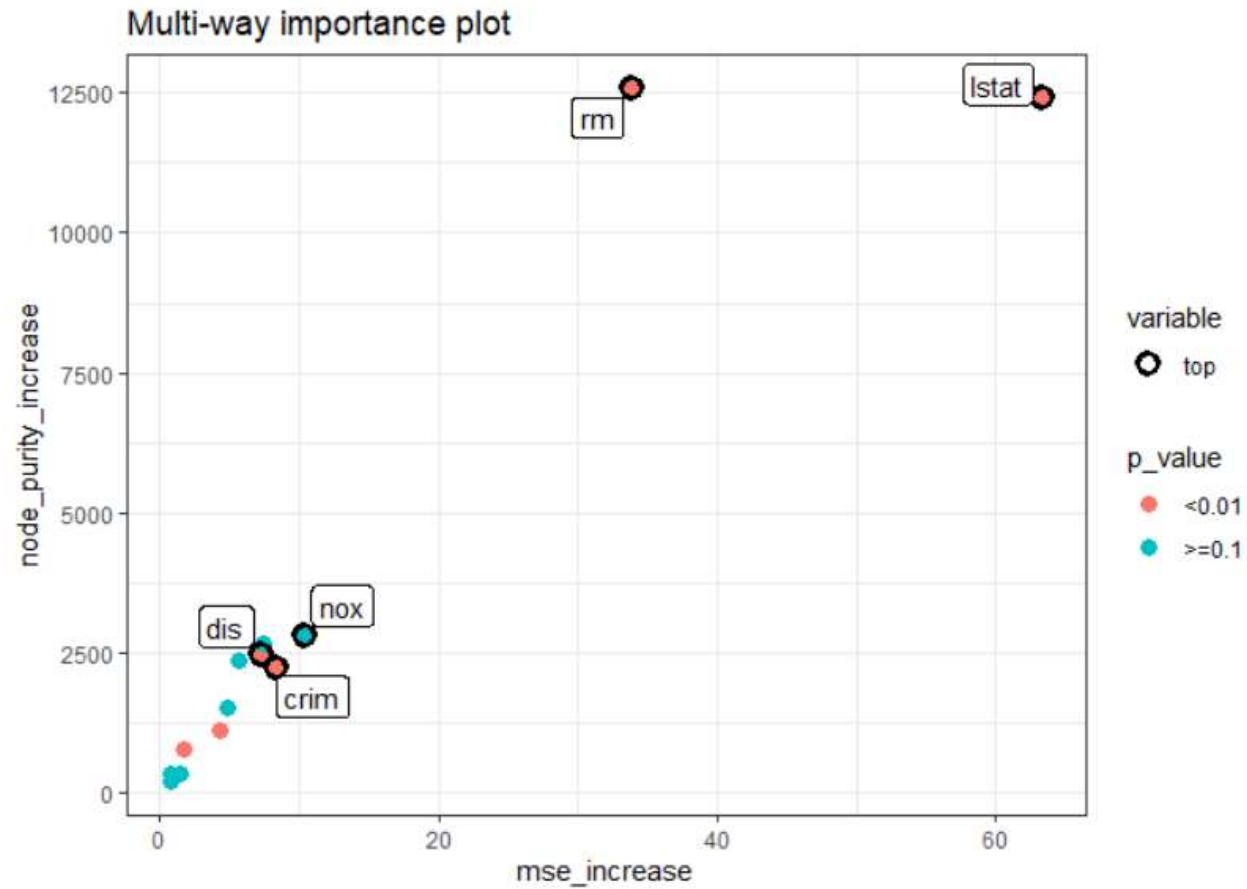
randomForestExplainer

Why highlight VIP?

Permutation

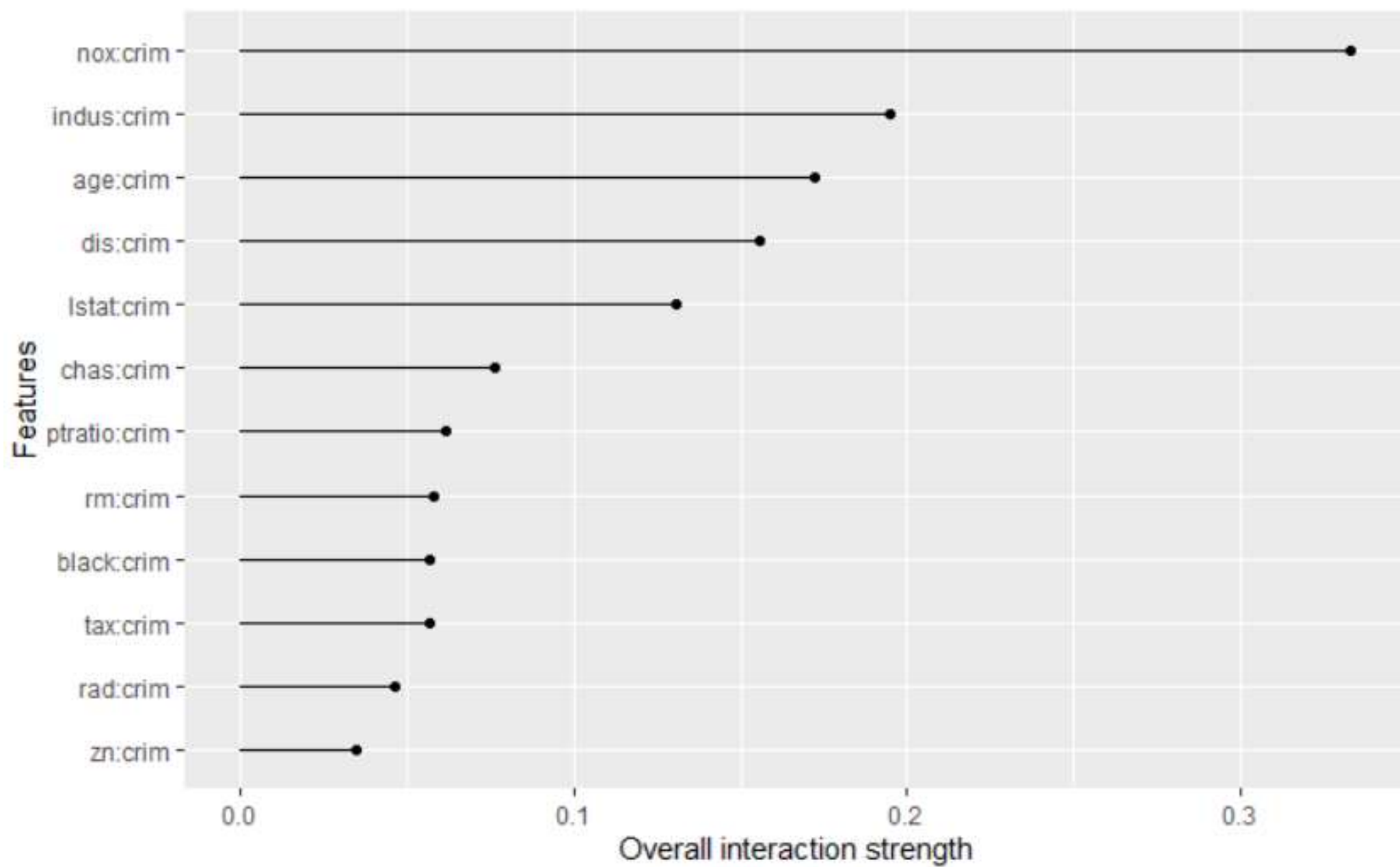


randomForestExplainer



What about Interactions?

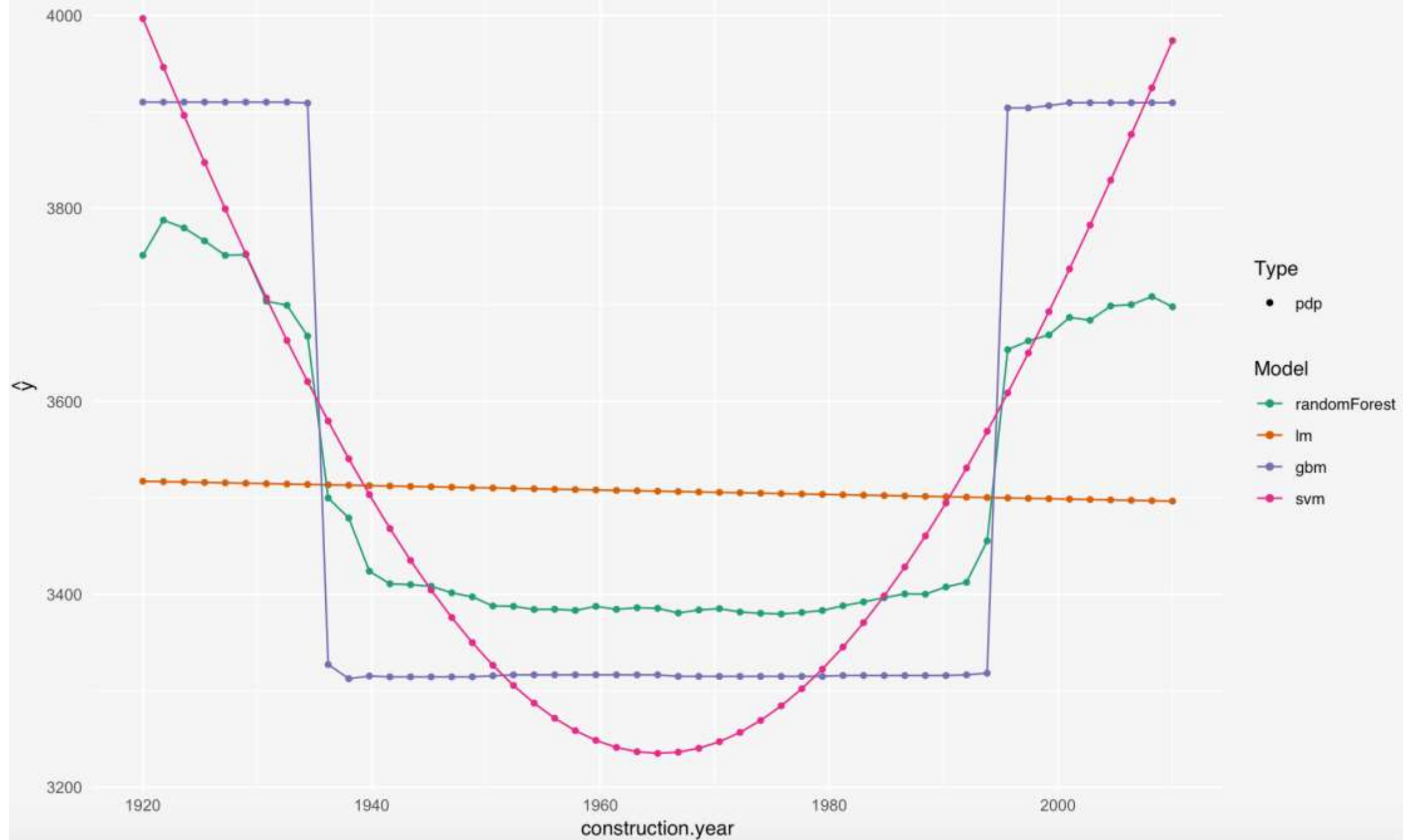
IML



Partial Dependence Plot

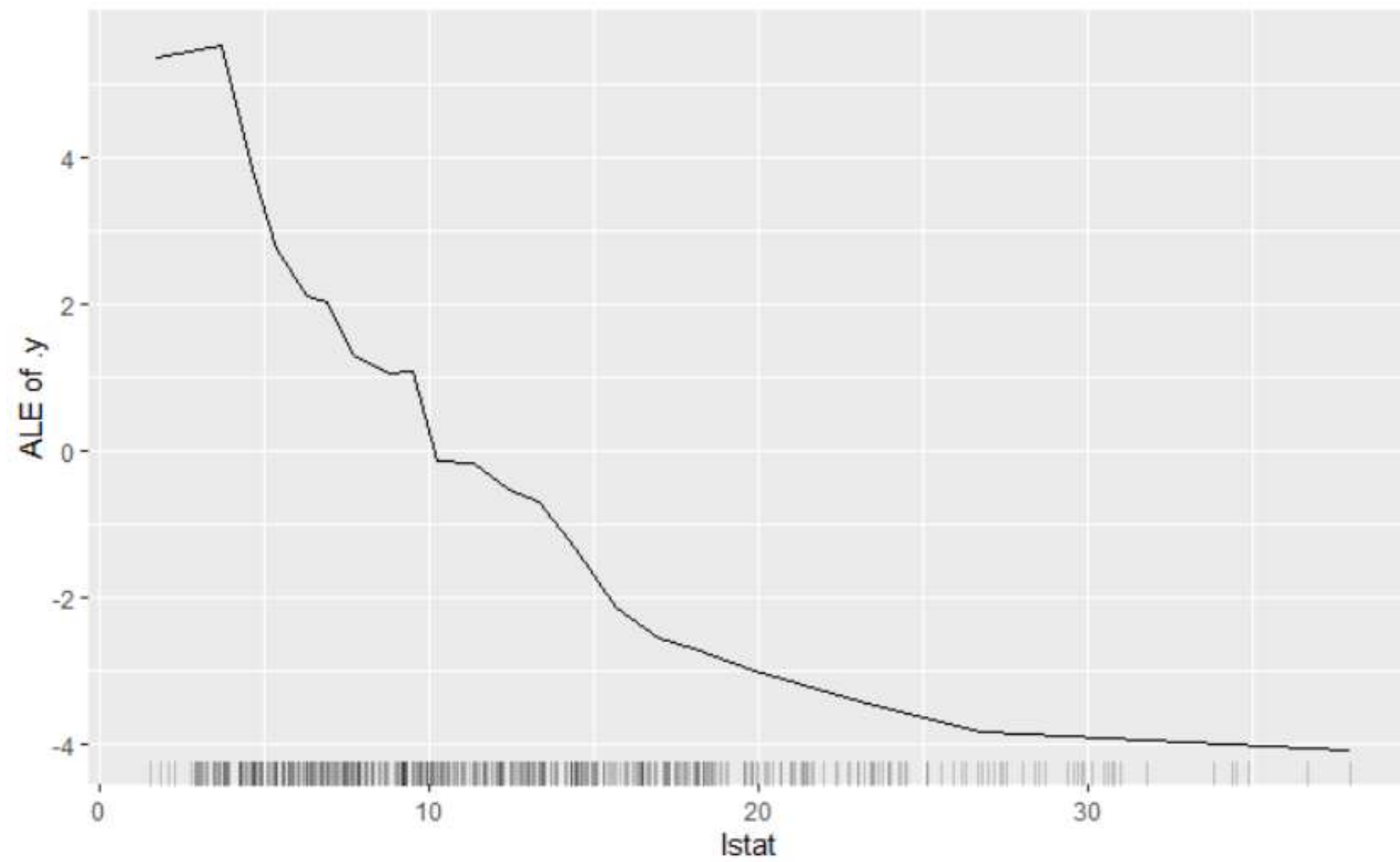
DALEX/IML/plotmo

Variable response



Accumulated Local Effects Plot

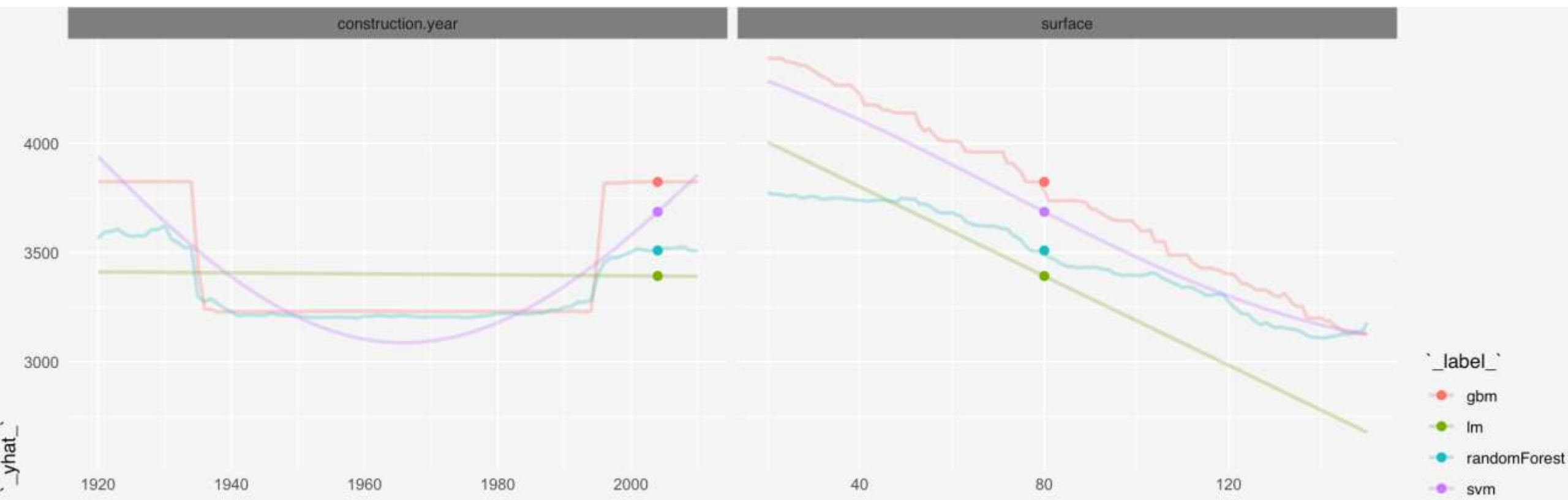
DALEX/IML



What If?

Ceteris Paribus Plot

DALEX



#2: Local

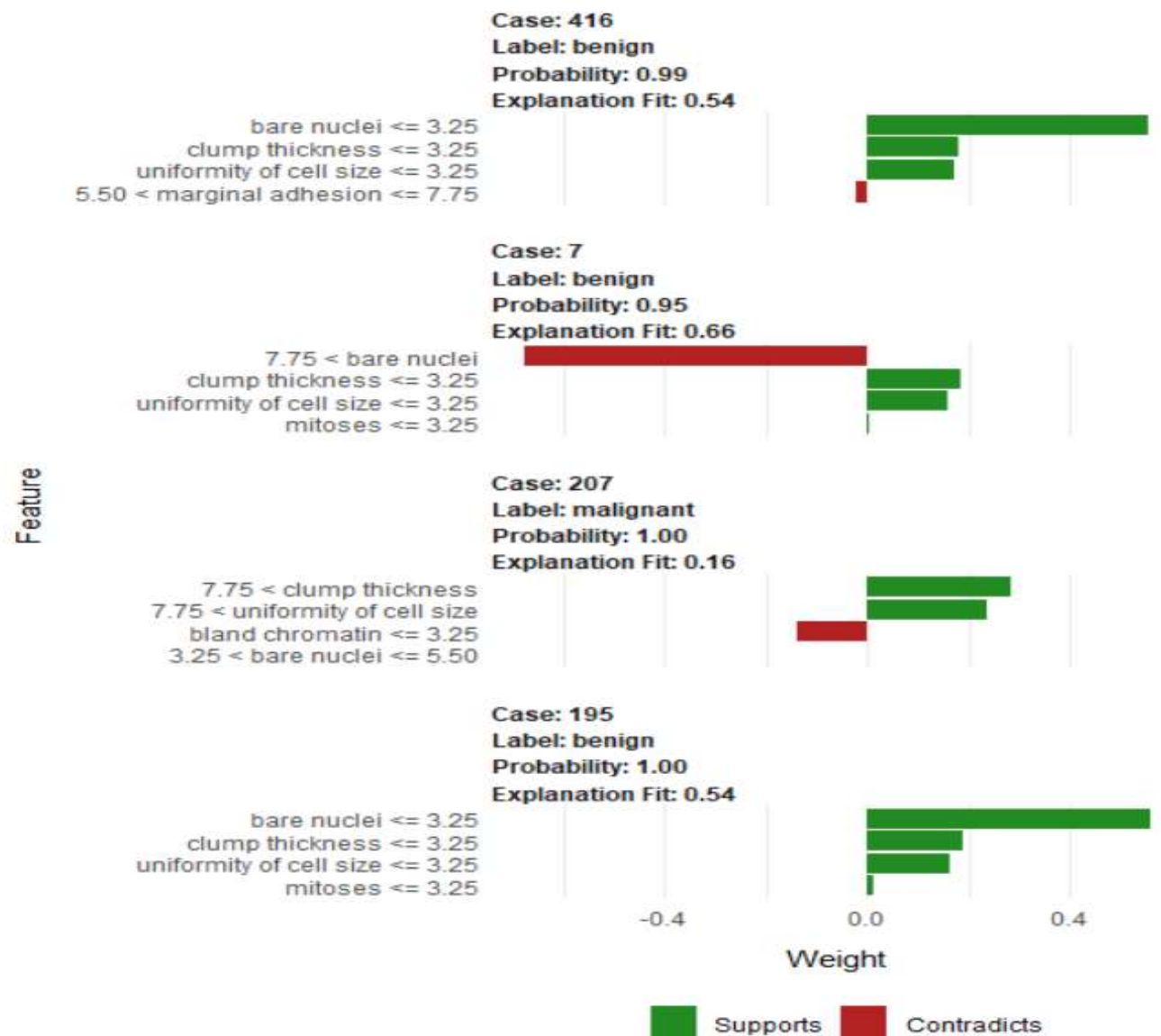
Local Interpretable Model-Agnostic Explanations

lime

Generate data points based on training data

Compute complex model predictions from the generated data to find the ‘most useful features’

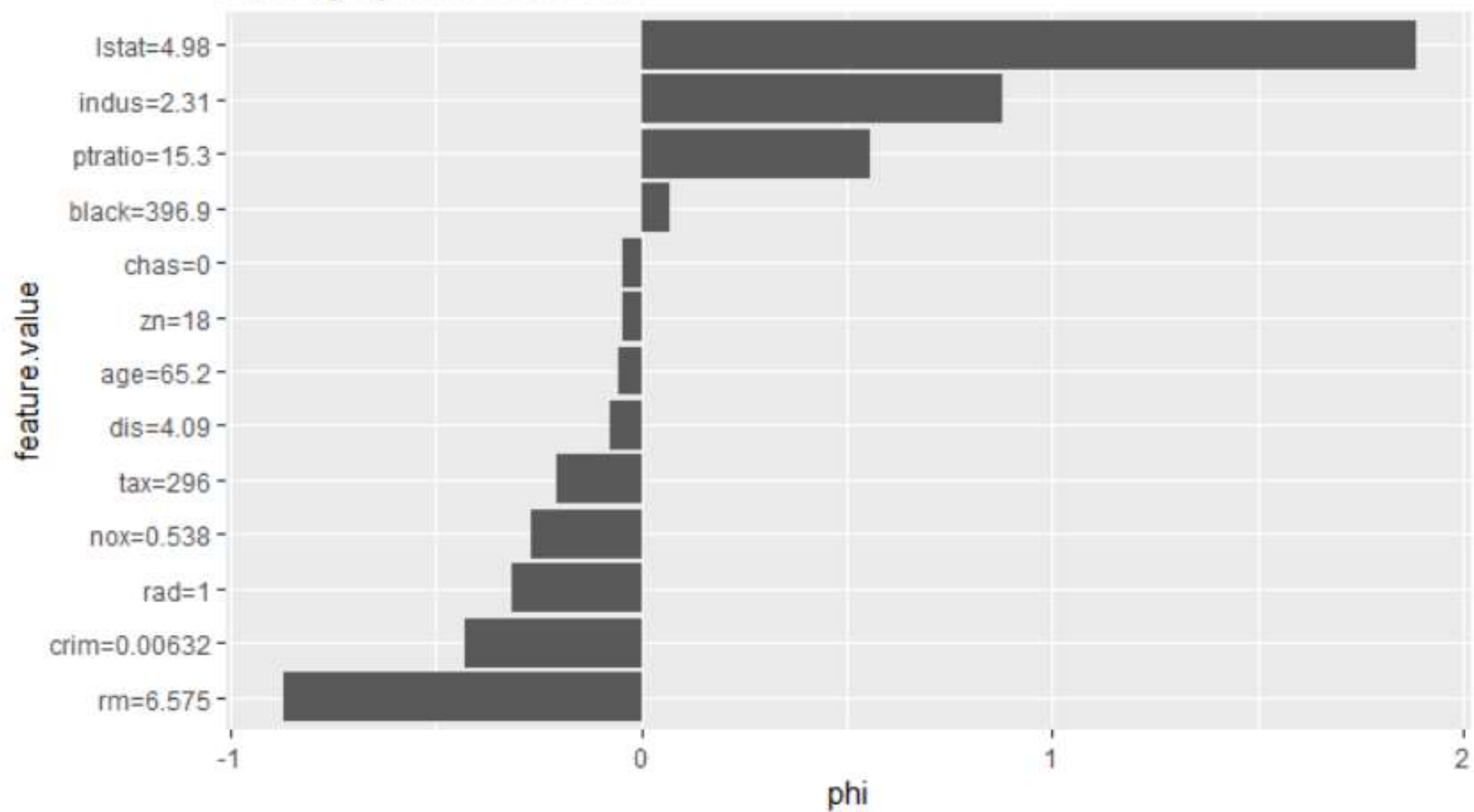
Fit a local linear model for the ‘most useful features’ and use the feature coefficients as reason codes



Shapley Values (Coalition Attribution)

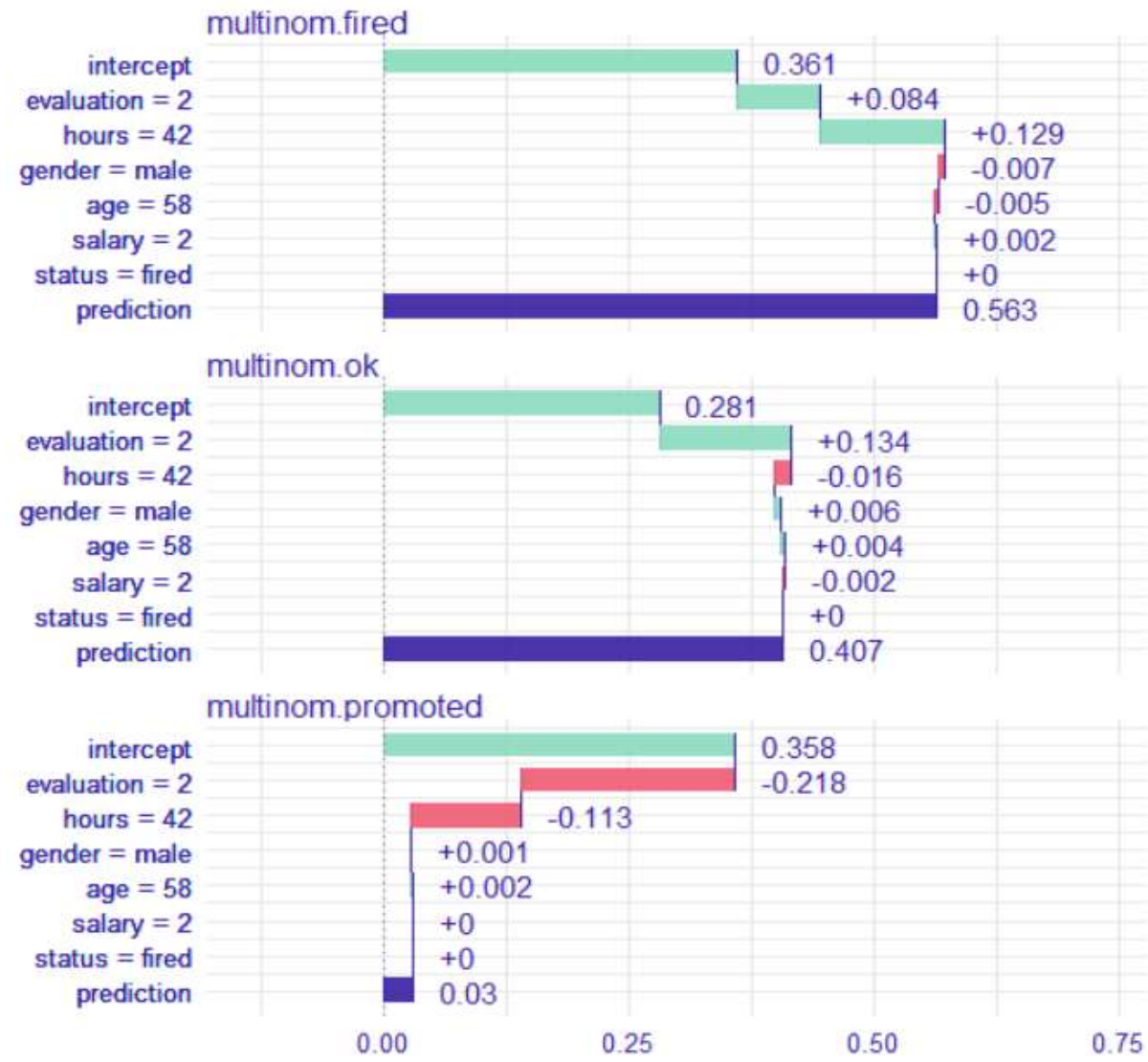
IML/ShapleyR/Shapper/
ExplainPrediction

Actual prediction: 25.75
Average prediction: 22.56



BreakDown
(Feature Contributions)

iBreakDown*



Gotchas?

Stability?

Fidelity?

Scope?

Foundations

Interpretable Data

Performant Model

“Variable Importances”...

Linearity Assumptions...

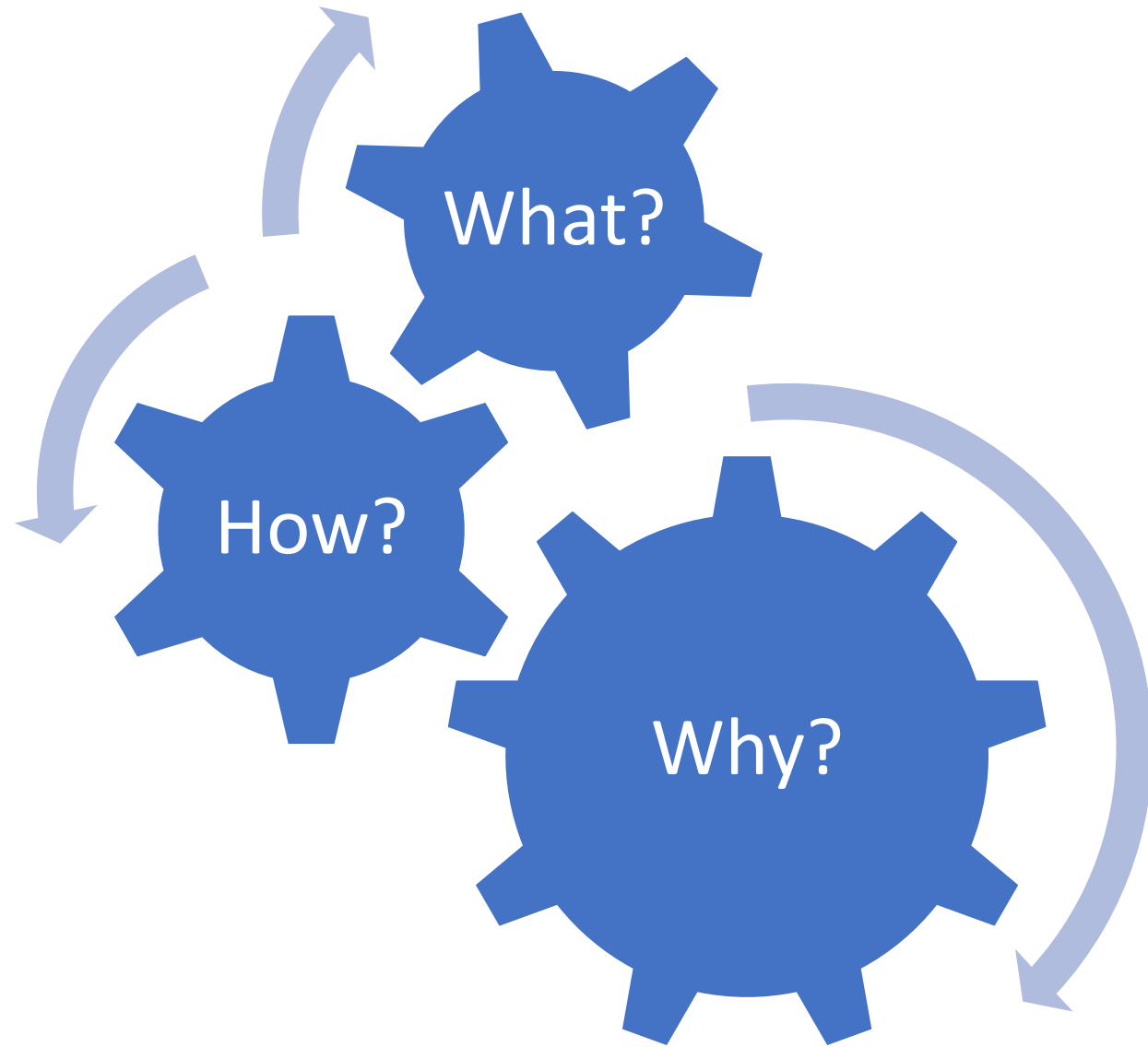
Collinearity

Computation Time

Conclusion

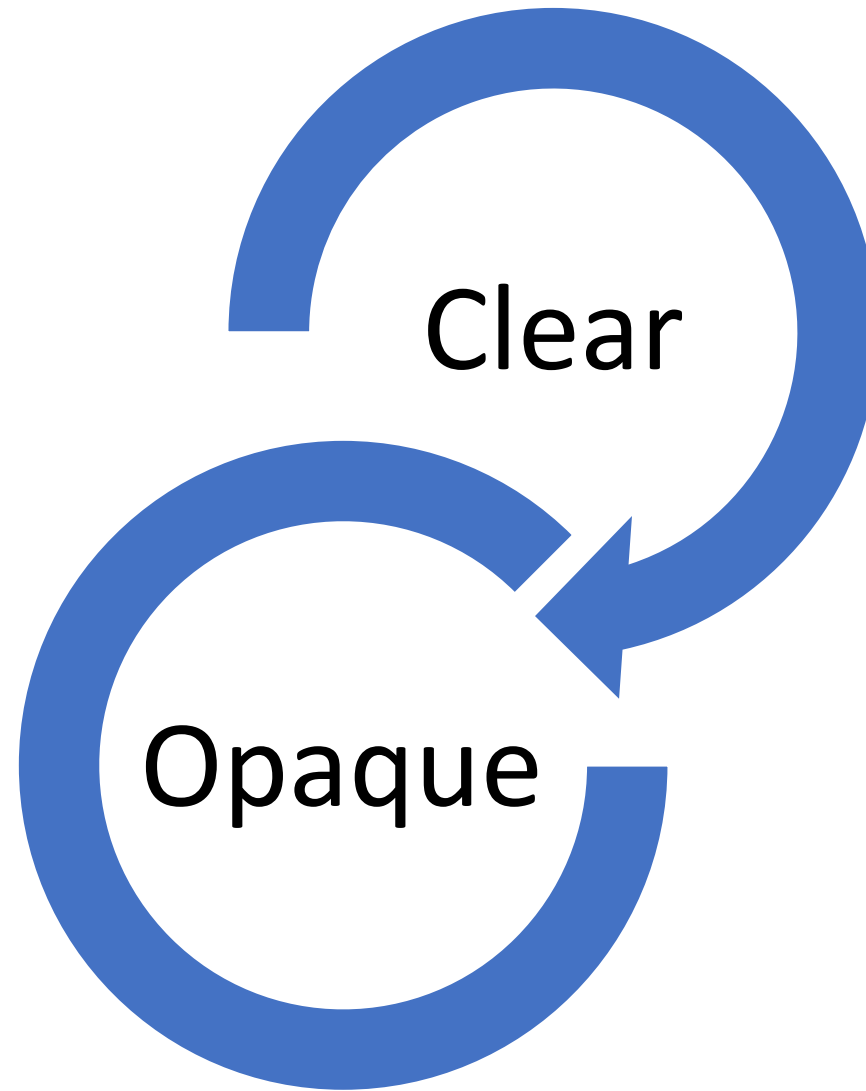
Compare different ML Models

Compare different Interpretation
Approaches



Consider Local and Global

Combine Multiple Perspectives

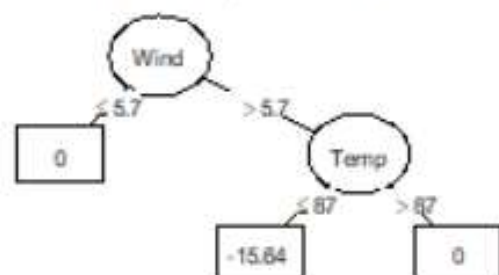


More transparent models?

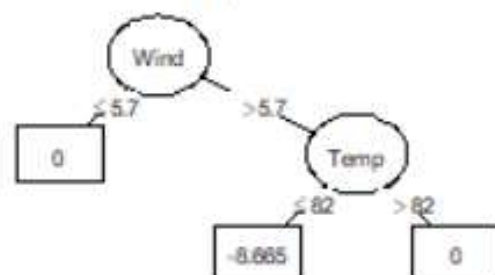
Rule Ensembles

pre

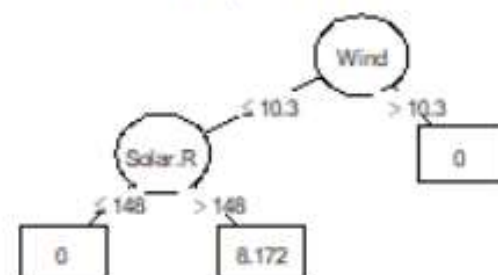
rule191: Importance = 6.263



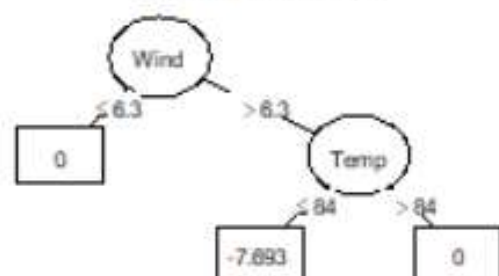
rule173: Importance = 4.074



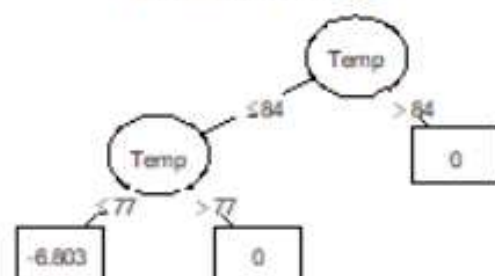
rule204: Importance = 4.056



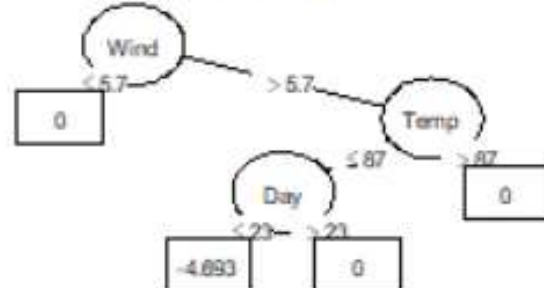
rule42: Importance = 3.591



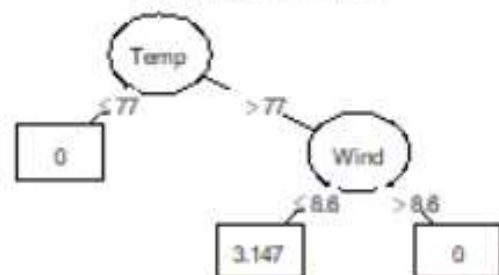
rule10: Importance = 3.4



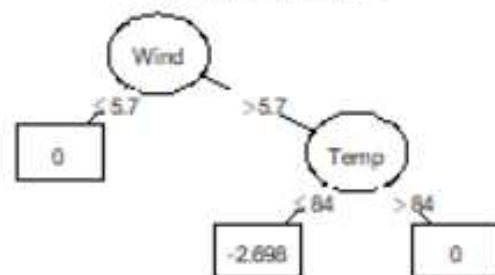
rule192: Importance = 2.275



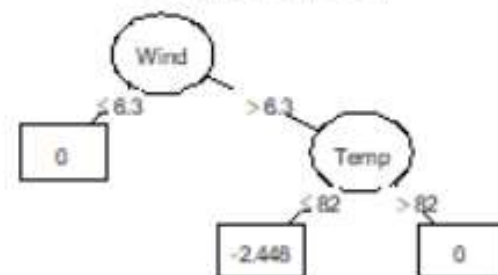
rule93: Importance = 1.457



rule51: Importance = 1.228



rule25: Importance = 1.167



Recommendations:

Data - [Data Explorer]

Local - Shapley Values [IML]

Global - Variable Importance [DALEX]

Plots - ALE, Ceteris Paribus [DALEX]

Plots – ALE, PDP, ICE [IML]

Thanks for Listening!

dean_allsopp@hotmail.com