

# Introduction to Interpretable Machine Learning

Dean Allsopp

Global AI Bootcamp

December 15, 2018

# Machine Learning?

Iteratively derive structure from data

(Patterns, Rules, Clusters, Policies)

ML

```
graph TD; ML[ML] --- Supervised[Supervised]; ML --- Unsupervised[Unsupervised]; ML --- Reinforcement[Reinforcement];
```

A hierarchical diagram with 'ML' at the top level. A vertical line descends from 'ML' and connects to a horizontal line. From this horizontal line, three vertical lines descend to three separate boxes below: 'Supervised', 'Unsupervised', and 'Reinforcement'.

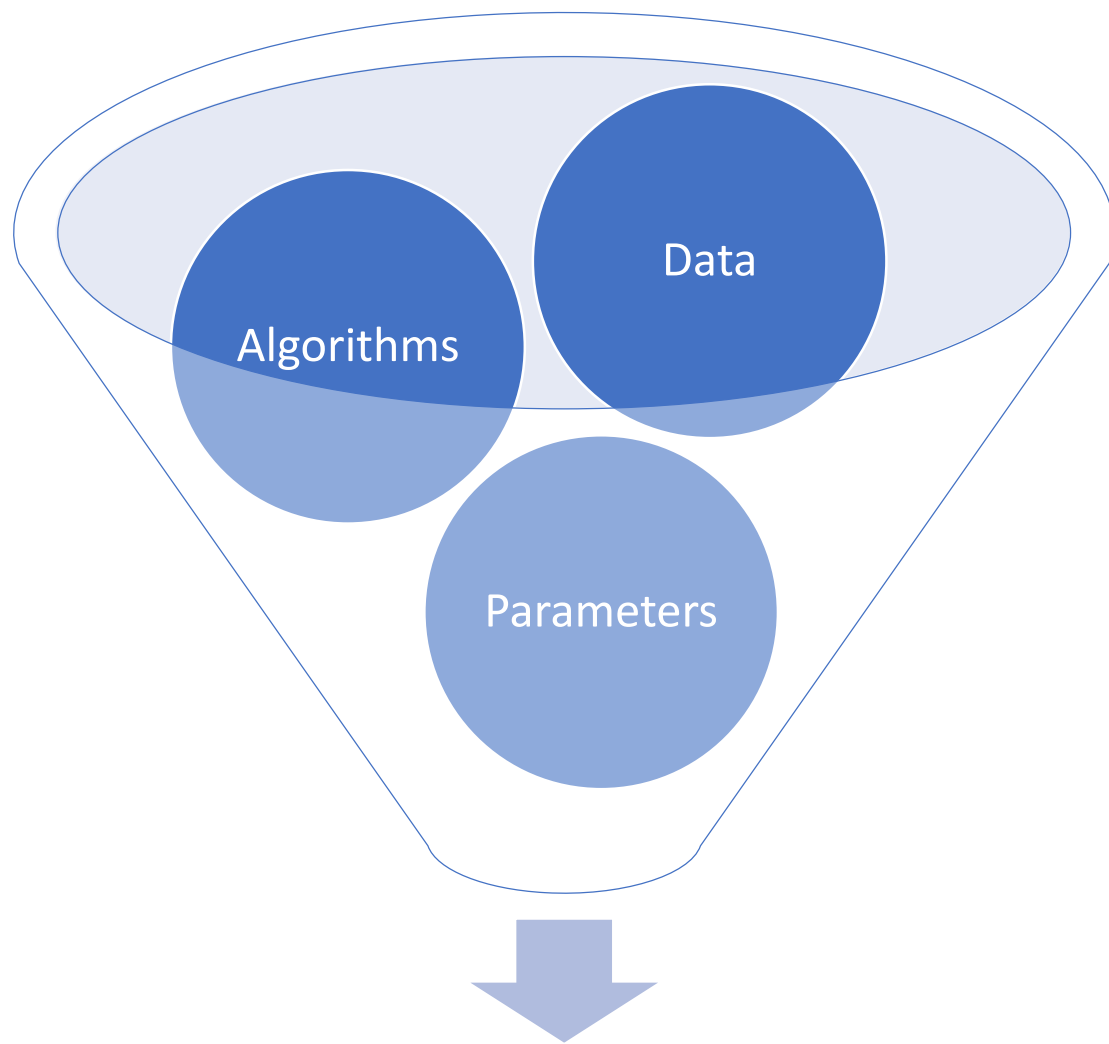
Supervised

Unsupervised

Reinforcement

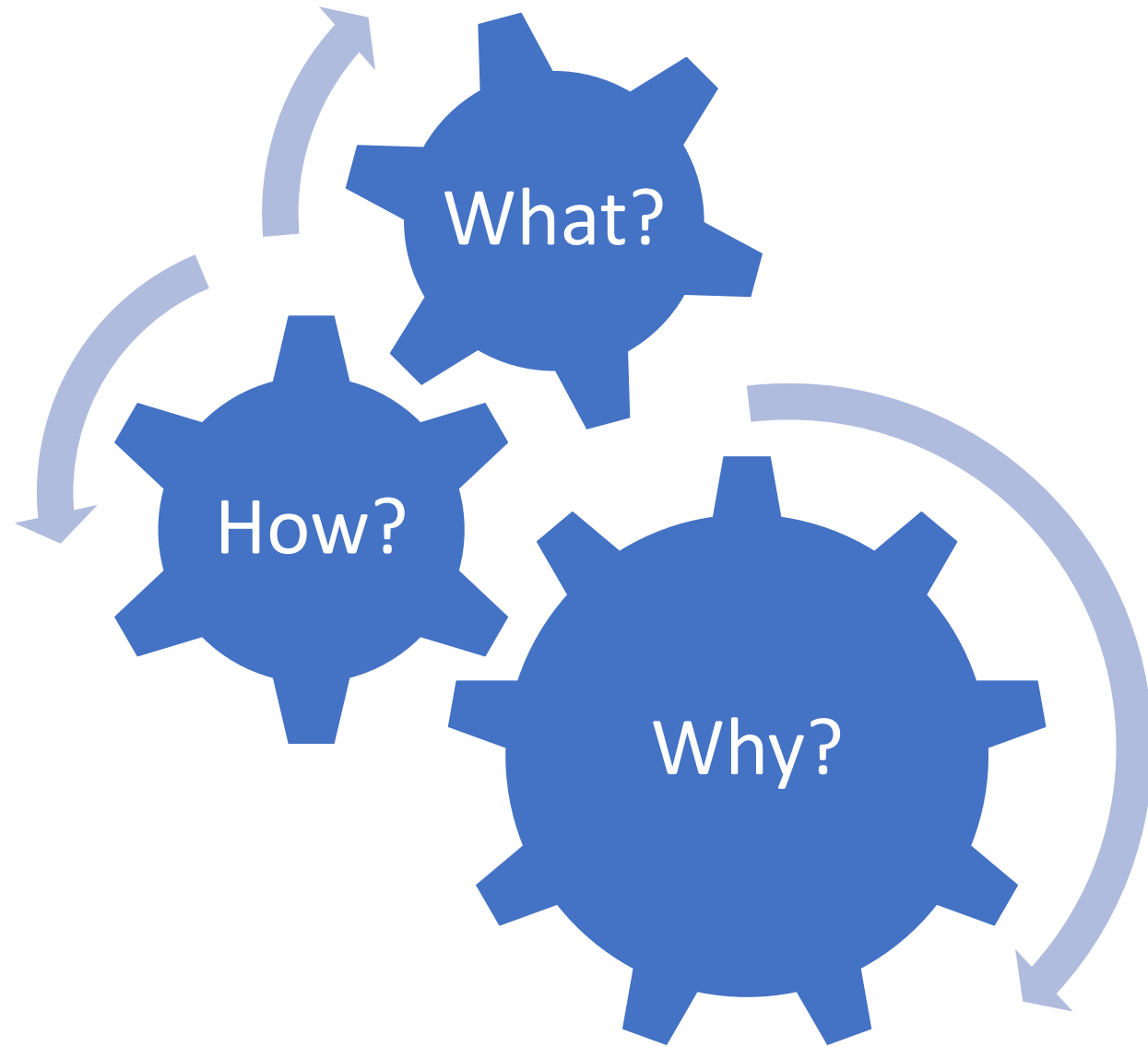
# Supervised Machine Learning?

(We have labels)



Predictions

# Interpretable Machine Learning?





What? – Data (Features)

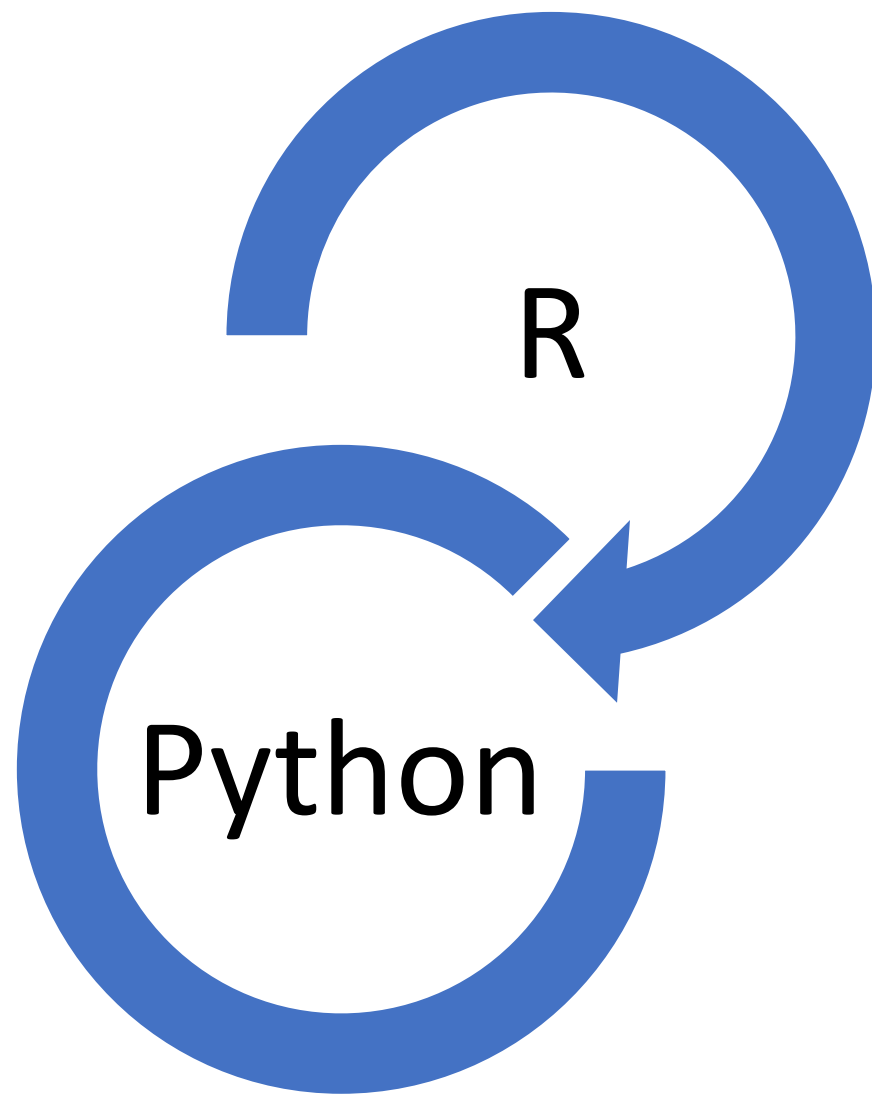
How? – Model (Global)

Why? – Prediction (Local)

How well do you understand  
the data?



Tools?



EDA

Data Explorer (R)

Pandas Profiling (Py)

```
create_report([dataset])
```

# Data Profiling Report

- [Basic Statistics](#)
  - [Raw Counts](#)
  - [Percentages](#)
- [Data Structure](#)
- [Missing Data Profile](#)
- [Univariate Distribution](#)
  - [Histogram](#)
  - [Bar Chart \(by frequency\)](#)
  - [QQ Plot](#)
- [Correlation Analysis](#)
- [Principle Component Analysis](#)

## Basic Statistics

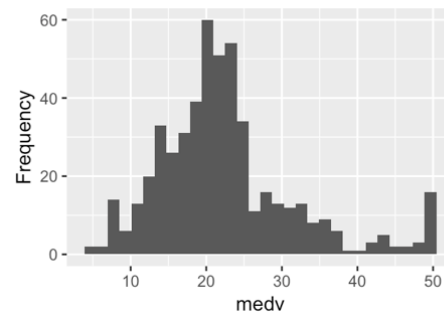
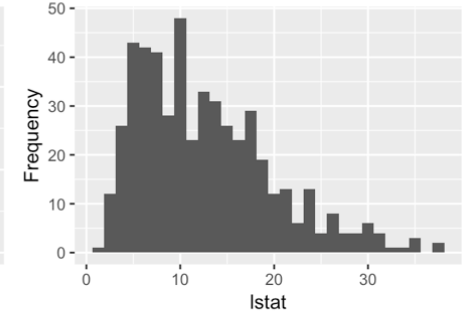
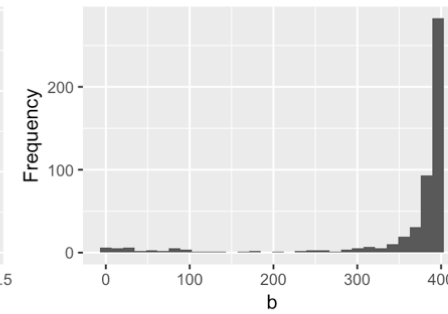
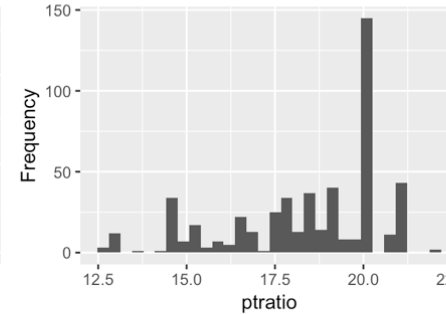
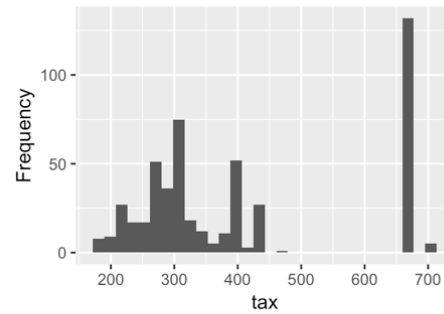
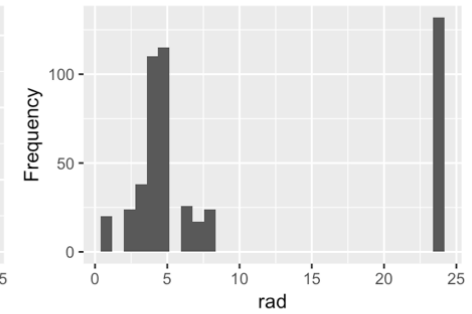
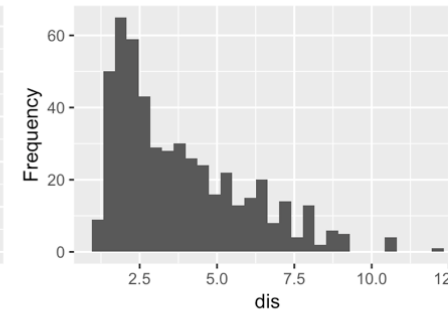
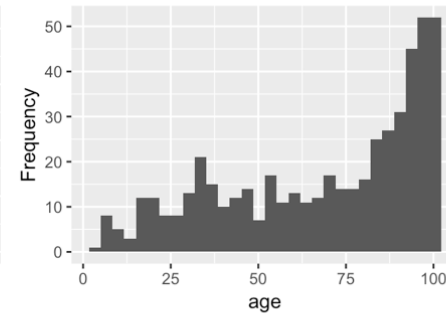
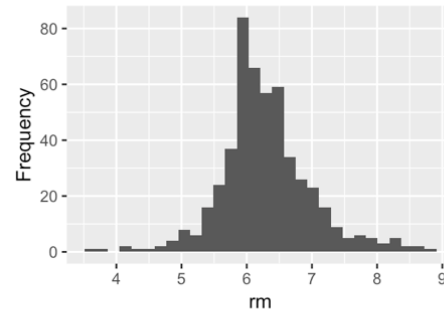
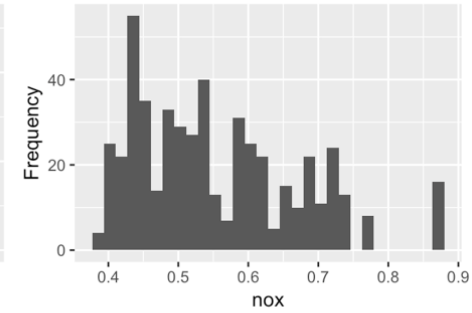
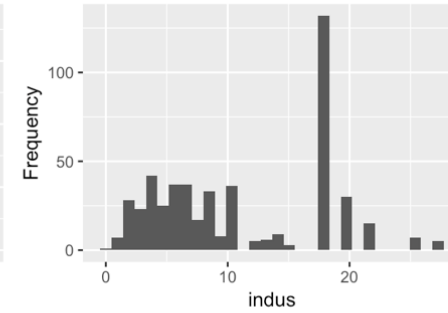
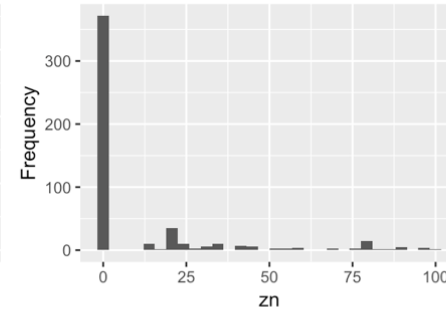
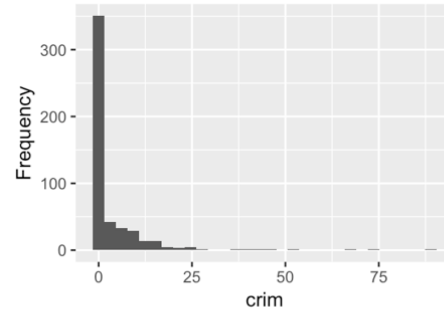
### Raw Counts

Name	Value
Rows	506
Columns	14
Discrete columns	1
Continuous columns	13
All missing columns	0
Missing observations	0
Complete Rows	506
Total observations	7,084
Memory allocation	57.4 Kb

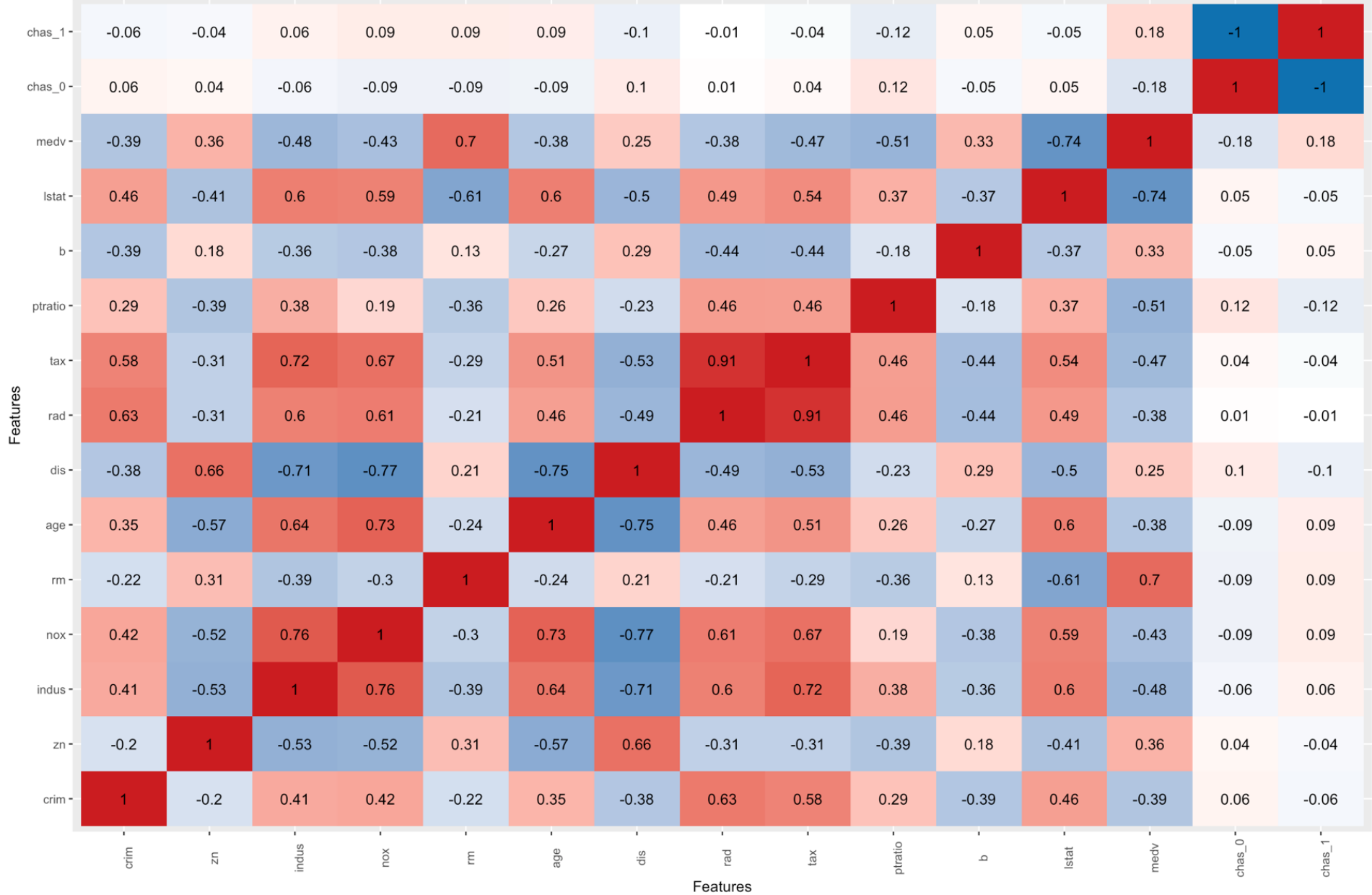


# Univariate Distribution

## Histogram



# Correlation Analysis

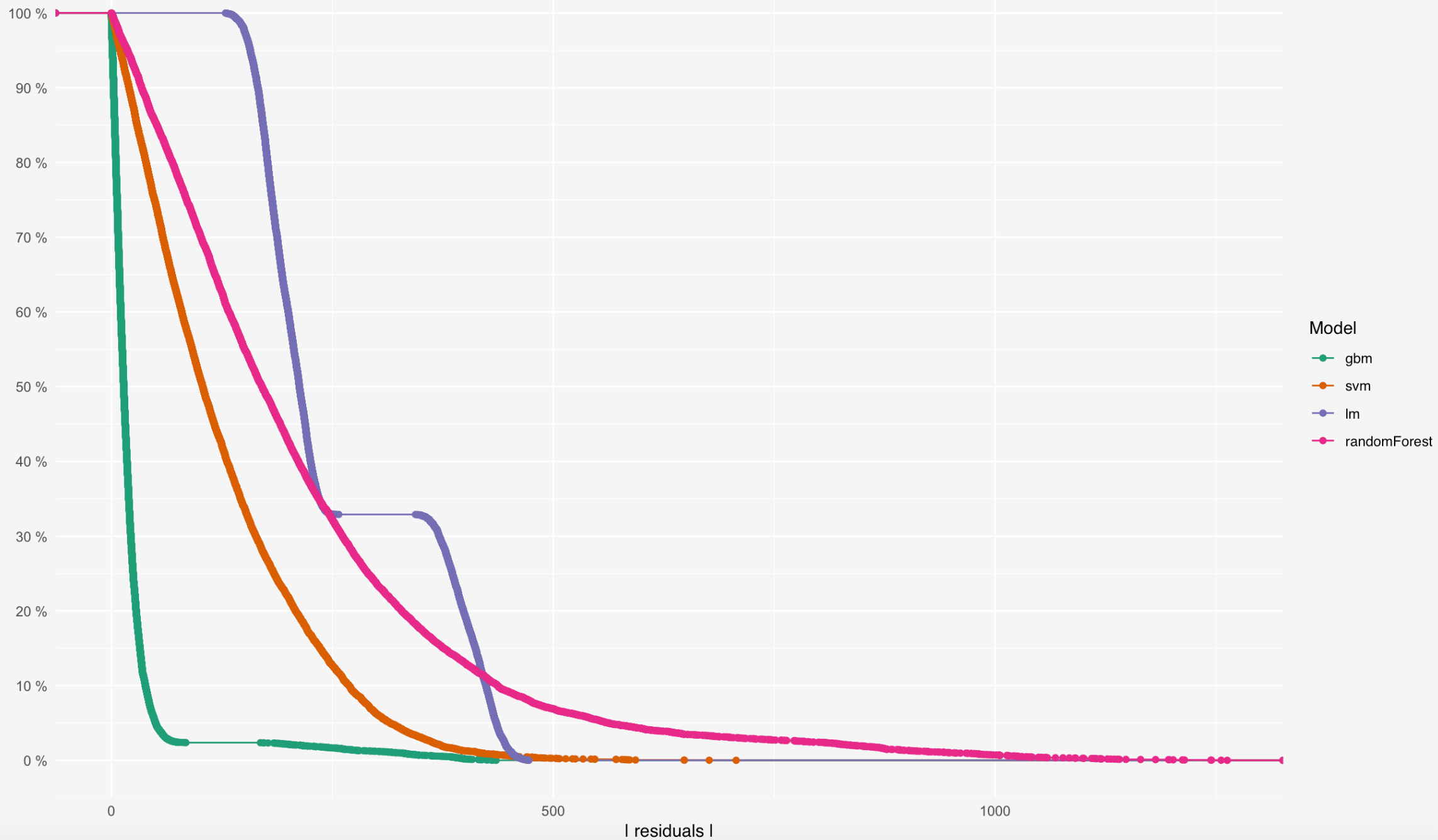


# Model Residuals Plot

DALEX (R)

YELLOWBRICK (Py)

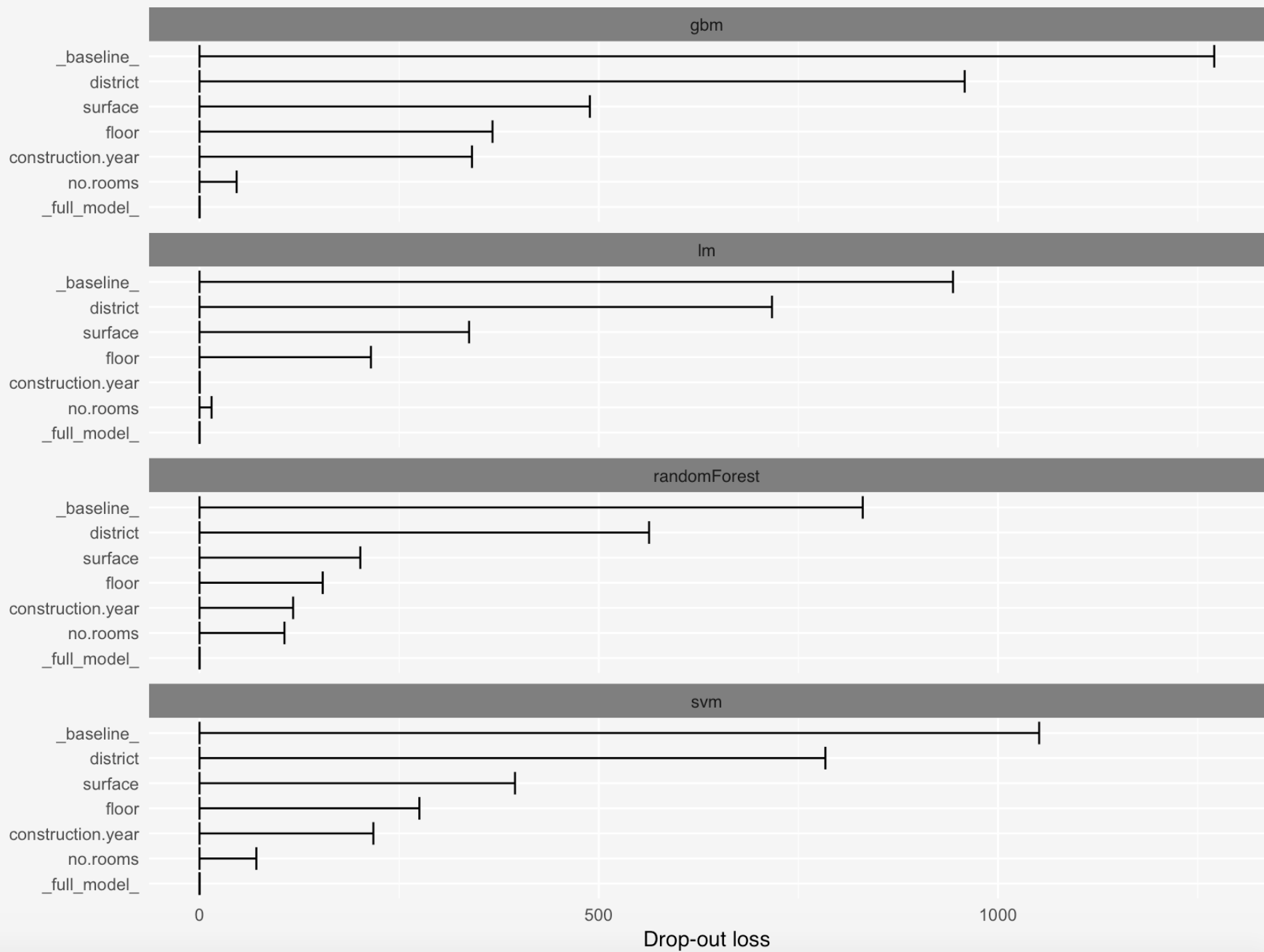
Distribution of  $|$  residuals  $|$



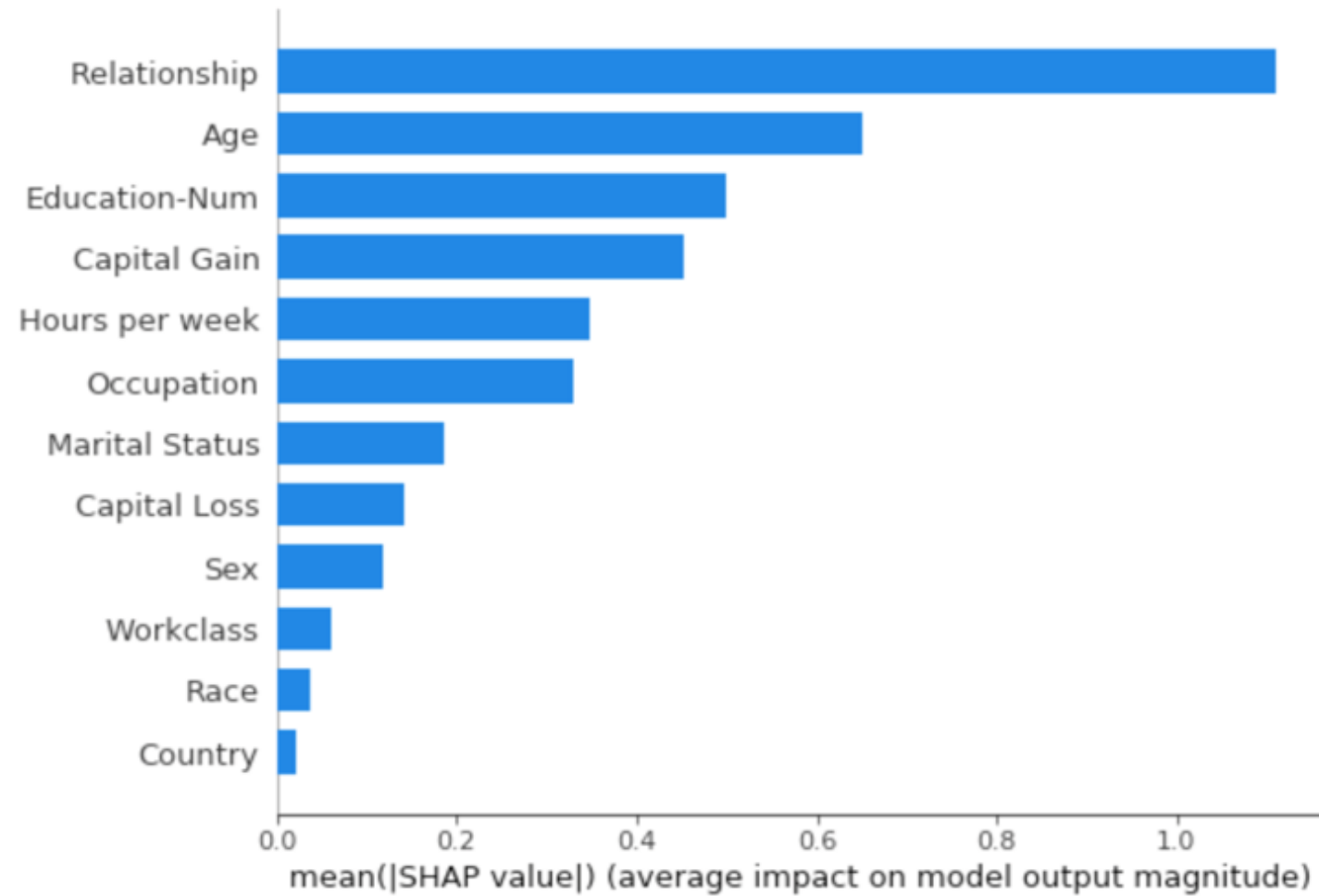
# Variable Importance Plot

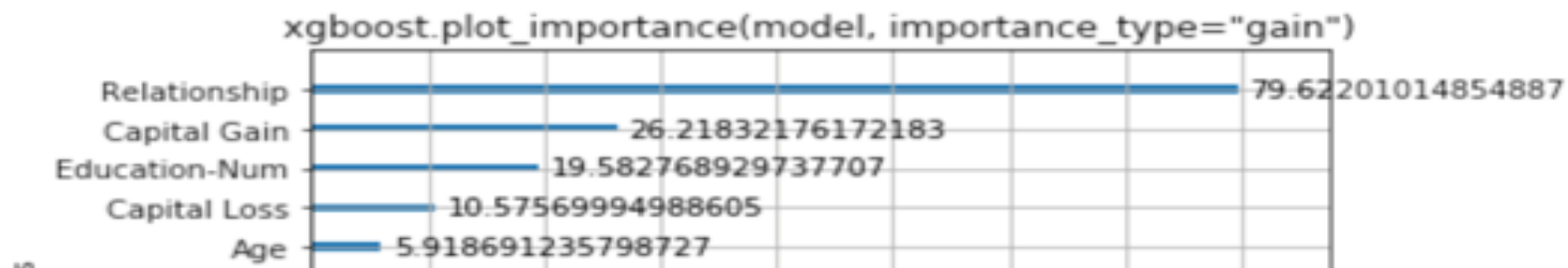
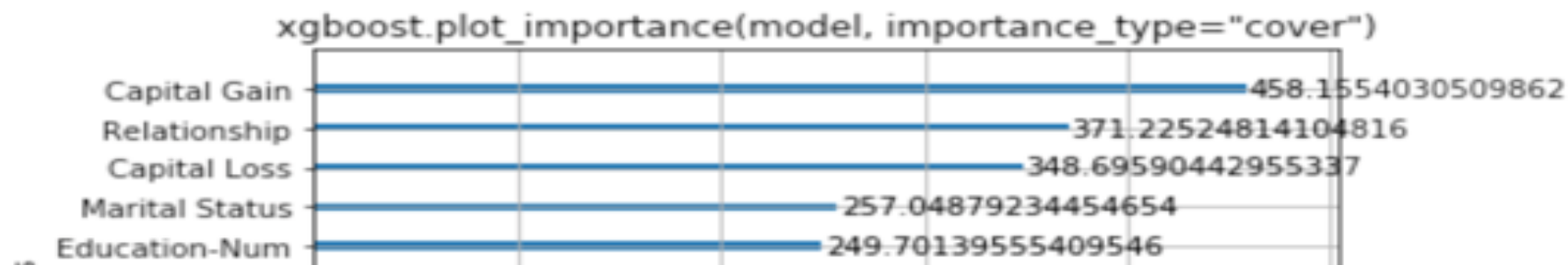
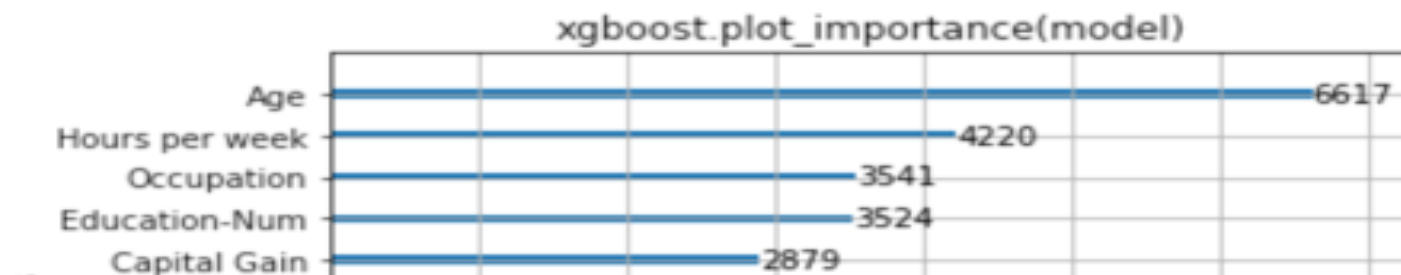
DALEX (R)

SHAP (PY)



```
In [11]: shap.summary_plot(shap_values, X_display, plot_type="bar")
```





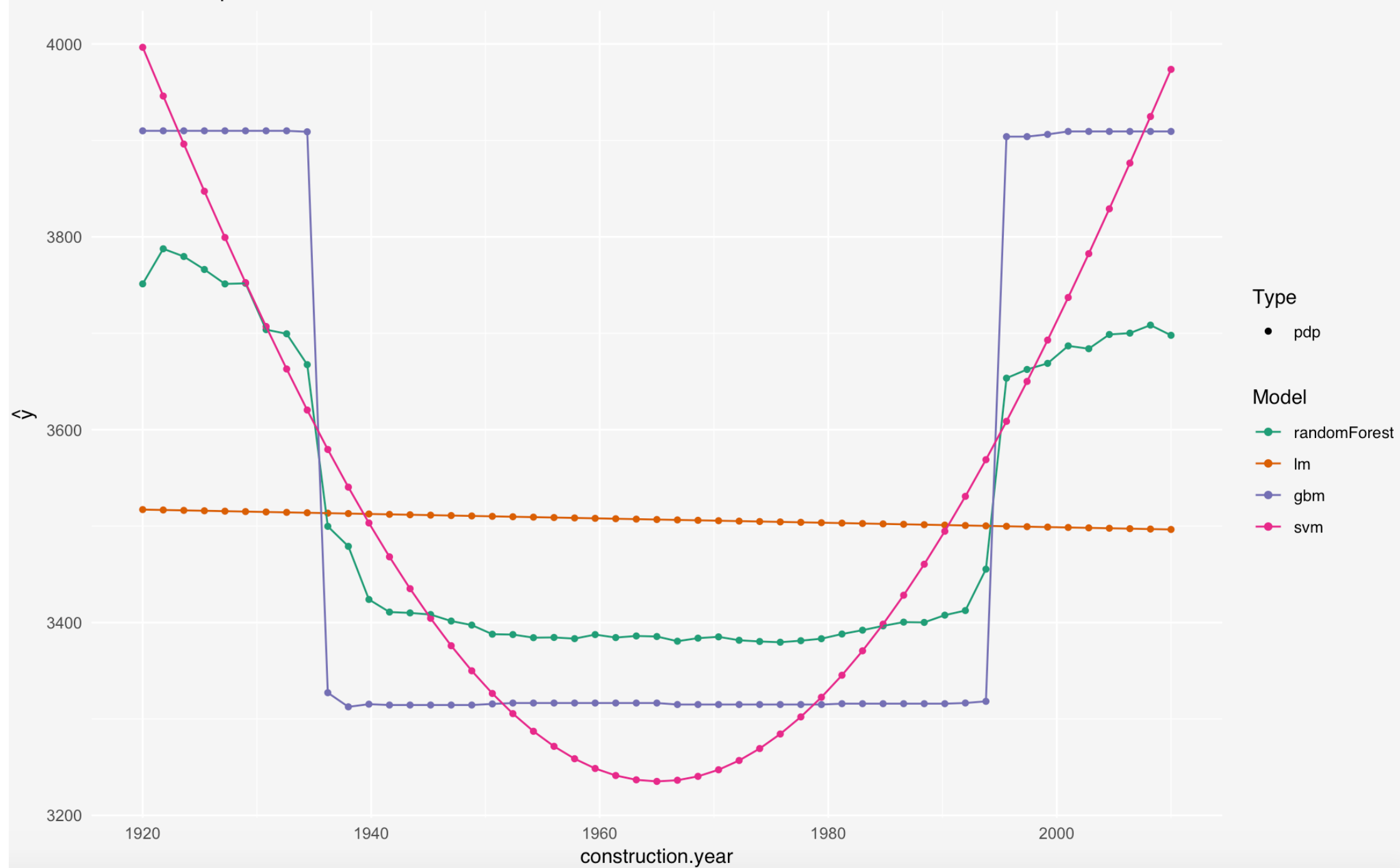


# Partial Dependence & Individual Conditional Plots

DALEX/IML (R)

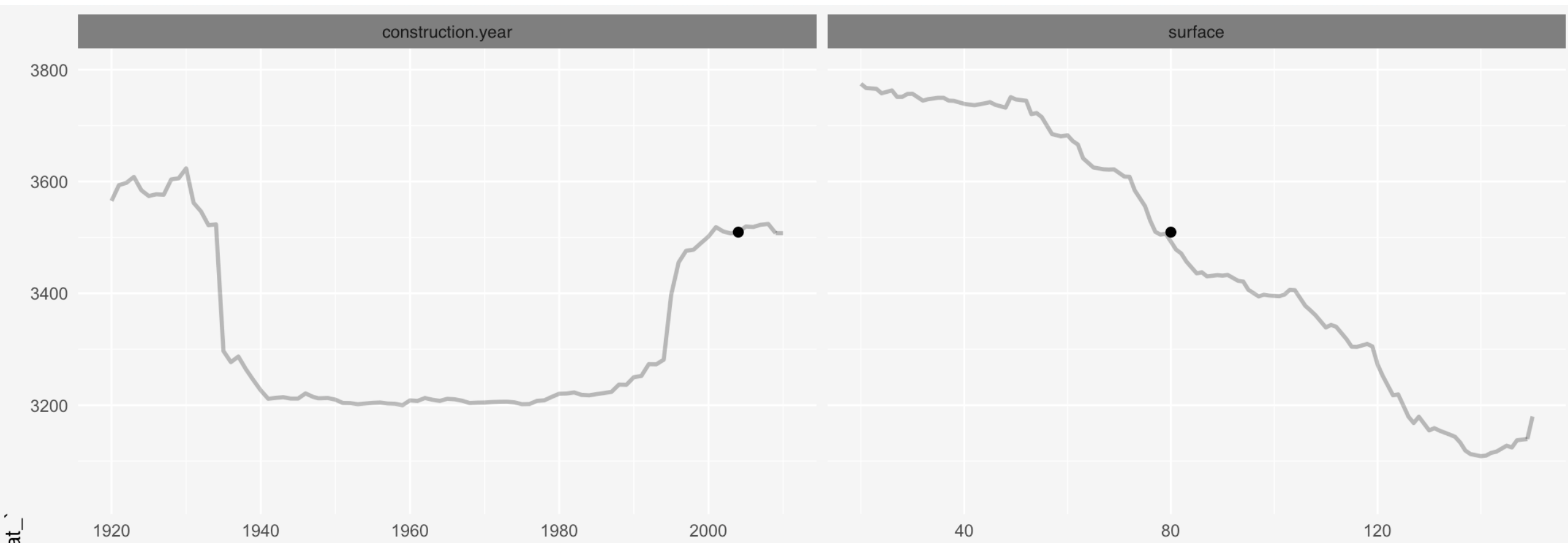
PDPBOX (Py)

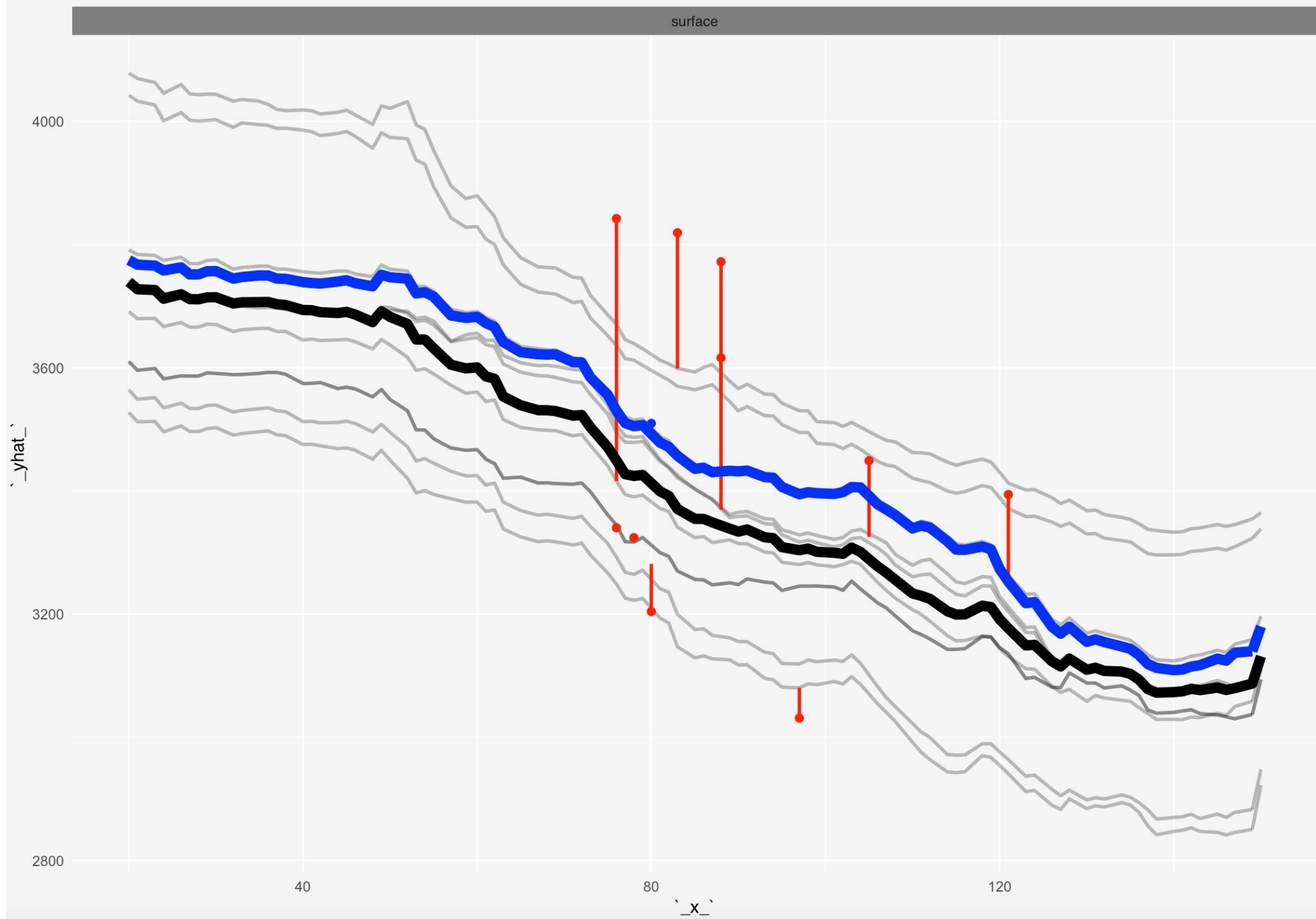
Variable response

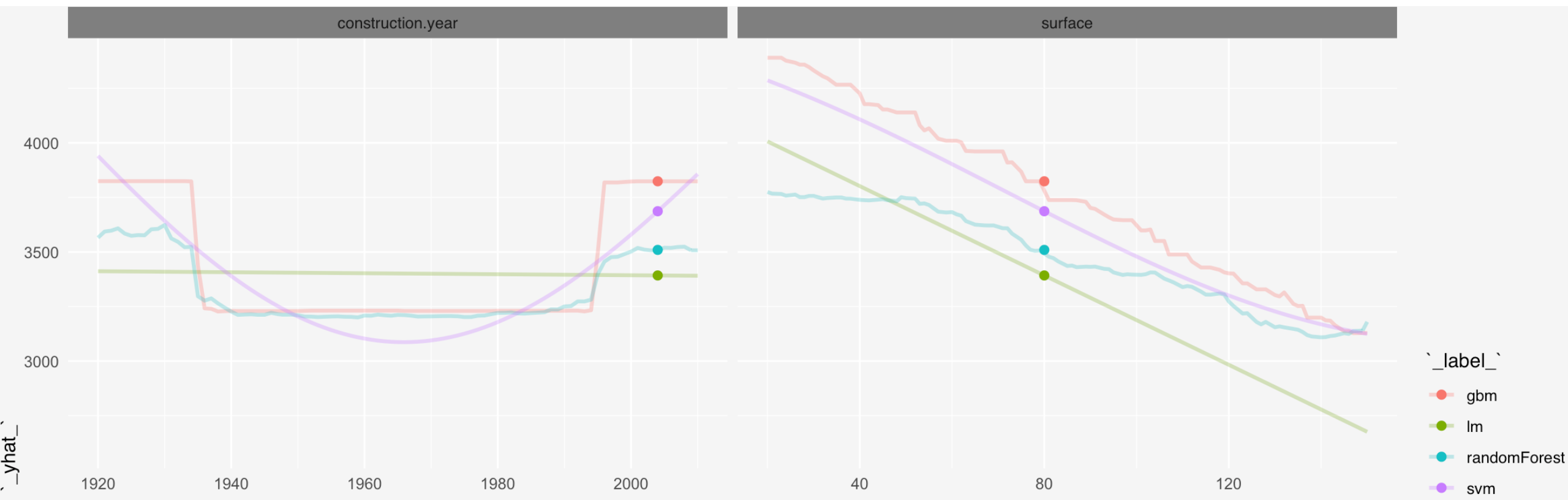


# Ceteris Paribus Plot

DALEX (R)







LIME (R)

LIME (Py)

```

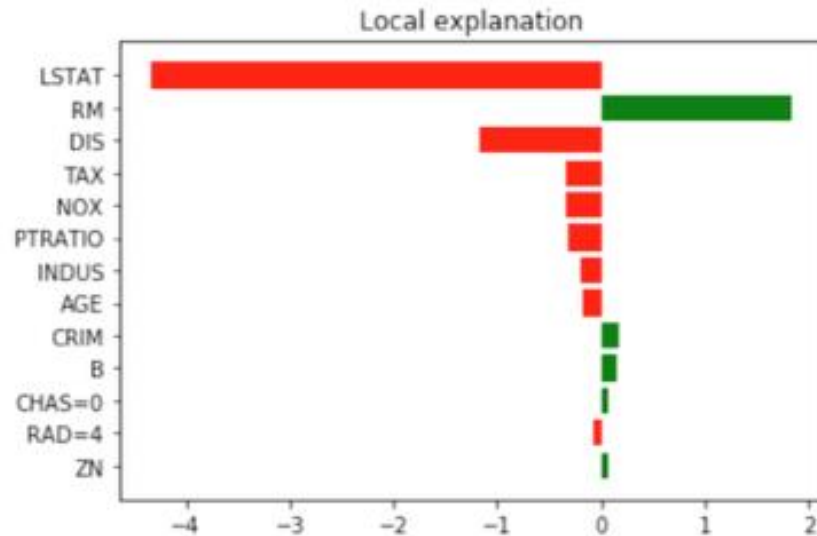
In [2]: # import lime tools
import lime
import lime.lime_tabular

# generate an "explainer" object
categorical_features = np.argwhere(np.array([len(set(boston.data[:,x])) for x in range(boston.data.shape[1])]) <= 10).
explainer = lime.lime_tabular.LimeTabularExplainer(train, feature_names=boston.feature_names, class_names=['price'], ca

In [3]: #generate an explanation
i = 13
exp = explainer.explain_instance(test[i], rf.predict, num_features=14)

In [4]: %matplotlib inline
fig = exp.as_pyplot_figure();

```





*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*

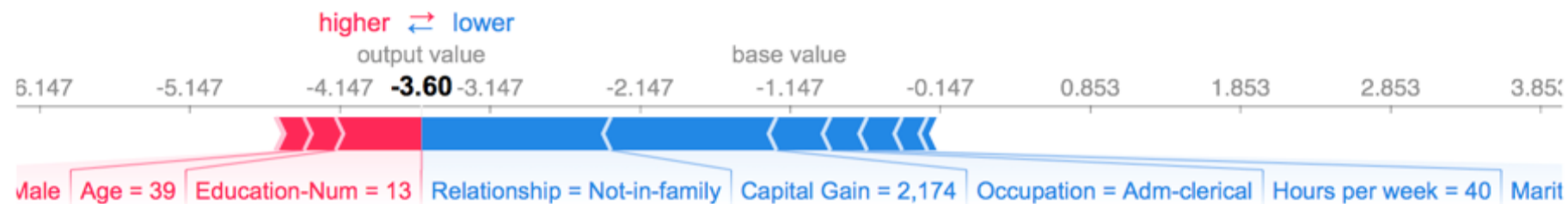
# Shapley Values (Coalition Attribution)

IML (R)

SHAP (Py)

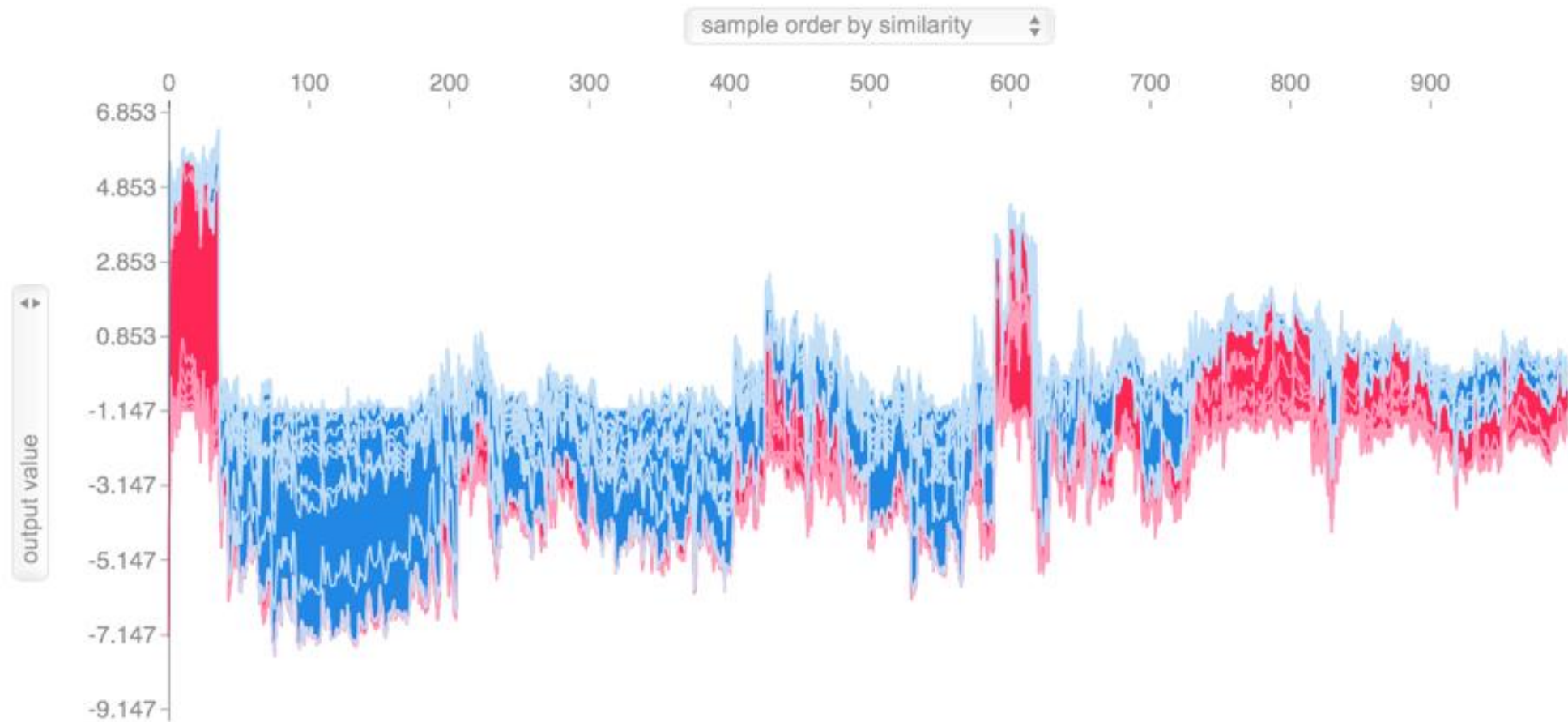
```
In [9]: shap.force_plot(explainer.expected_value, shap_values[0,:], X_display.iloc[0,:])
```

Out[9]:



```
In [10]: shap.force_plot(explainer.expected_value, shap_values[:1000,:], X_display.iloc[:1000,:])
```

Out[10]:



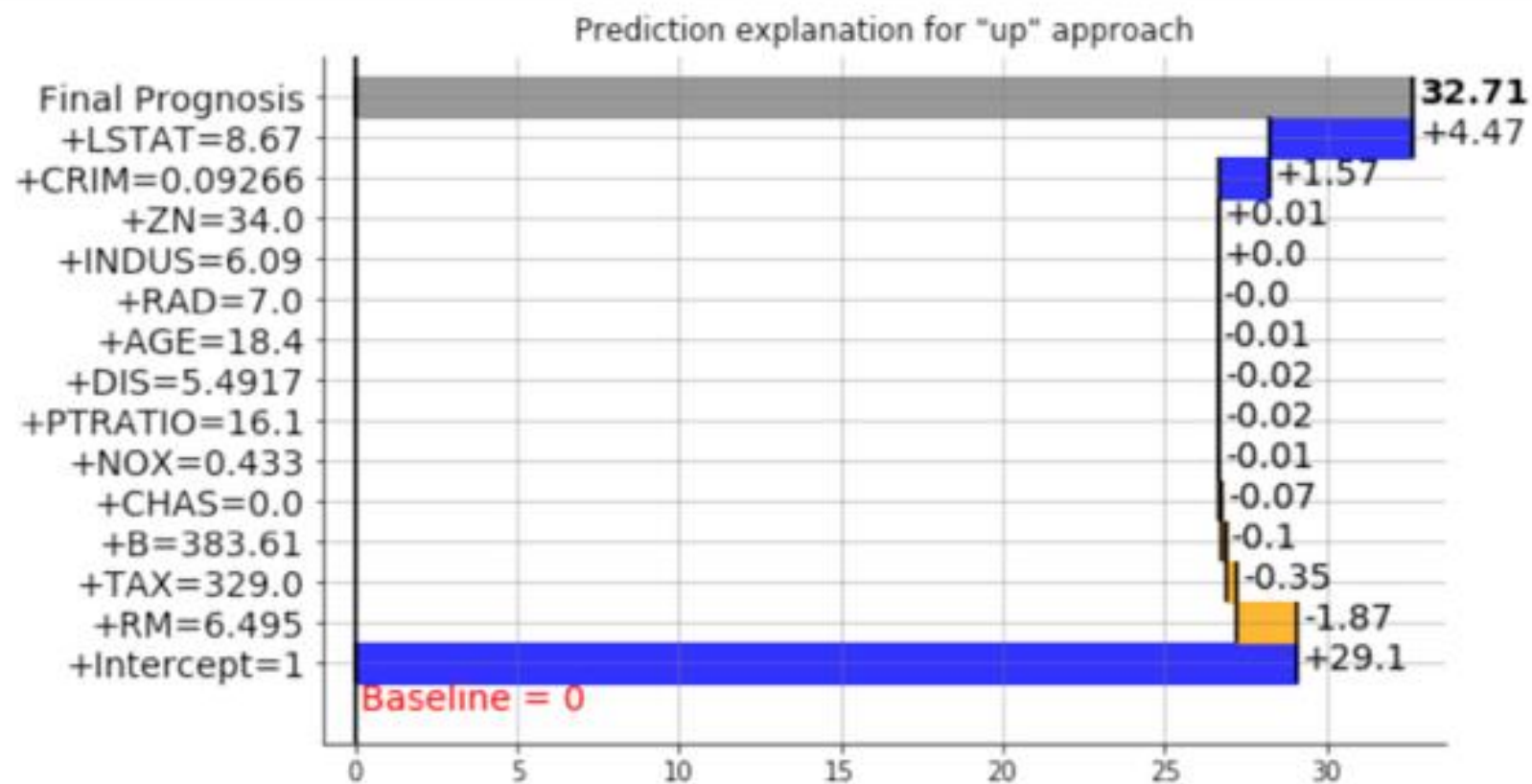
BreakDown

DALEX (R)

PyBreakdown (Py)

```
In [7]: # visualisation
```

```
explanation.visualize(figsize=(8,5),dpi=100)
```



# Surrogate Trees

IML (R)

SKATER (Py)

```
In [49]: surrogate_explainer_reg.decisions_as_txt('local', X_test.iloc[sample_index])
```

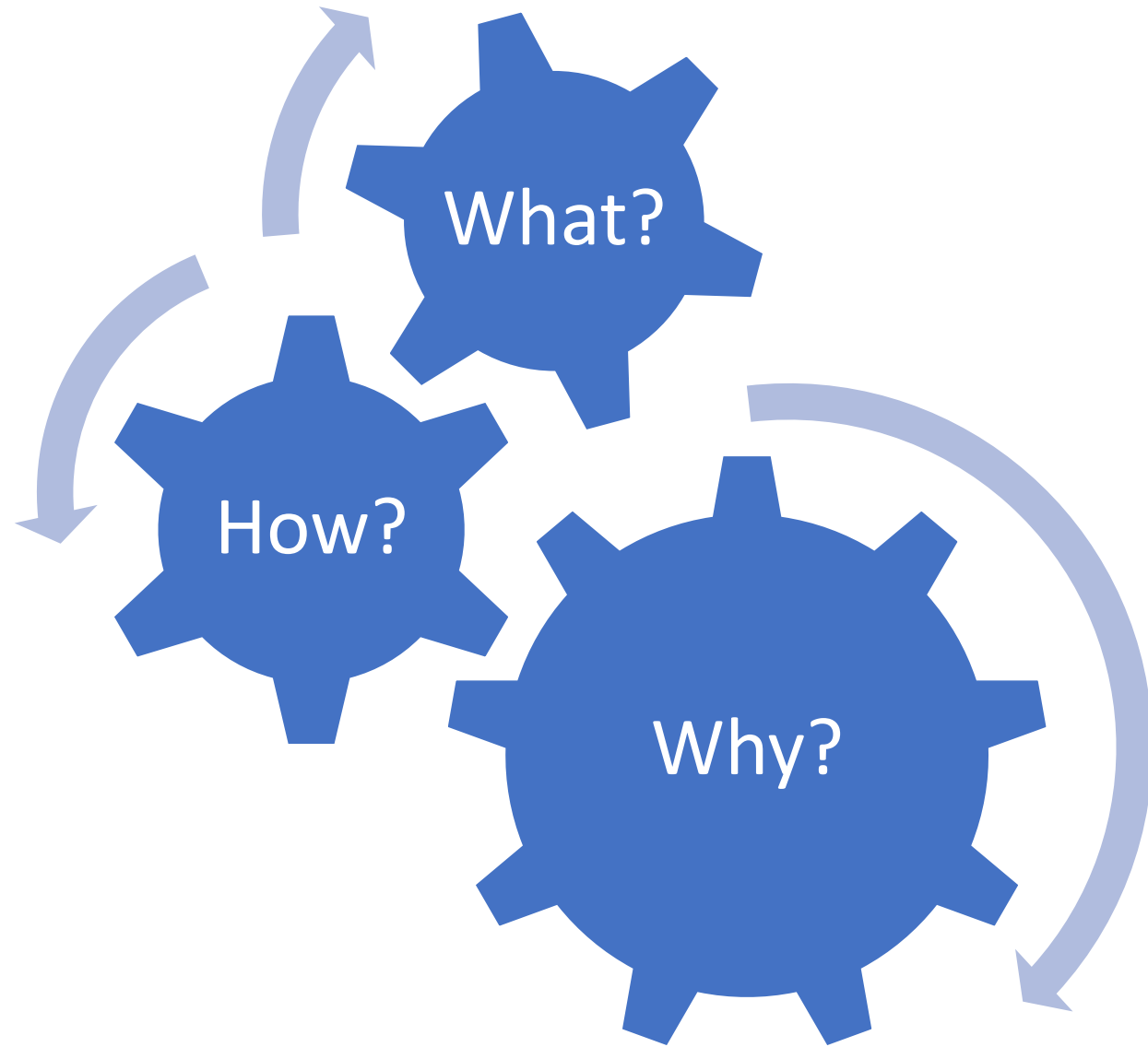
```
As RM[5.914] <= 6.837500095367432 then,  
  As LSTAT[18.33] > 14.399999618530273 then,  
    As CRIM[0.31827] <= 7.084139823913574 then,  
      As CRIM[0.31827] <= 0.6147900223731995 then,  
        As DIS[3.9986] > 1.9799000024795532 then,  
          As AGE[83.2] > 73.30000305175781 then,  
            As CRIM[0.31827] > 0.17127001285552979 then,  
              As CHAS[0.0] <= 0.5 then,  
                As B[390.7] <= 391.875 then,  
                  As NOX[0.544] > 0.48649999499320984 then,  
                    As B[390.7] > 389.96002197265625 then,  
                      Value: [[19.36]]
```



Conclusion

Multiple ML Models

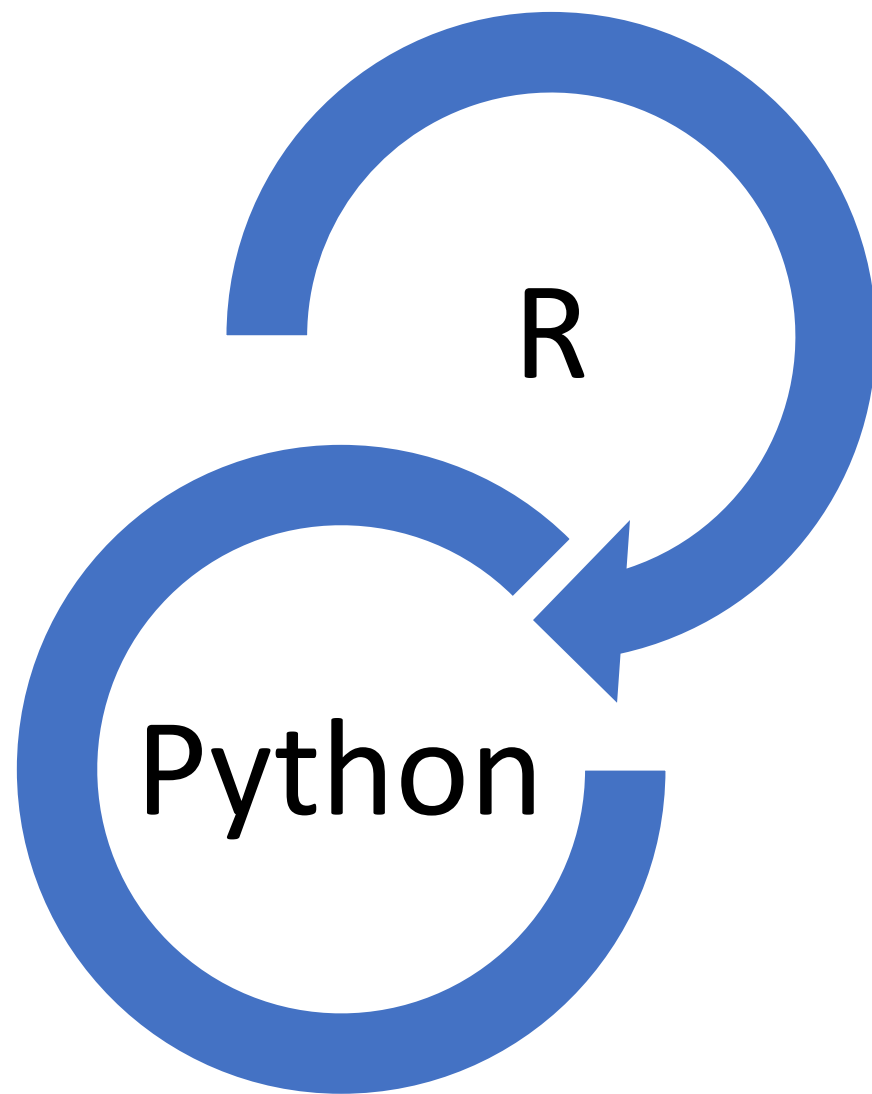
Multiple Interpretation Approaches



Data Understanding

Consider Local and Global

Combine Multiple Perspectives



# Recommendations (Python):

Local - [SHAP]

Global - Surrogate Trees [SKATER]

Data -[Pandas Profiling]

Visual - [Yellowbrick], [PDPBox]

## Recommendations (R):

Local - SHAP [IML]

Global - Surrogate Trees [IML], Variable Importance [DALEX]

Data - [Data Explorer]

Visual - Breakdown, ALE, Ceteris Paribus [DALEX]

Visual - PDP, ICE [IML]

# Thanks for Listening!

[dean\\_allsopp@hotmail.com](mailto:dean_allsopp@hotmail.com)