

Seeing the Wood for the Trees: Interrogating the Structure of Random Forests

Chris Harbron

Discovery Statistics

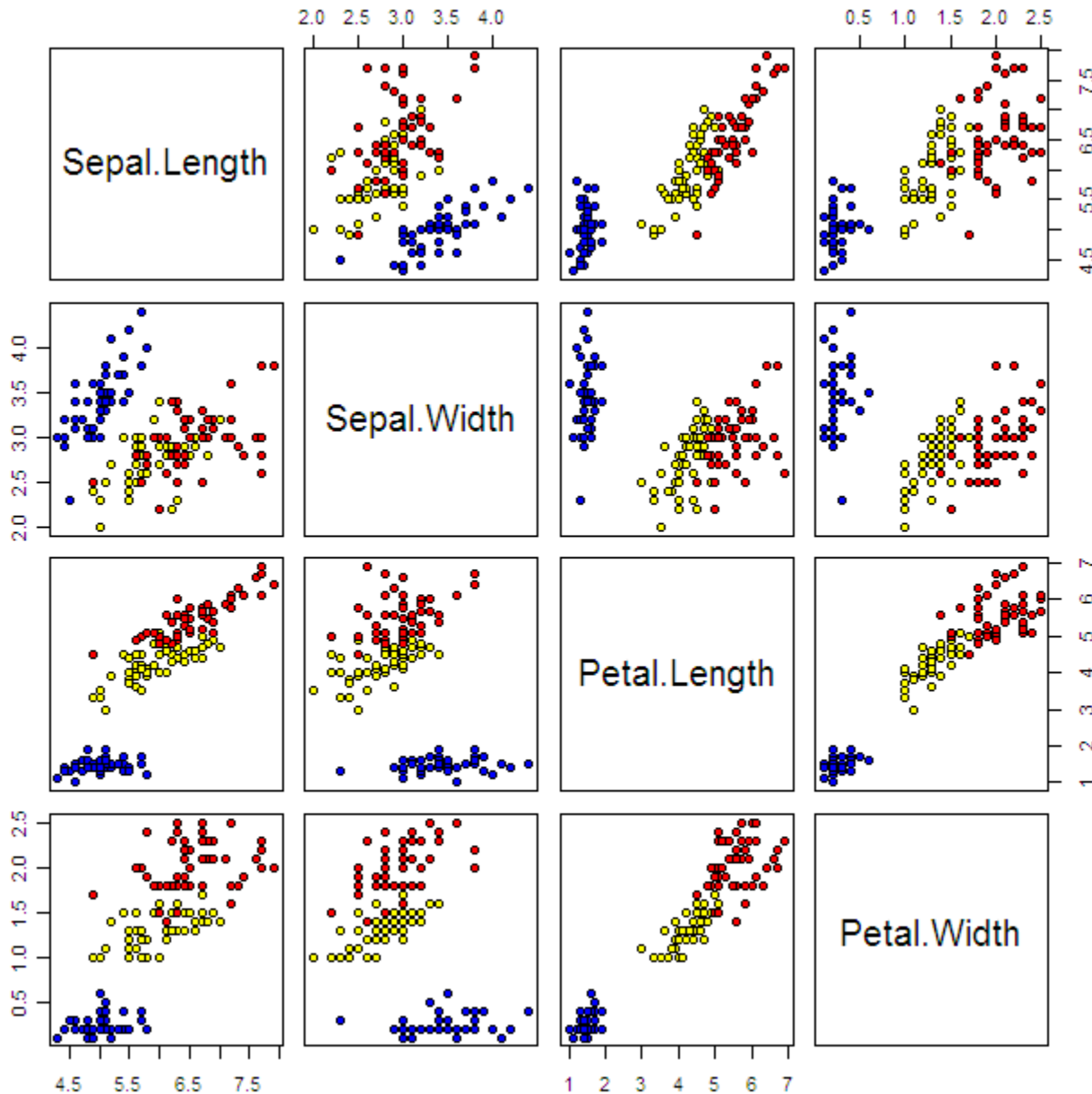
AstraZeneca



Random Forests

- First published by Breimann and Cutler in 2001
- Popular modelling tool for classification, regression and survival analysis with (highly) multivariate data
- Powerful, robust, easy to apply with many attractive properties
- Breimann was passionate that modelling was about more than just predicting:
 - “With scientific data sets more is required than an accurate prediction”
 - “Looking inside the black box is necessary”
 - Variable Importance Measures, Proximities
- SURF : An extension of Variable Importance Measures

Fisher's Iris Data



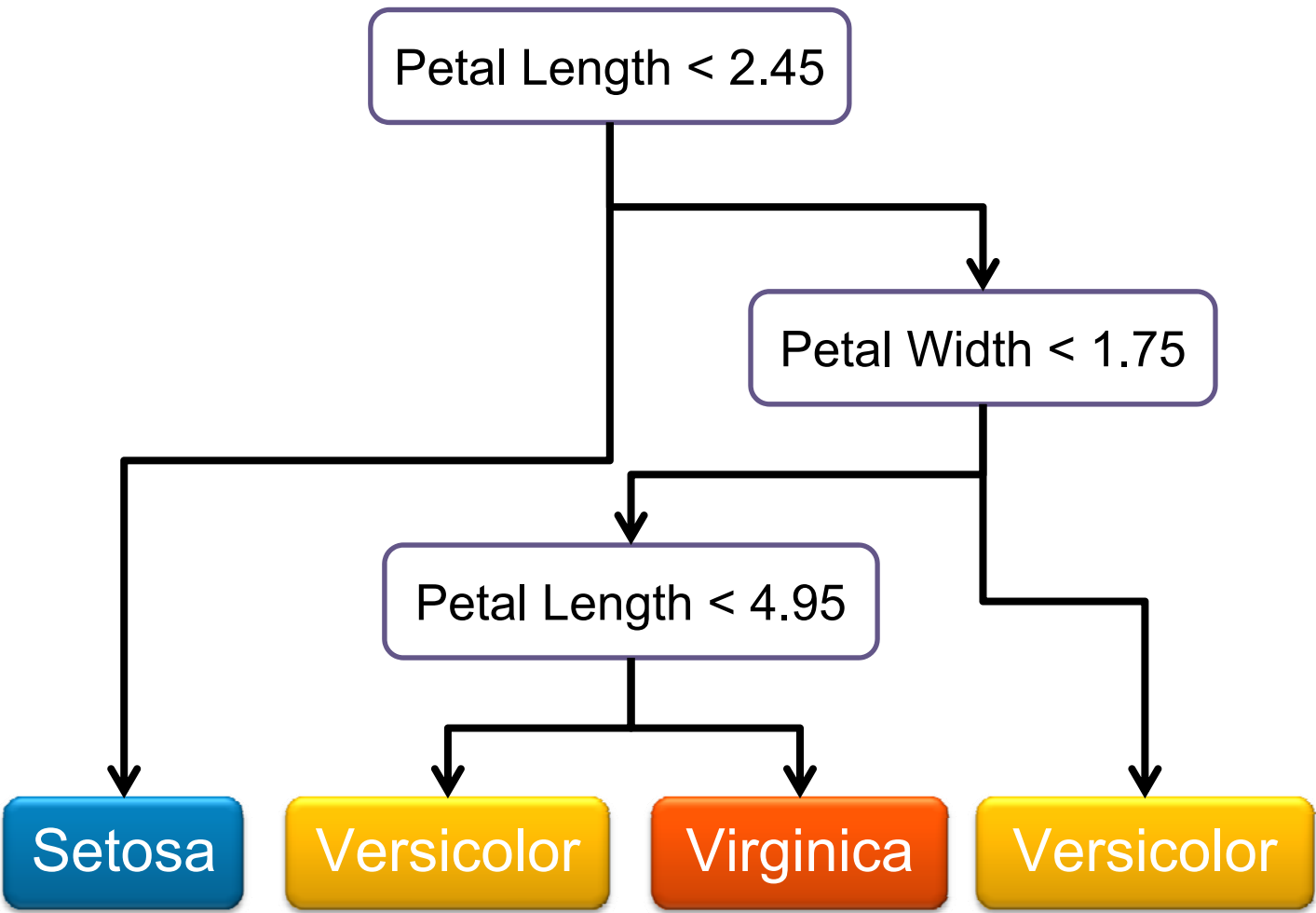
Setosa

Versicolor

Virginica



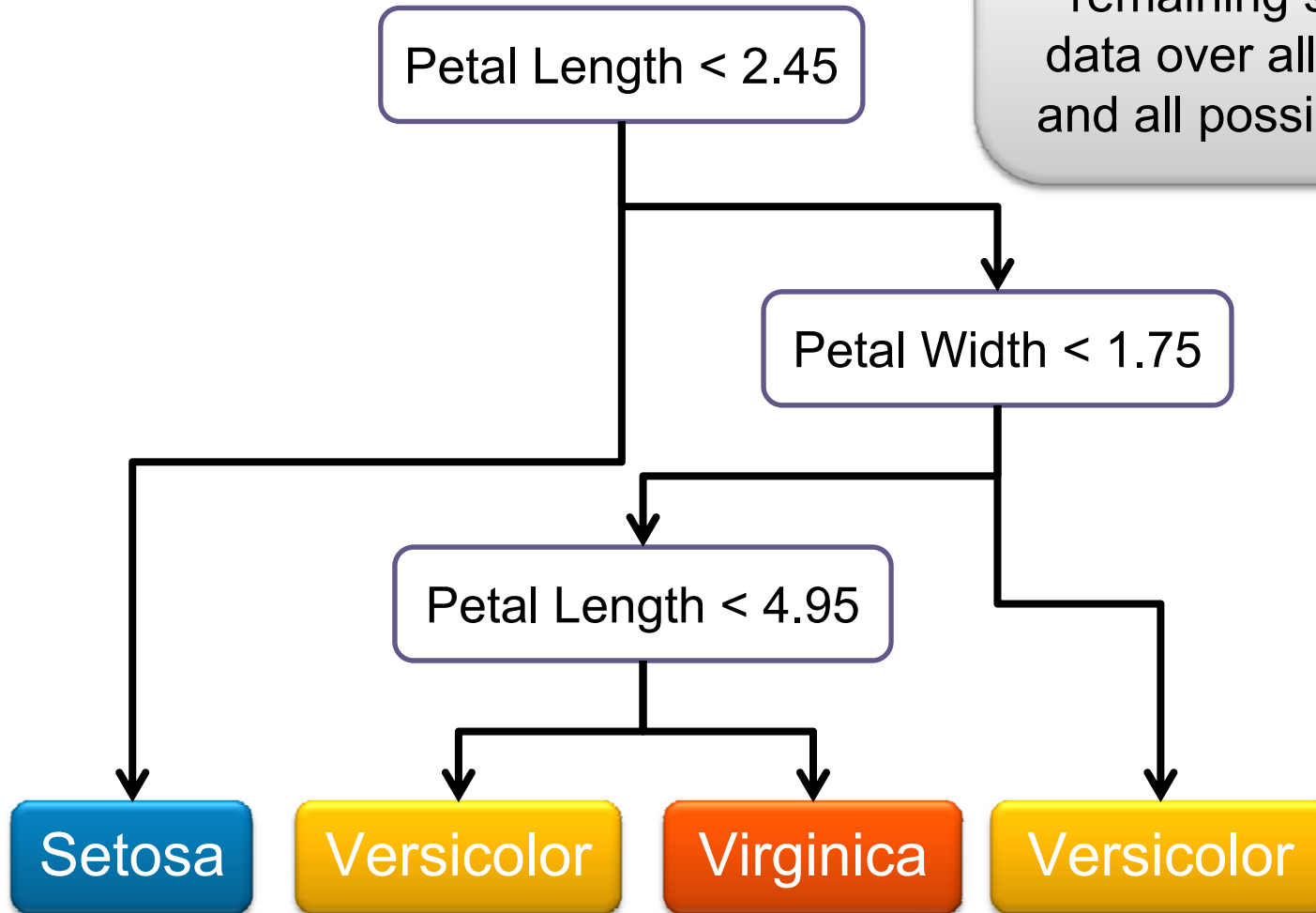
Decision Trees





Decision Trees

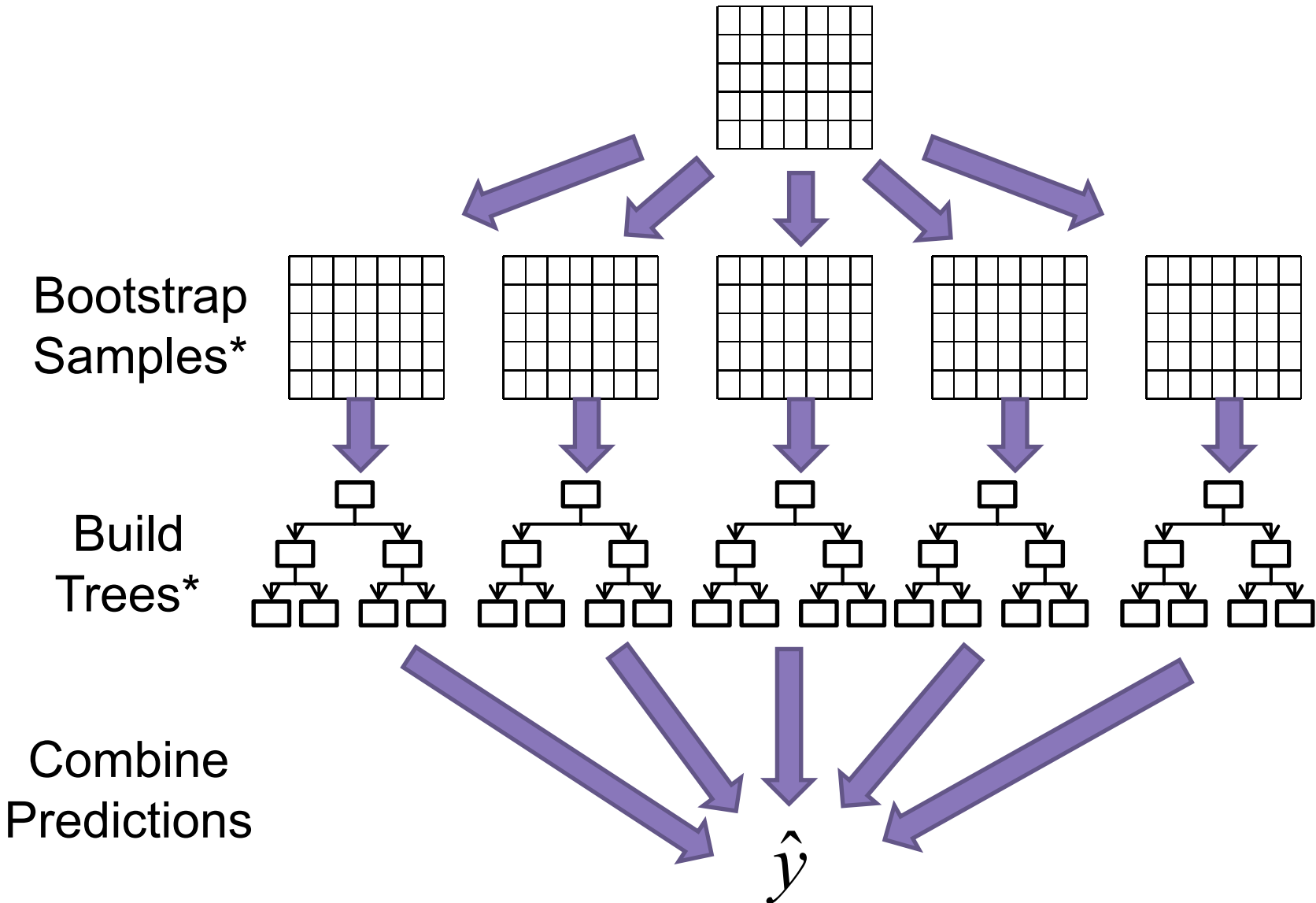
At each split, maximise the improvement in prediction with the remaining subset of data over all variables and all possible values





Random Forests

Lots of trees with 2 perturbations* of the data



Variable Importance Measures

2 Standard Measures

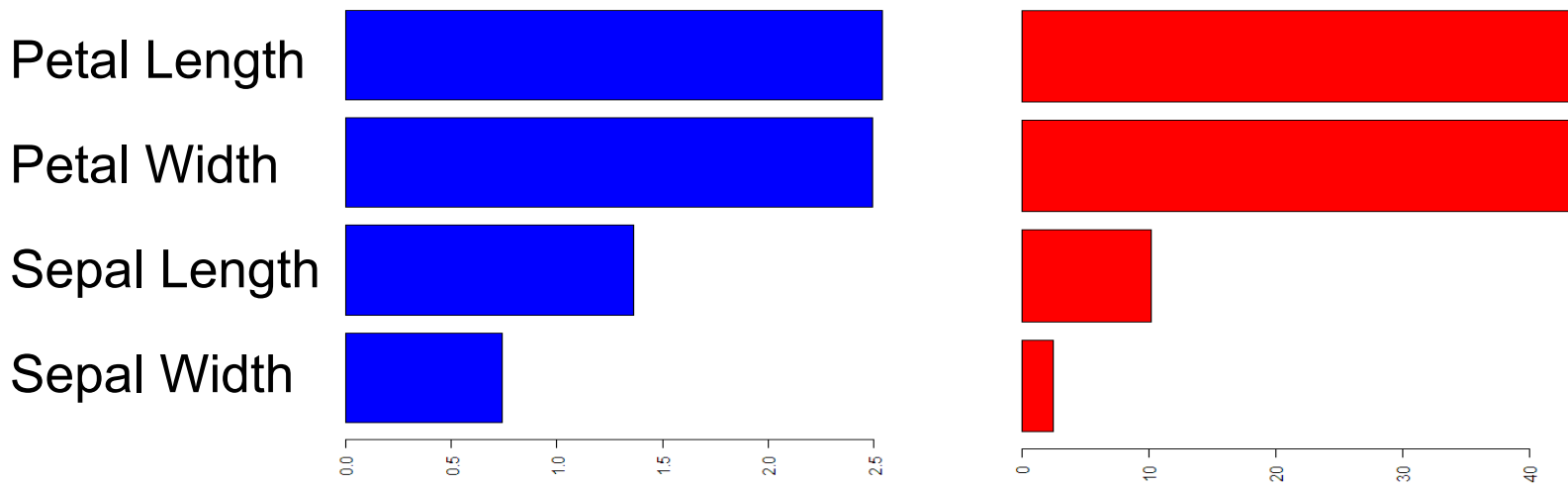
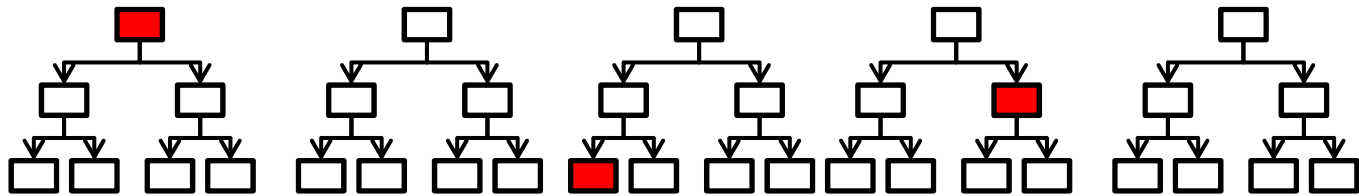


Permutation

- Change in accuracy of predictions when permuting each variable in turn

Reduction in Impurities

- Sum of the reduction in Gini index / MSE over all splits by that variable in the forest





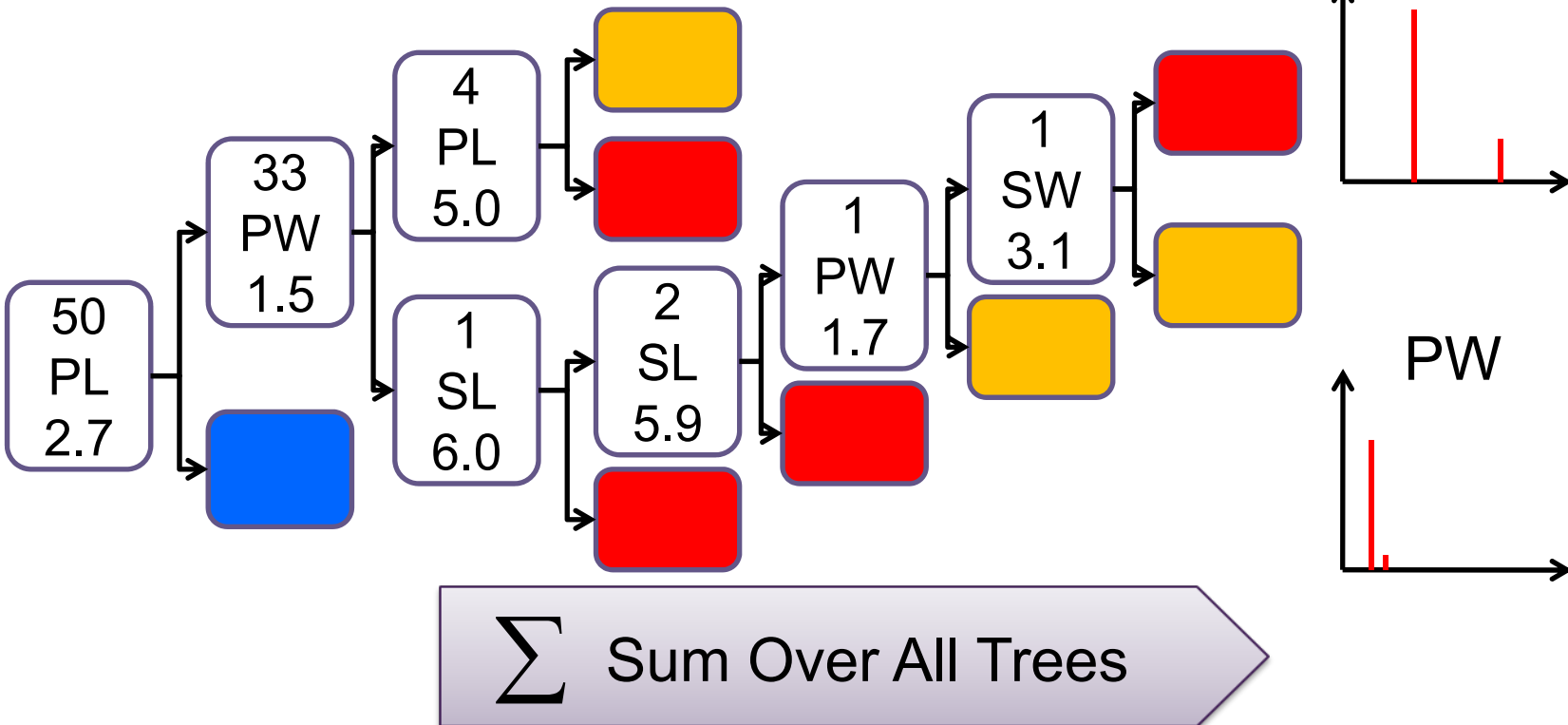
SURF :

Understanding Splits In Random Forests

Split Total Gini Score By Predictor Variables

Split Total Gini Score By Predictor Variables and their Split Values

Example Tree in Forest

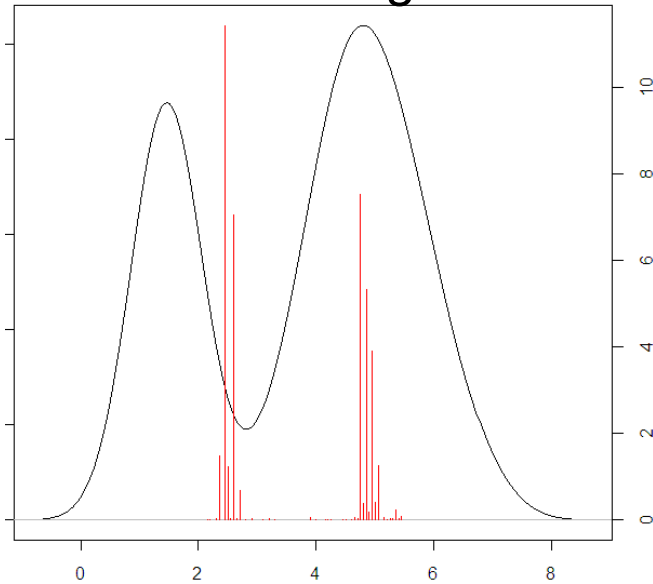


Σ Sum Over All Trees

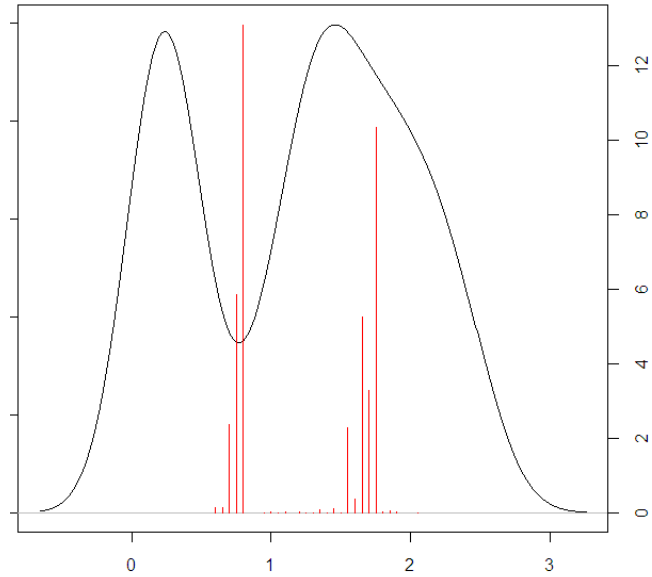
Univariate Results



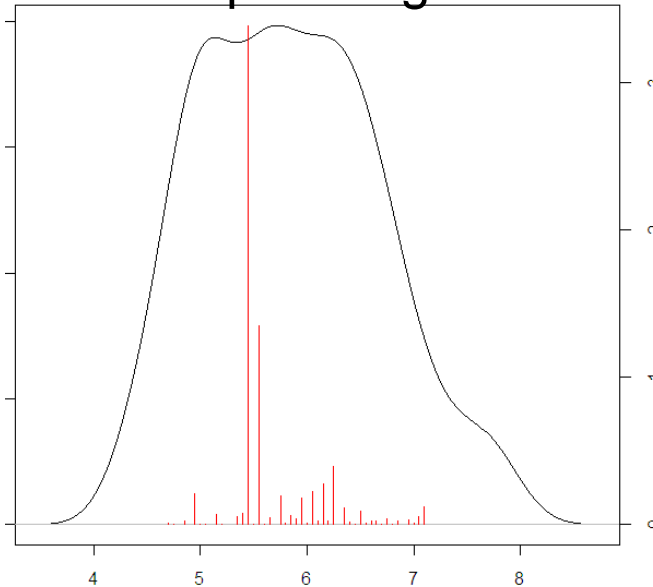
Petal Length



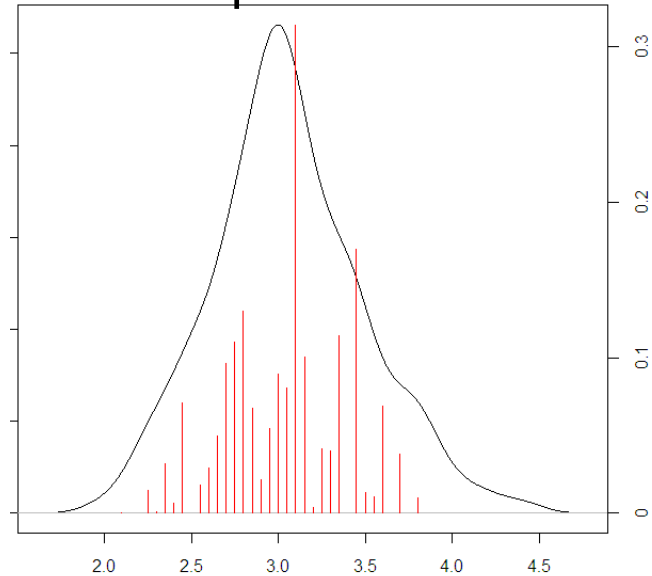
Petal Width



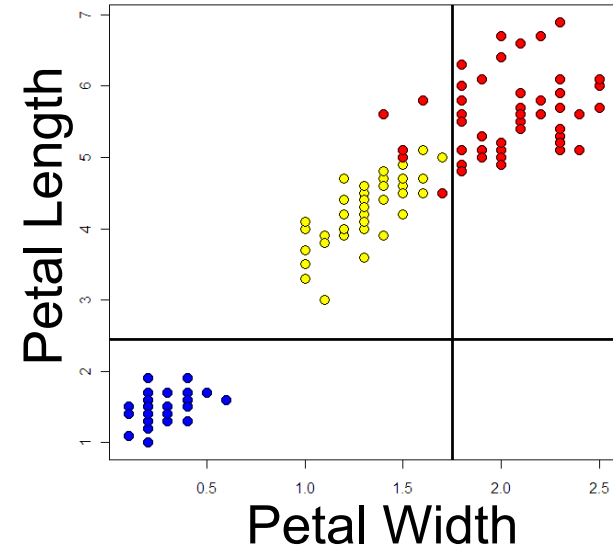
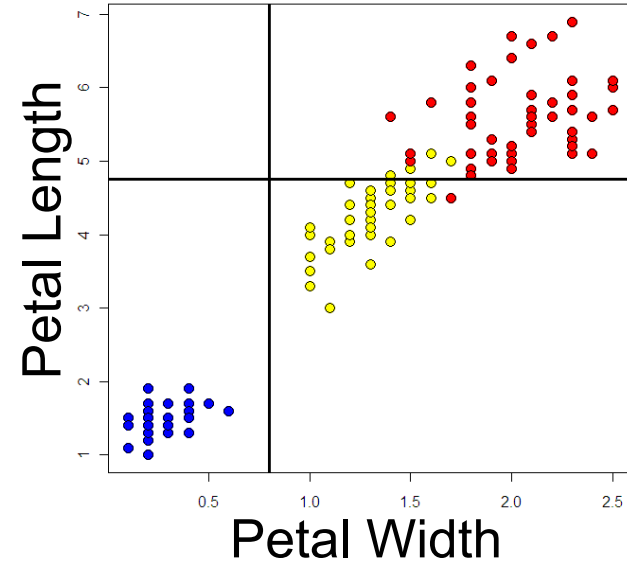
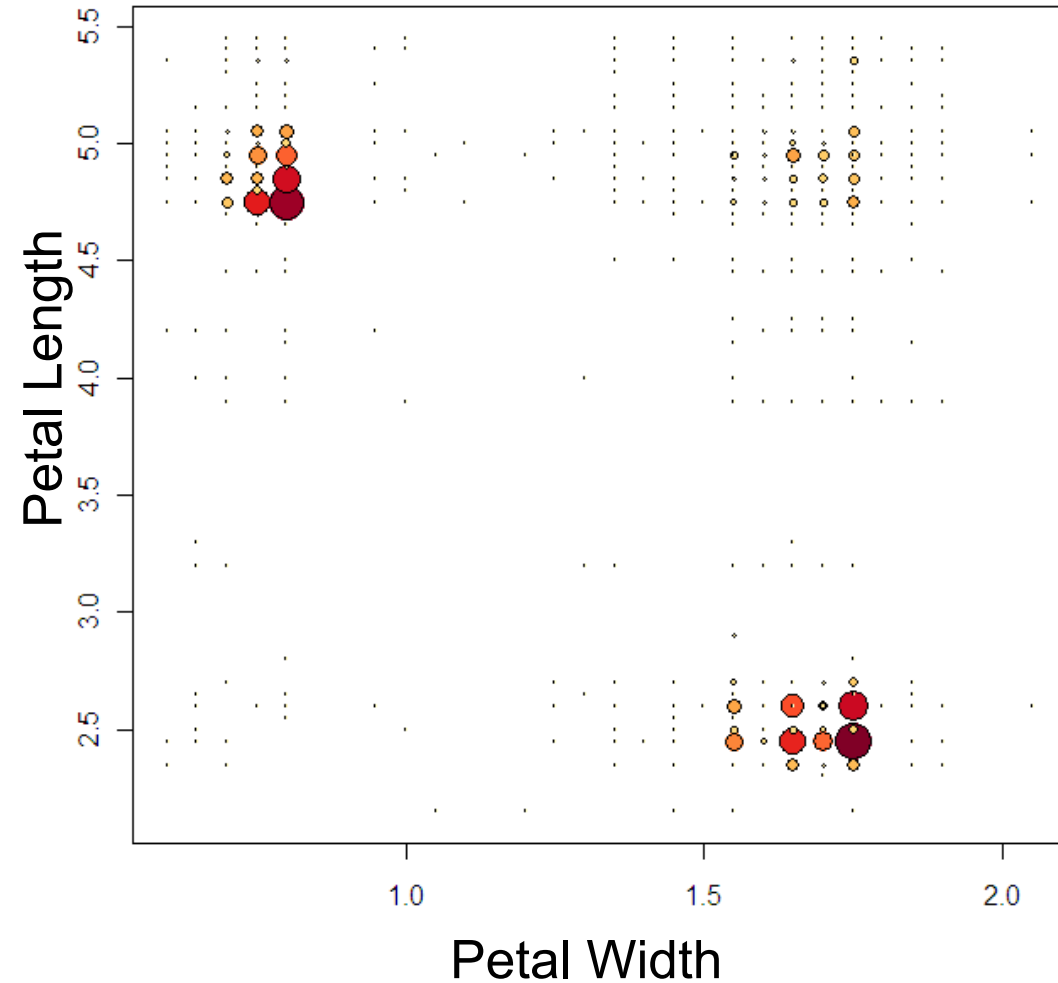
Sepal Length



Sepal Width

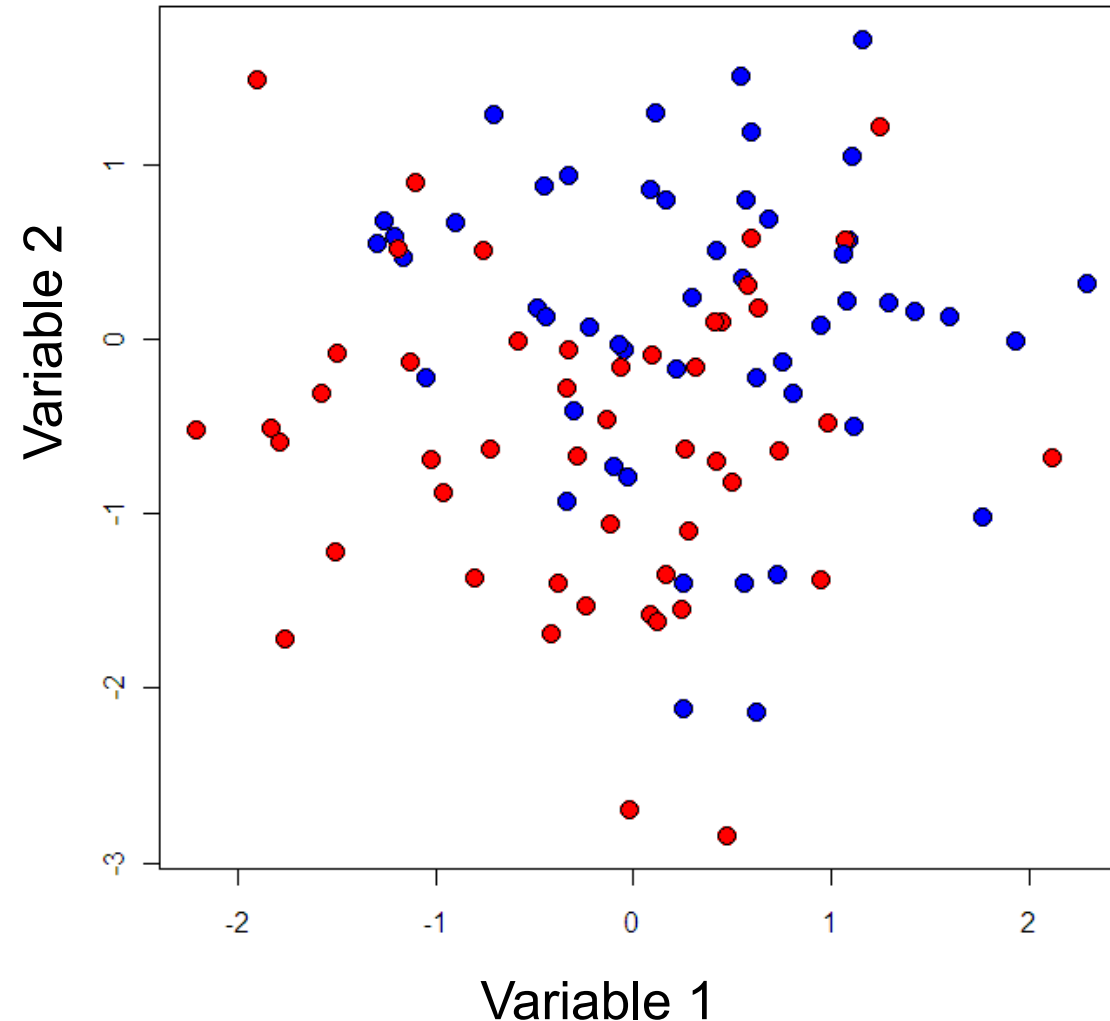


Pairwise Results



Simulated Data

Continuous Response

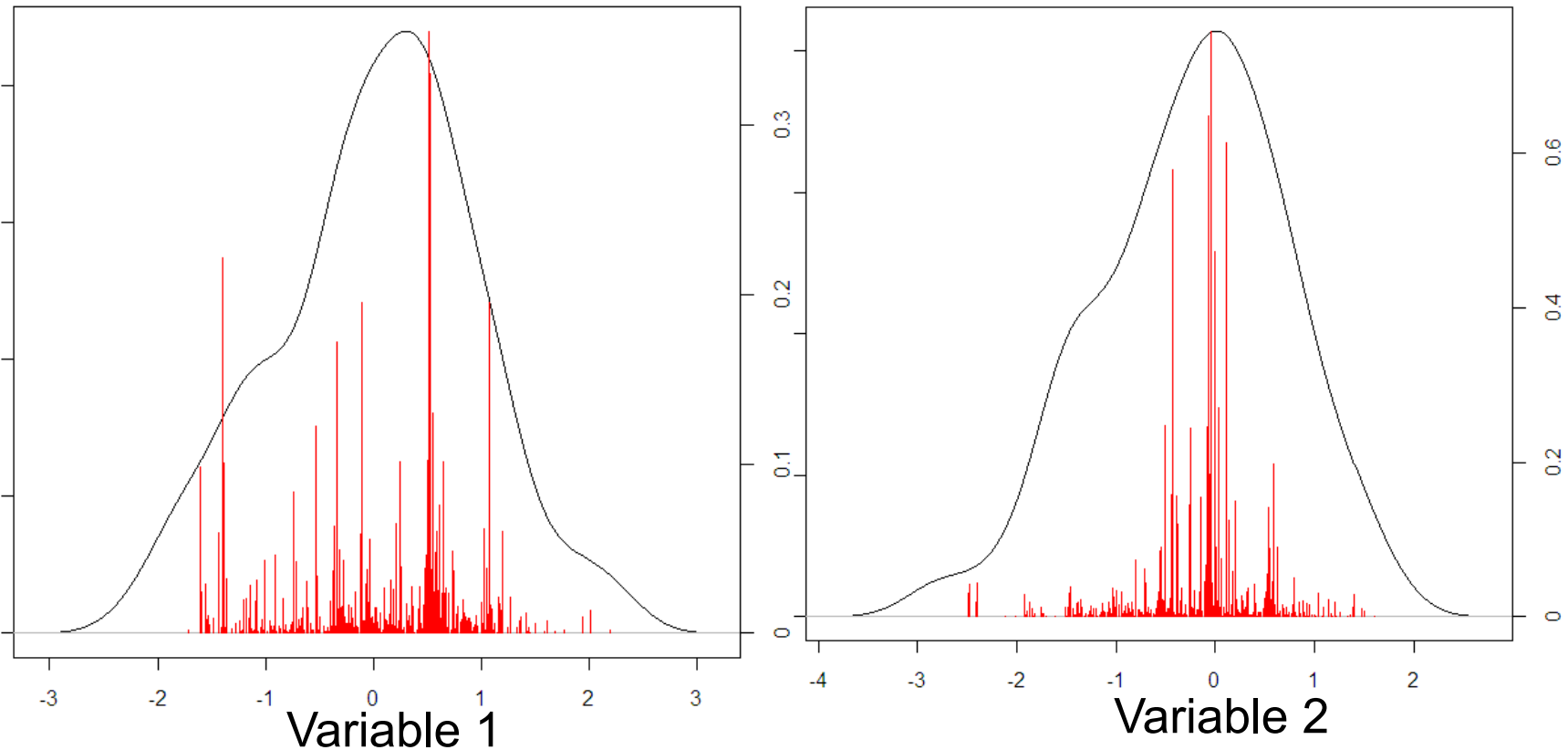


$$y \sim \text{Bin}(1, p)$$

$$\text{Logit}(p) = \text{Var1} + \text{Var2}$$

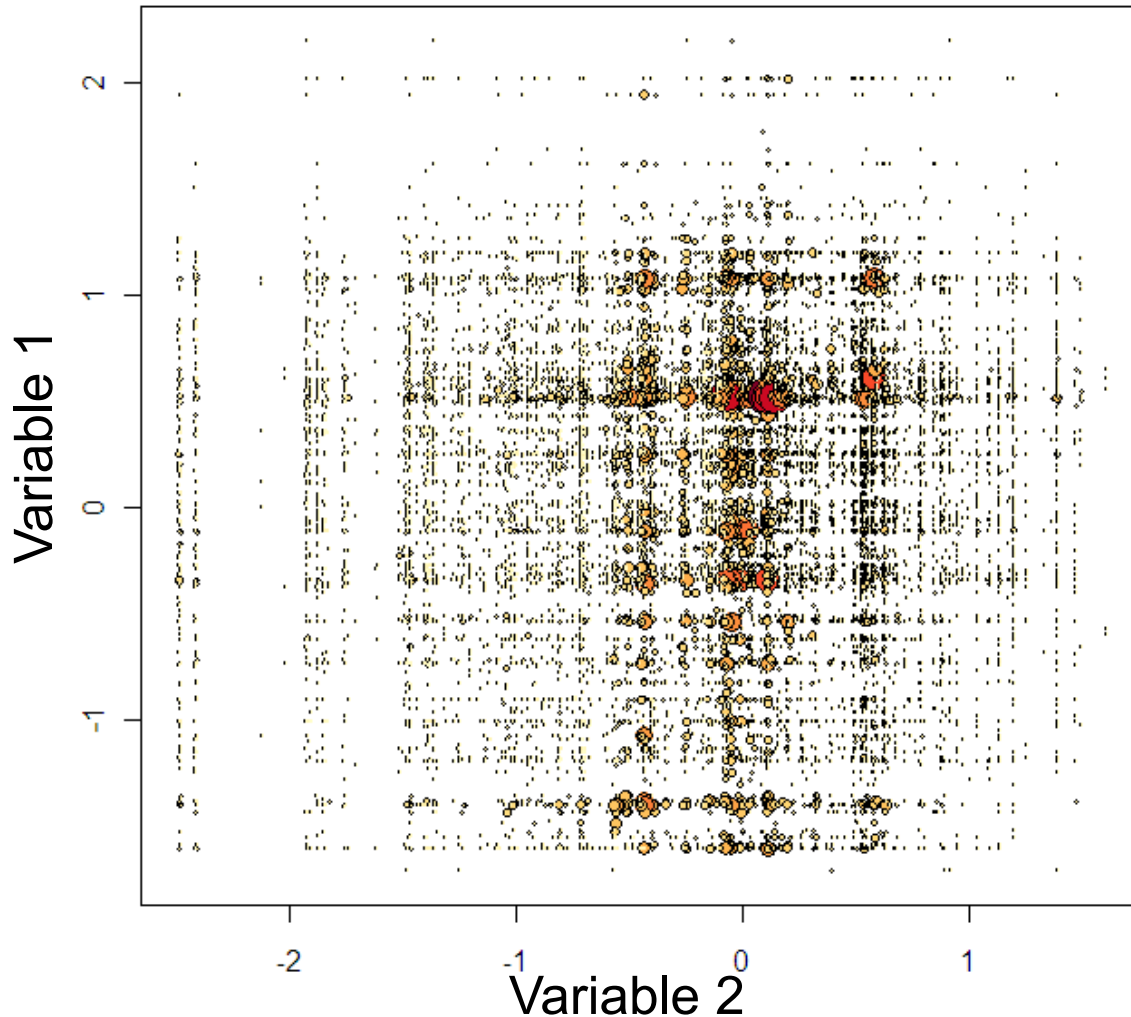
Four additional
non-correlated
noise variables
included in model

Simulated Data



Splits spread over a much more diffuse set of values

Simulated Data

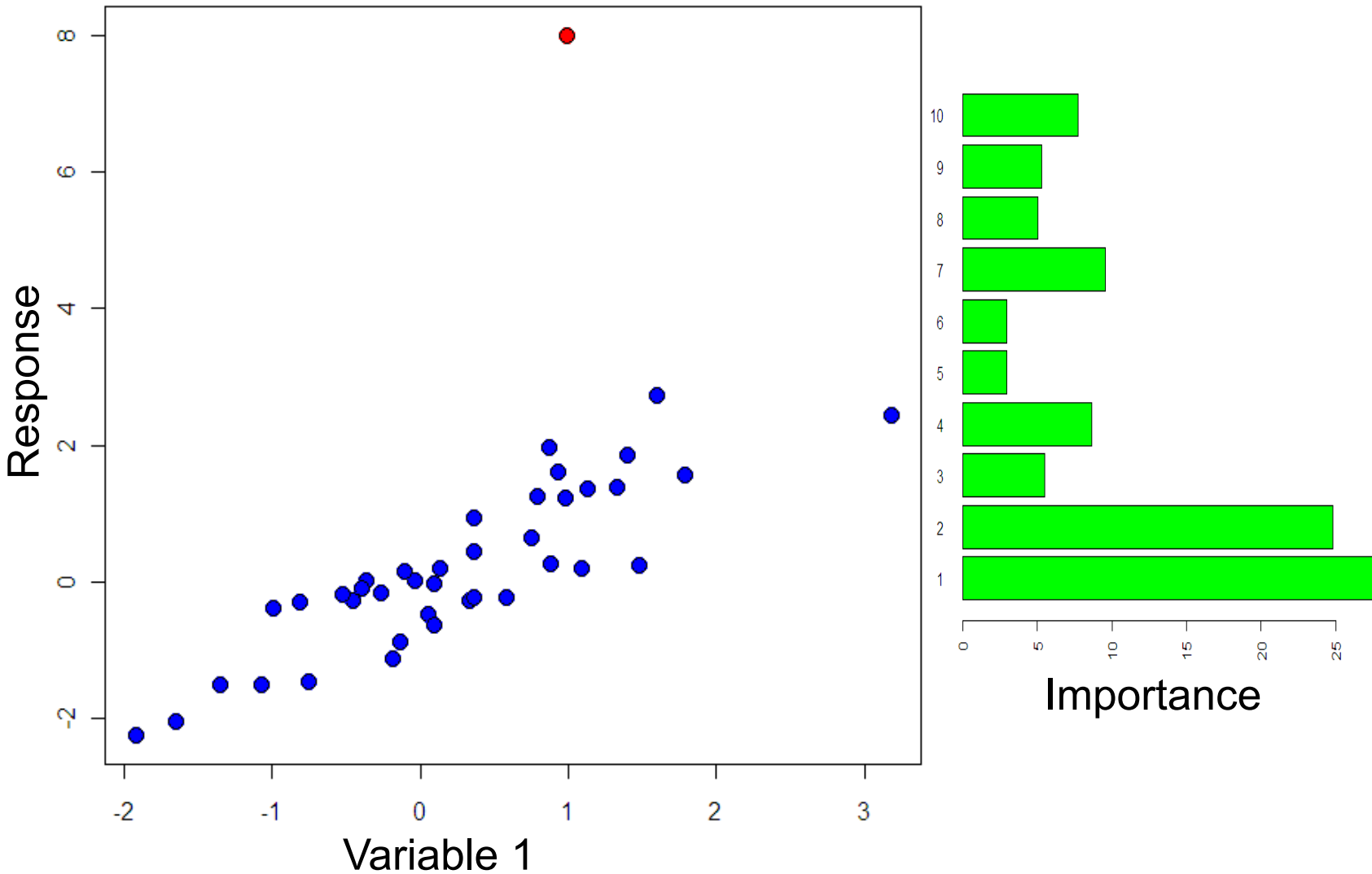


Splits spread over a much more diffuse set of values
Splits in the 2 variables acting independently from each other



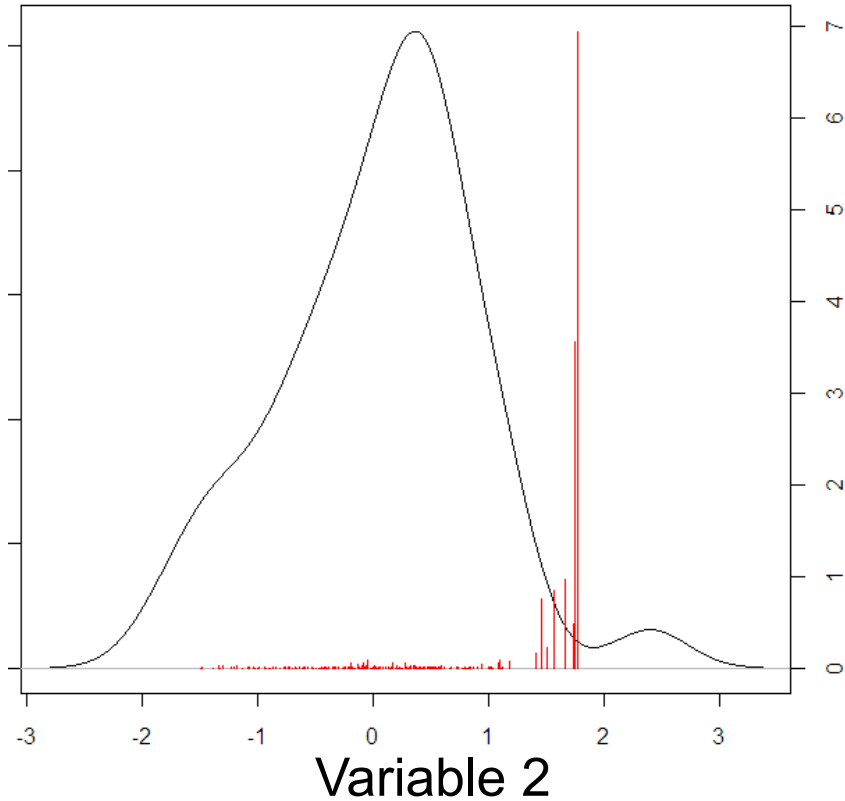
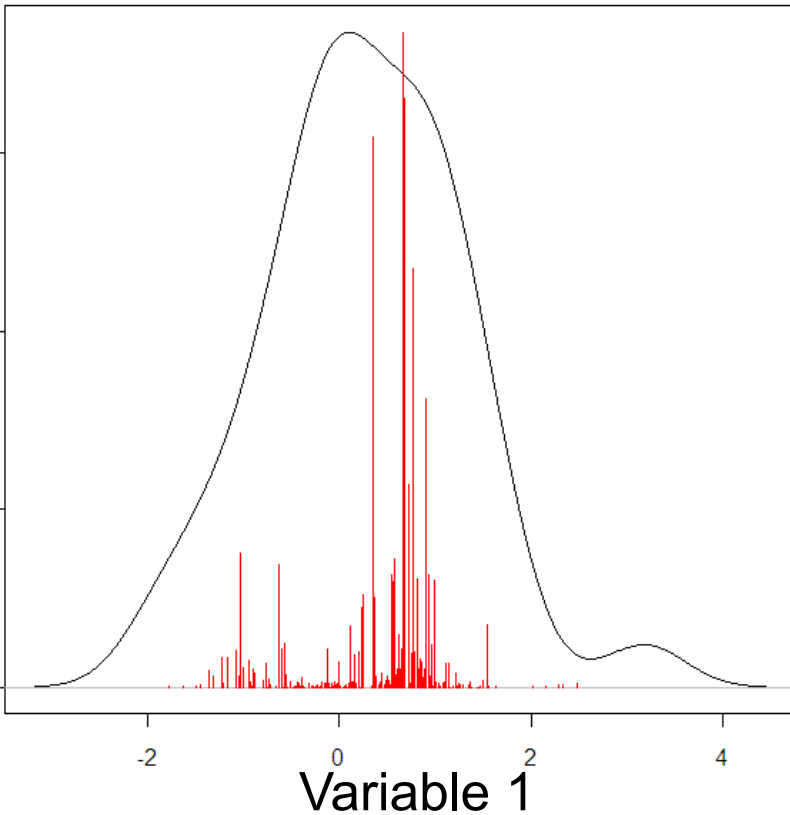


Simulated Data - Outliers





Simulated Data - Outliers



**Diffuse set of splits
centered around median**

**All splits separating
single outlying value**

Summary

- Random Forests are a powerful and robust multivariate modelling technique
- Beyond acting as a black-box predictor, random forests give insights into the underlying structure of the data through variable importance scores
- This can be taken a step further by considering the locations of the splits
- Also gives a valuable insight into the quality of the model

