

# Trees and Random Forests

## Trees

Decision trees may be familiar to many with a background in the social or medical sciences as convenient ways to represent data and can assist in decision making. Morgan and Sonquist (1963) derived a way for constructing trees motivated by the specific feature space of data collected from interviews and surveys. Unlike, say agricultural data which involves mostly numerical variables like rainfall, the data collected from interviews is mostly categorical. On top of this issue, the datasets Morgan and Sonquist dealt with had few participants,  $n$ , and much data collected on them,  $p$ . To continue with their list of difficulties, there was reason to believe that there were lurking errors in the variables that would be hard to identify and quantify. Lastly, many of the predictors were correlated and Morgan and Sonquist doubted that the additive assumptions of many models would be appropriate for this data. Morgan and Sonquist noted that while many statistical methods would have a difficulty accurately parsing this data, a clever researcher with quite a lot of time could create a suitable model simply by grouping values in the feature space and predicting that the response corresponding to these values would be the mean of the observed responses given the grouped conditions. Their formalization of this procedure in terms of “decision rules” laid the ground work for future research on decision trees.

Later researchers proposed new methods for creating trees that improved upon the Morgan and Sonquist model. Leo Breiman et al 1984 proposed an algorithm called CART, *classification and regression trees*, that would allow trees to be fit on various types of data. An alternative to this method is conditional inference trees. Torsten Hothorn, Kurt Hornik, Achim Zeileis argue in their 2006 paper **Unbiased Recursive Partitioning: A Conditional Inference Framework**, CART has a selection bias toward variables with either missing values or a great number of possible splits. This bias can effect the interpretability of all tree models fit using this method. As an alternative to CART and other algorithms, Hothorn et al propose a new method, conditional inference trees.

There is a limit to the predictive capabilities of a single tree as they suffer from high variance. To alleviate this, random forests are often used instead. They function by enlisting the help of many trees, and then by aggregating the responses over all of them but with a subtle trick that ensures the trees will be independent of each other. At each split only  $m$  variables are considered as possible candidates.

## What We Mean When We Talk About Inference

### Inferential vs Descriptive Statistics

A note should be made of the difference between inferential and descriptive statistics. This paper’s aim is to describe a process of making inferential claims using random forests, not descriptive ones. Descriptive statistics describe the data at hand without making any reference to a larger data generating system that they come from. It follows that inferential statistics then make claims about the data generating system given the data at hand.

—Frequentist vs Bayesian—

—There is some debate about interpreting inferential statistics. On one hand, we have the Bayesian model—

*Need a better way to discuss inference than Bayes/frequentist*

## Permutations and Populations

As stated in the introduction of the **Chronical of Permutations Statistical Methods** by KJ Berry et al, 2014, there are two models of statistical inference. One is the population model, where we assume that the data was randomly sampled from one (or more) populations. Under this model, we assume that the data

generated follows some known distribution. “Under the population model, the level of statistical significance that results from applying a statistical test to the results of an experiment or a survey corresponds to the frequency with which the null hypothesis would be rejected in repeated random samplings from the same specified population(s)”, (Berry et al, 2014).

The permutation family of methods, on the other hand, only assumes that the observed result was caused by experimental variability. The test statistics is first calculated for the observed data, then the data is permuted a number of times. The statistic is calculated after each permutation to derive a distribution of possible values. Then the original test statistic is tested against this distribution. If it is exceptionally rare, then there is evidence that our observation was not simply experimental variability.

## Inference on Random Forests

### The Problem

Random forests create models with great predictive-, but poor inferential capabilities.

### Proposed solutions to this problem

Statisticians Breiman and Cutler proposed a method of permuted variable importance to answer this problem. Their method compares the variable importance for each variable in a tree-wise manner. For each tree, the permuted variable importance of the variable  $X_j$  is:

$$PV^t(x_j) = \frac{\sum_{i \in |B|} y - \hat{y}^t}{|B|} - \frac{\sum_{i \in |*B|} y - \hat{*y}^t}{|*B|}$$

Where  $B$  is the matrix representing the feature space,  $|B|$  is the number of observations,  $*B$  is the matrix of predictors but with  $X_j$  permuted,  $\hat{y}$  is the predicted outcome, and  $\hat{*y}^t$  is the predicted outcomes after variable  $X_j$  has been permuted. This value is averaged over all the trees. It’s important to note that if the variable  $X_j$  is not split on in the tree  $t$ , the tree-wise variable importance will be 0.

Creating a permutation-based method is certainly an attractive solution to our problem. One, it allows us to estimate the distribution of variable importance and generate a Z score under the null hypothesis that  $PV = 0$ .

$$PV(x_j) = \frac{\sum_1^n treePV^t(x_j)}{\frac{\hat{\sigma}}{\sqrt{ntree}}}$$

Strobl et al from the University of Munich criticize this method in their 2008 technical report, *Danger: High Power! – Exploring the Statistical Properties of a Test for Random Forest Variable Importance*. One, this method has the downside of increasing power with increasing numbers of trees in the forest. This is a more or less arbitrary parameter which we would hope would not affect our importance estimates. Secondly, the null hypothesis under Breiman and Cutler’s strategy is that the variable importance  $V$  for any variable  $X_j$  is not equal to zero given  $Y$ , the response. Because random forests are most often used in situations with multicollinearity that would make other methods like the linear model difficult, Strobl argues that any variable importance measure worth its salt should not be mislead by correlation within the predictors.

The researchers at the University of Munich published a fully fleshed response to the Breiman and Cutler method in 2008, titled *Conditional Variable Importance for Random Forests* that address these issues. Strobl et al propose restructuring the Breiman and Cutler algorithm to account for conditional dependence among the predictors. Their algorithm looks like this:

1. Fit a random forest to the model,  $R_0$ , and calculate base variable importance for each variable  $V$

2. For every predictor  $X_j \in X_1, \dots, X_n$ :
  - 2a. Conditionally permute  $X_j$  given the splits found in  $R_0$
  - 2b. Fit a new random forest  $R_j$  with the permuted data
  - 2c. Calculate a new variable importance  $\hat{V}_j$
3. For every variable  $X_1, \dots, X_n$ ,

$$CV(X_j) = \hat{V}_j - V_j$$

The null hypothesis is that  $CV(X_j) = 0$  given the predictor  $Y$  and all other predictors  $X_1, \dots, X_n$ . This accounts for interactions between  $X_j$  and the other predictors. Using the simulated data from the previous example, here's an implementation of the algorithm outlined here as it is in the **party** package.

This paper aims to provide a response to this method. One the conditional permutation algorithm is notoriously slow with any dataset of a size that is appropriate for a random forest. Two, the partitions are made from the random forest corresponding to the formula of  $Y \sim X_1, \dots, X_n$  instead of a model of  $X_j \sim X_1, \dots, X_n$ .