

# Statistical Inference on Random Forests

---

A Thesis  
Presented to  
The Division of Mathematics and Natural Sciences  
Reed College

---

In Partial Fulfillment  
of the Requirements for the Degree  
Bachelor of Arts

---

Aurora Owens

May 2017



Approved for the Division  
(Mathematics)

---

Andrew Bray



# Acknowledgements

I want to thank a few people.



# Preface

This is an example of a thesis setup to use the reed thesis document class.





# Table of Contents

|   |           |
|---|-----------|
| <b>Chapter 1: Introduction</b>                                  | <b>1</b>  |
| 1.1 Trees and Random Forests                                    | 1         |
| 1.1.1 Trees   | 1         |
| 1.2 What We Mean When We Talk About Inference                   | 2         |
| 1.2.1 Inferential vs Descriptive Statistics                     | 2         |
| 1.3 Permutations and Populations                                | 2         |
| 1.4 Inference on Random Forests                                 | 3         |
| 1.4.1 The Problem   | 3         |
| 1.4.2 Proposed solutions to this problem                        | 3         |
| <b>Chapter 2: Simulations and Comparisons</b>                   | <b>5</b>  |
| 2.1 Simulated Data  | 5         |
| 2.1.1 Figure 8:   | 12        |
| 2.2 Models and Comparisons                                      | 13        |
| 2.2.1 Figure 9:   | 13        |
| 2.2.2 Figure 10:  | 14        |
| 2.3 Bagged Forests  | 15        |
| 2.3.1 Figure 11:  | 15        |
| 2.4 Random Forests  | 16        |
| 2.4.1 Figure 12:  | 16        |
| 2.4.2 TO DO   | 17        |
| <b>Chapter 3: Random Forest Variable Importance</b>             | <b>19</b> |
| 3.1 Breiman et al Introduce Permuted Variable Importance (1984) | 19        |
| 3.1.1 Variable Importance on a Single Tree                      | 19        |
| 3.1.2 Variable Importance for a Random Forest                   | 20        |
| 3.2 Strobl et al Respond (2008)                                 | 20        |
| 3.3 Inferential Variable Importance                             | 23        |
| <b>Chapter 4: INFTrees and INFforests Variable Importance</b>   | <b>25</b> |
| 4.1 Theory  | 25        |
| 4.2 Implementation  | 25        |
| <b>Conclusion</b>   | <b>27</b> |

|   |           |
|---|-----------|
| <b>Appendix A: The First Appendix . . . . .</b>         | <b>29</b> |
| <b>Appendix B: The Second Appendix: CTree . . . . .</b> | <b>31</b> |
| B.0.1 Conditional Inference Trees . . . . .             | 31        |
| <b>References . . . . .</b>                             | <b>33</b> |

# List of Tables



# List of Figures



# Abstract

The preface pretty much says it all.

Second paragraph of abstract starts here.





# Dedication

You can have a dedication here if you wish.



# Chapter 1

## Introduction

### 1.1 Trees and Random Forests

#### 1.1.1 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 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. Random forests and their algorithms will be discussed at length in Chapter 2.

## 1.2 What We Mean When We Talk About Inference

### 1.2.1 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*

## 1.3 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.

## 1.4 Inference on Random Forests

### 1.4.1 The Problem

Random forests create models with great predictive-, but poor inferential capabilities. After Morgan and Sonquist initial development of decision trees, they quickly moved to the domain of machine learning and away from statistics, thus, researchers focused on bettering predictions and improving run times and less on the statistics behind them. Inferential statistics with random forests is usually treated as a variable selection problem, and generally falls behind the predictions in importance. This has limited the applications of random forests in certain fields, as to many the question of “why” the data is the way it is, is just, if not more, important as the predictions. There are several means of performing descriptive statistics with random forests that could be interpreted incorrectly as attempting to answer this, namely base variable importance, but without a statistically backed method for performing variable importance, the use of random forest is limited to prediction-only settings.

### 1.4.2 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 misled 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:

---

**Algorithm 1** Conditional Variable Importance for Random Forests

---

- 1: Fit a random forest to the model,  $R_0$ , and calculate base variable importance for each variable  $V$
- 2: **for** every predictor  $X_j \in X_1, \dots, X_n$  **do**
- 3:     Conditionally permute  $X_j$  given the splits found in  $R_0$
- 4:     Fit a new random forest  $R_j$  with the permuted data
- 5:     Calculate a new variable importance  $\hat{V}_j$
- 6: **end for**
- 7: 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  $X_j \sim X_1, \dots, X_n$ .

# Chapter 2

## Simulations and Comparisons

### 2.1 Simulated Data

Tree-based methods shine in situations with correlated predictors, although these situations can pose problems for inference. In a situation with correlated predictors  $X_1$  and  $X_2$ , and the model we are considering is  $Y \sim X_1 + X_2$ , it is difficult to say how much of the modeled effect on  $Y$  is due to  $X_1$  or  $X_2$ . To illustrate this idea, compare a few existing methods, and explore methods of inference on tree based models three datasets will be simulated with different correlation structures. We will be focused more on the correlation structure between the predictors than on their relationships with the response and this will be reflected in the simulations.

To aid in comparisons between the methods, one of the simulated datasets considered in this paper will be generated from the same method as used in (Strobl et al, 2008???). Under this method, the 13 x 1000 data set,  $D_1$ , has 12 predictors,  $V_1, \dots, V_{12}$ , where  $V_j \sim N(0, 1)$ . The first four are, however, block correlated to each other with  $\rho = .9$ . They are related to  $Y$  by the linear equation:

$$Y = 5 \cdot V_1 + 5 \cdot V_2 + 2 \cdot V_3 + 0 \cdot V_4 + -5 \cdot V_5 + -5 \cdot V_6 + 0 \cdot V_7 + 0 \cdot \dots + E, E \sim N(0, \frac{1}{2})$$

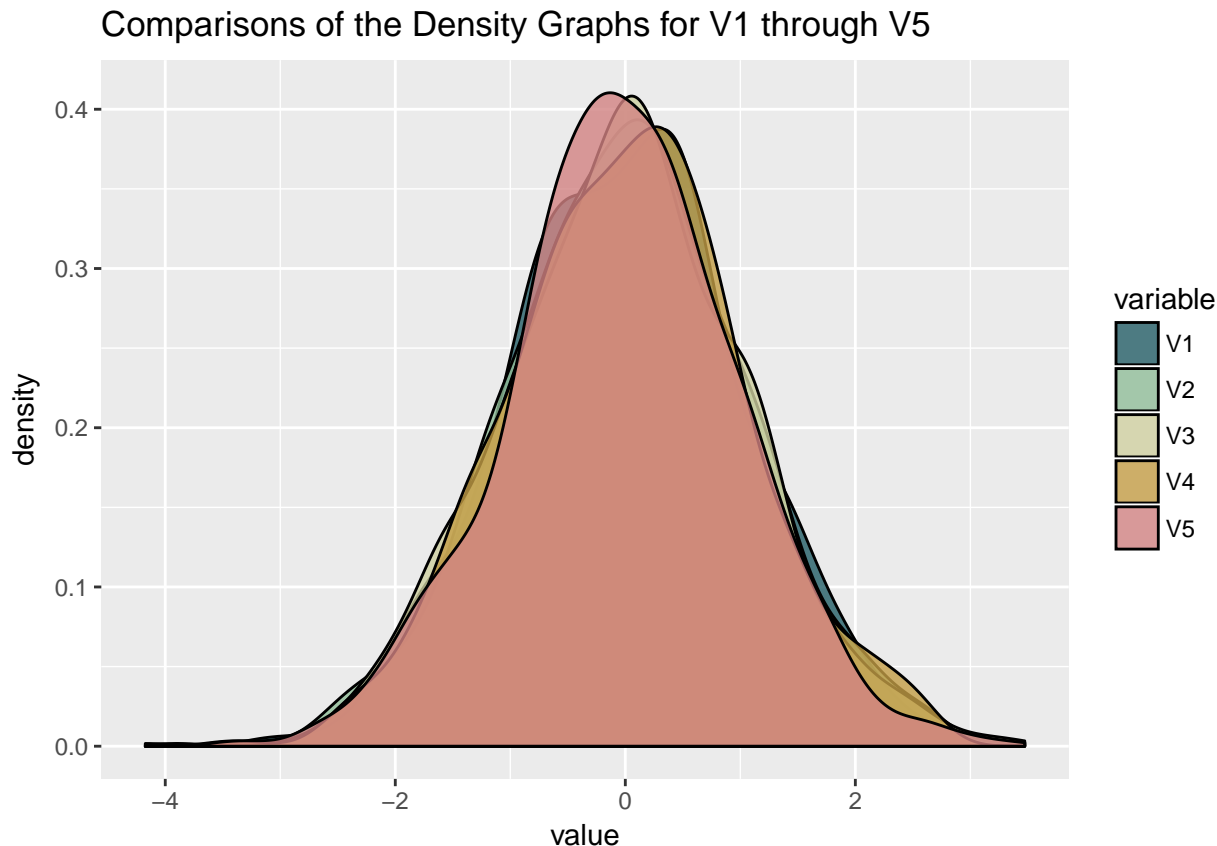
Note that the coefficients for  $V_7, \dots, V_{12}$  are all zero.

**Table 1: Correlation of  $V_1, \dots, V_7$  and  $Y$**

|    | V1     | V2     | V3     | V4     | V5     | V6     | V7     | y      | beta |
|----|--------|--------|--------|--------|--------|--------|--------|--------|------|
| V1 | 1.000  | 0.915  | 0.908  | 0.907  | -0.034 | 0.006  | 0.012  | 0.829  | 5    |
| V2 | 0.915  | 1.000  | 0.914  | 0.914  | -0.020 | -0.001 | -0.001 | 0.830  | 5    |
| V3 | 0.908  | 0.914  | 1.000  | 0.903  | -0.017 | -0.007 | 0.007  | 0.808  | 2    |
| V4 | 0.907  | 0.914  | 0.903  | 1.000  | -0.002 | -0.015 | 0.023  | 0.789  | 0    |
| V5 | -0.034 | -0.020 | -0.017 | -0.002 | 1.000  | 0.044  | 0.005  | -0.388 | -5   |
| V6 | 0.006  | -0.001 | -0.007 | -0.015 | 0.044  | 1.000  | -0.005 | -0.364 | -5   |
| V7 | 0.012  | -0.001 | 0.007  | 0.023  | 0.005  | -0.005 | 1.000  | -0.141 | -2   |

As can be seen from the last column in the table, “beta”, although  $V_4$  was not included in the model  $Y \sim V_1, \dots, V_{12}$ , its’ strong correlation with more influential predictors  $V_1, \dots, V_3$  insures that it still shows a strong linear correlation with  $Y$ . A linear model would likely *overstate* the effect of  $V_4$  on  $Y$ .<sup>12</sup>

**Figure 1:**



As can be seen above in Figure 1 the densities of  $V_1, \dots, V_5$  are all very similar due to the way they were generated.

<sup>1</sup>A brief note on uncertainty is needed here. It’s true that in this setting we can say that  $V_4$  is actually unimportant to understanding  $Y$ , but in situations with real data this is profoundly more difficult to parse. Often like in the social science situations that Morgan and Sonquist encountered, the real relationship between correlated predictors is complicated and often there is some theoretical backing or other insight that is gained to include variables that may not be important to the model.

<sup>2</sup>Another point that could be said is that, no  $V_4$  is not unimportant,  $V_1, V_2$ , and  $V_3$  are just stand ins for the real star,  $V_4$ , as they are nearly the same ( $\rho \sim 1$ ). Then the real relationship represented here is  $Y \sim (5 + 5 + 2) \cdot V_4 + -5 \cdot V_5 + -5 \cdot V_6 + -2 \cdot V_7$ . This model is not unsuccessful in capturing the structure of the data, and this is typically the practice used to model data with highly correlated predictors. If this seems philosophically satisfying to you, the rest of this thesis may seem a bit inconsequential. I apologize.



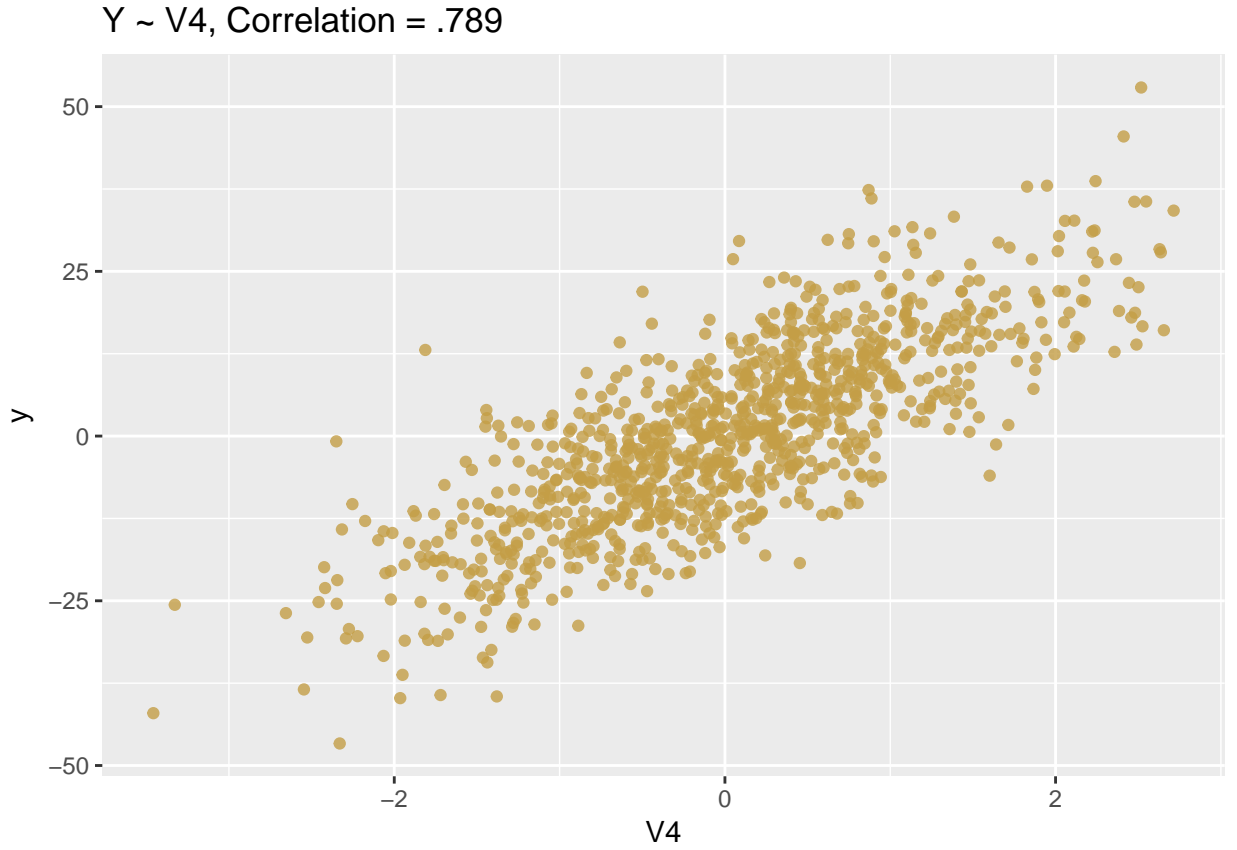
**Figure 2:**

Figure 2 is an illustration of the relationship between  $Y \sim V_4$  with linear correlation of .789.

While  $D_1$  represents a situation with linear correlation between the predictors,  $D_2$  does not. Here, the model is the same,  $Y \sim X_1, \dots, X_{12}$  where  $Y$  is generated according to the equation:

$$Y = 5 \cdot X_1 + 5 \cdot X_2 + 2 \cdot X_3 + 0 \cdot X_4 + -5 \cdot X_5 + -5 \cdot X_6 + 0 \cdot X_7 + 0 \cdot \dots + E, E \sim N(0, \frac{1}{2})$$

However, instead of block correlation with  $\rho = .9$ , four variables are related to each other by the equations below. Note that  $X_1, X_5, \dots, X_{12} \sim N(0, 1)$

$$X_2 = X_1 + E, E \sim \text{Exponential}(1)$$

$$X_3 = X_2 + E, E \sim \text{Exponential}(1)$$

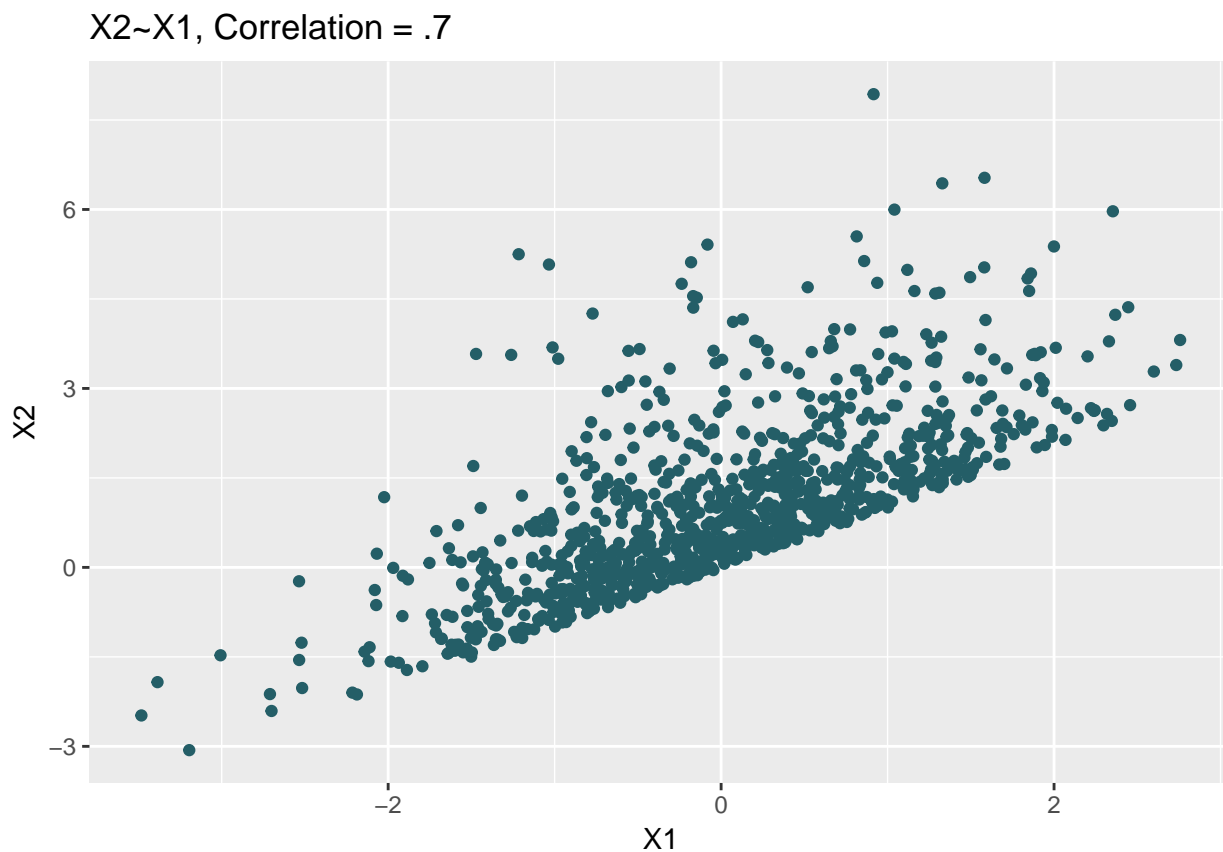
$$X_4 = X_3 + E, E \sim \text{Exponential}(1)$$

**Table 2: Correlation of  $X_1, \dots, X_7$  and  $Y$**

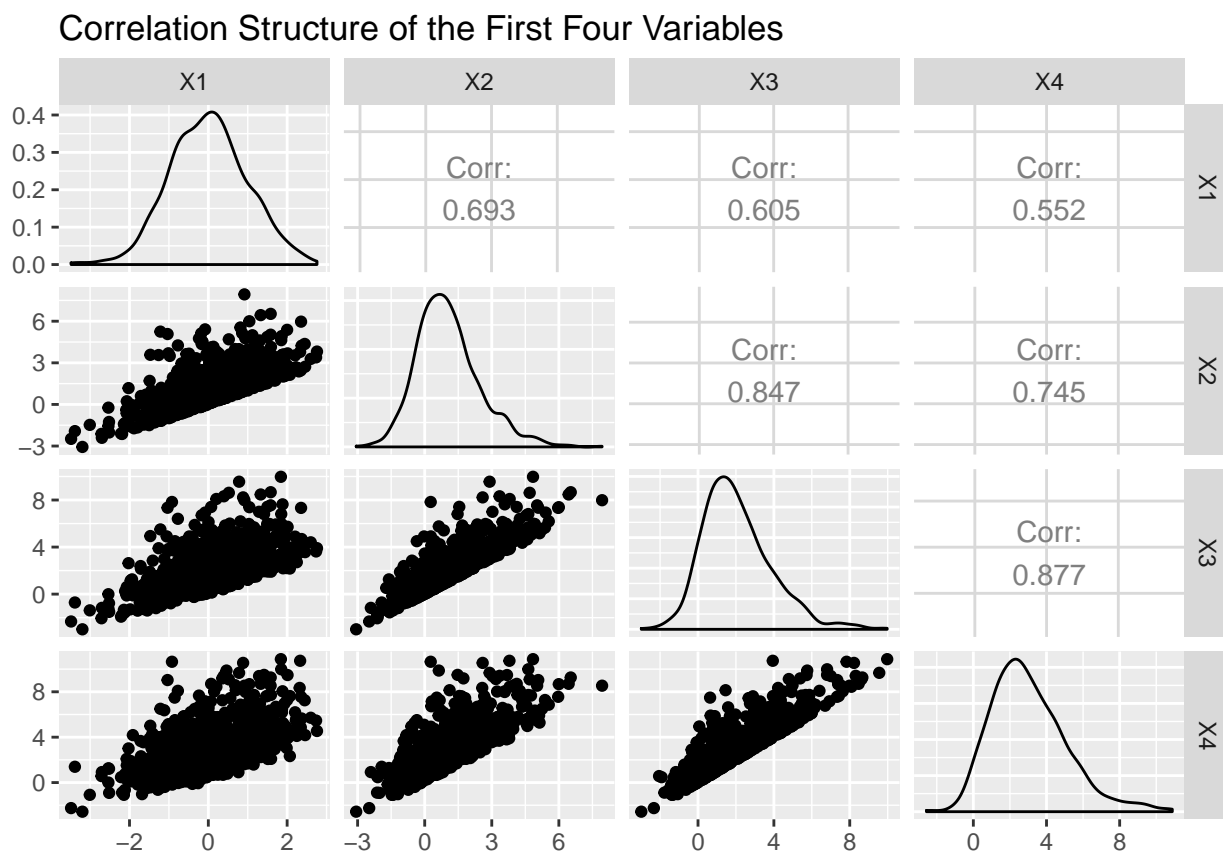
|    | X1     | X2     | X3     | X4     | X5     | X6     | X7     | y      | beta |
|----|--------|--------|--------|--------|--------|--------|--------|--------|------|
| X1 | 1.000  | 0.693  | 0.605  | 0.552  | -0.043 | 0.009  | -0.006 | 0.760  | 5    |
| X2 | 0.693  | 1.000  | 0.847  | 0.745  | 0.004  | 0.006  | -0.018 | 0.845  | 5    |
| X3 | 0.605  | 0.847  | 1.000  | 0.877  | 0.007  | 0.005  | -0.024 | 0.785  | 2    |
| X4 | 0.552  | 0.745  | 0.877  | 1.000  | 0.011  | 0.006  | -0.032 | 0.696  | 0    |
| X5 | -0.043 | 0.004  | 0.007  | 0.011  | 1.000  | -0.008 | 0.020  | -0.318 | -5   |
| X6 | 0.009  | 0.006  | 0.005  | 0.006  | -0.008 | 1.000  | -0.046 | -0.310 | -5   |
| X7 | -0.006 | -0.018 | -0.024 | -0.032 | 0.020  | -0.046 | 1.000  | -0.133 | -2   |

As one can see, Table 2 mirrors Table 1. For this dataset, however, the correlation structure is more complicated.  $X_1$  and  $X_2$  are highly correlated with  $\rho = .7$ .

**Figure 3:**



**Figure 4:**



As seen in Figure 4, the pattern observed between  $X_1$  and  $X_2$  does not carry over to the other correlated predictors.

**Figure 5:**

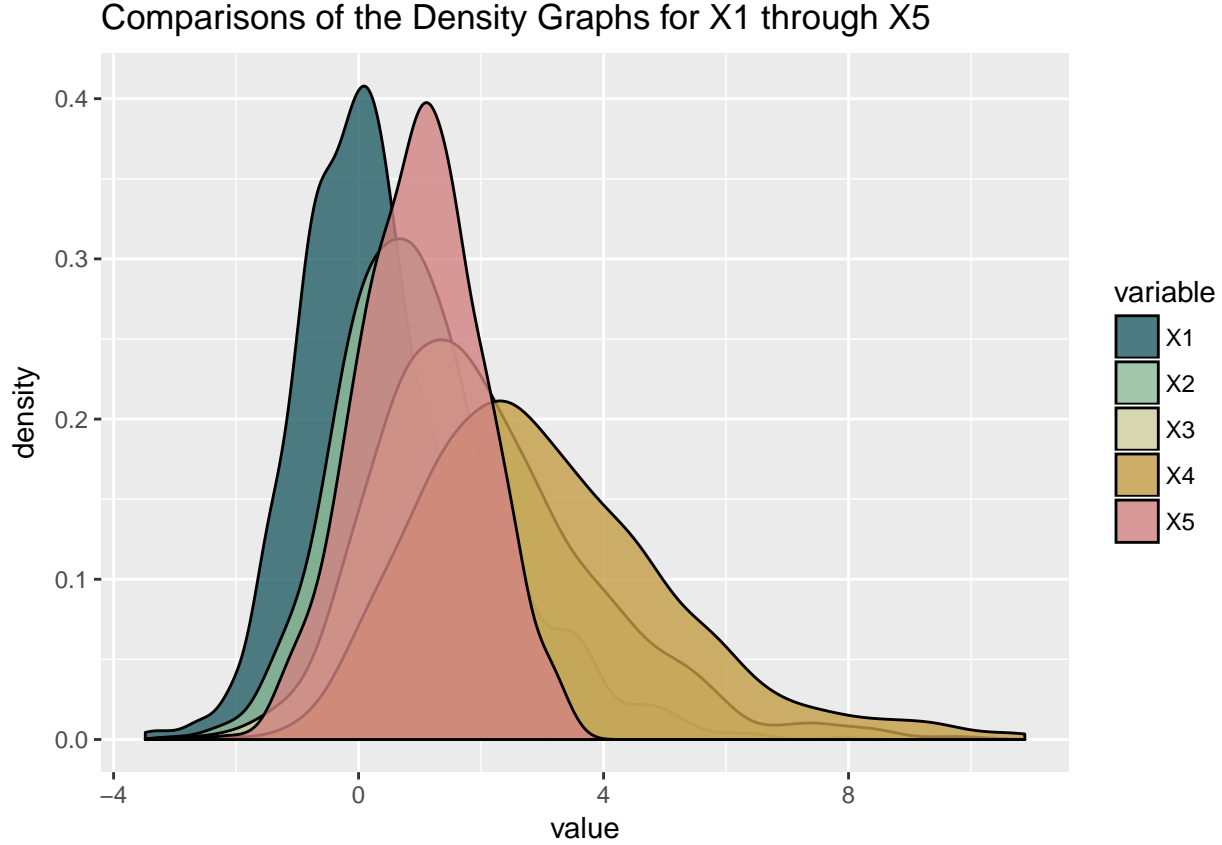


Figure 5 demonstrate how the correlation between a few of the predictors and  $Y$  may be effected by slope. Scale is much more a factor in this dataset, with some variables like  $X_3$  having a larger range than the variables  $X_1 \sim N(0, 1)$  or  $X_5, \dots, X_{12} \sim MVN()$ .

The last dataset we'll consider is  $D_3$ , a data set with even more non-linear relationships between the first four variables. Otherwise it is very similar to both  $D_1$  and  $D_2$ . The first four variables are generated as follows:

$$\begin{aligned}\omega_1 &\sim N(1, 0) \\ \omega_2 &= \log(\omega_1) + E, E \sim N(1, 0) \\ \omega_3 &= \log(\omega_2) + E, E \sim N(1, 0) \\ \omega_4 &= \log(\omega_4) + E, E \sim N(1, 0)\end{aligned}$$

**Table 3: Correlation of  $\omega_1, \dots, \omega_7$  and  $Y$**

|    | W1     | W2     | W3     | W4     | W5     | W6     | W7     | y      | beta |
|----|--------|--------|--------|--------|--------|--------|--------|--------|------|
| W1 | 1.000  | -0.056 | -0.040 | 0.041  | 0.002  | -0.034 | -0.028 | 0.322  | 5    |
| W2 | -0.056 | 1.000  | -0.533 | -0.279 | -0.002 | 0.049  | -0.003 | 0.668  | 5    |
| W3 | -0.040 | -0.533 | 1.000  | -0.002 | -0.019 | -0.031 | -0.010 | -0.096 | 2    |
| W4 | 0.041  | -0.279 | -0.002 | 1.000  | -0.007 | -0.008 | -0.079 | -0.223 | 0    |

|    | W1     | W2     | W3     | W4     | W5     | W6     | W7     | y      | beta |
|----|--------|--------|--------|--------|--------|--------|--------|--------|------|
| W5 | 0.002  | -0.002 | -0.019 | -0.007 | 1.000  | -0.012 | -0.019 | -0.382 | -5   |
| W6 | -0.034 | 0.049  | -0.031 | -0.008 | -0.012 | 1.000  | 0.004  | -0.358 | -5   |
| W7 | -0.028 | -0.003 | -0.010 | -0.079 | -0.019 | 0.004  | 1.000  | -0.159 | -2   |

The linear correlation structure in  $D_3$  is not as striking as in  $D_1$ . The two strongest linear relationships are between  $\omega_2$  and  $\omega_3$  with  $\rho = -.534$  and between  $Y$  and  $\omega_2$  with  $\rho = .700$ .

**Figure 6:**

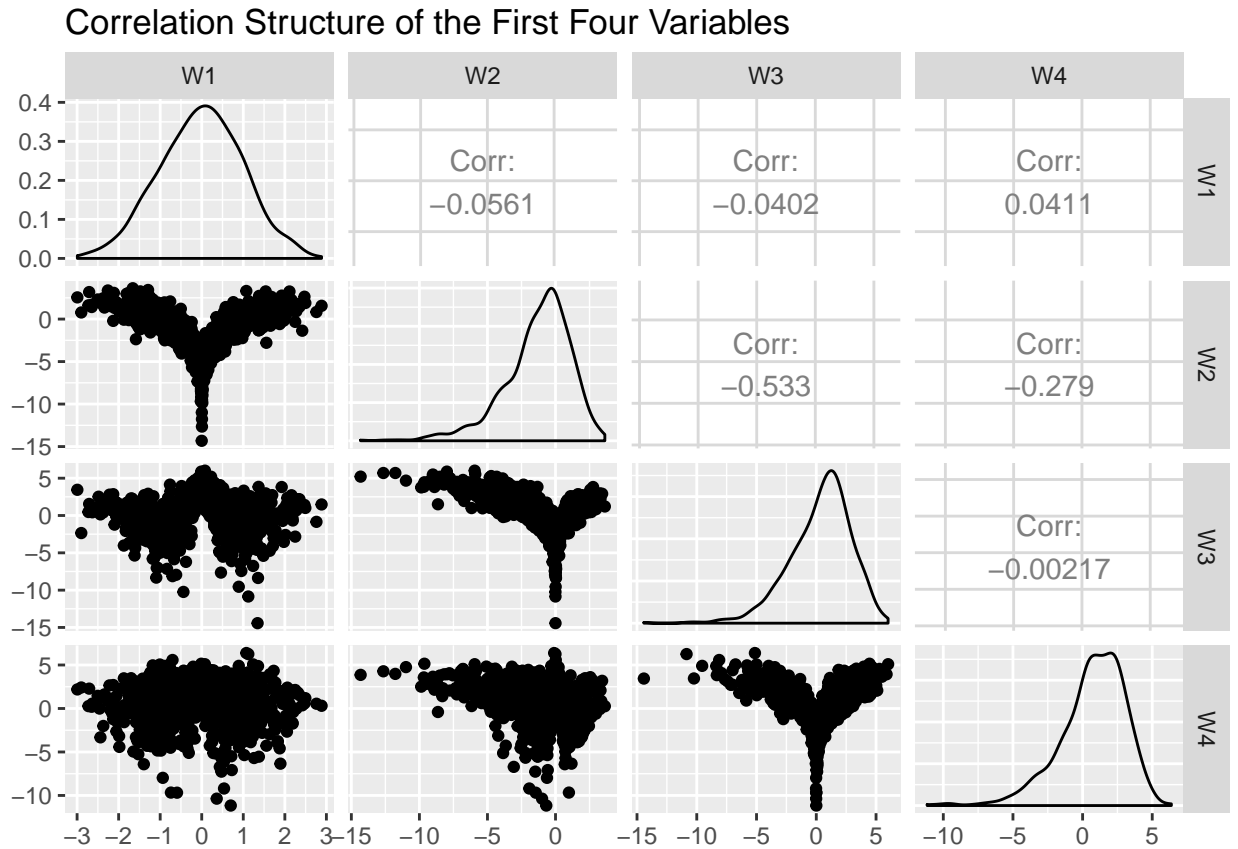
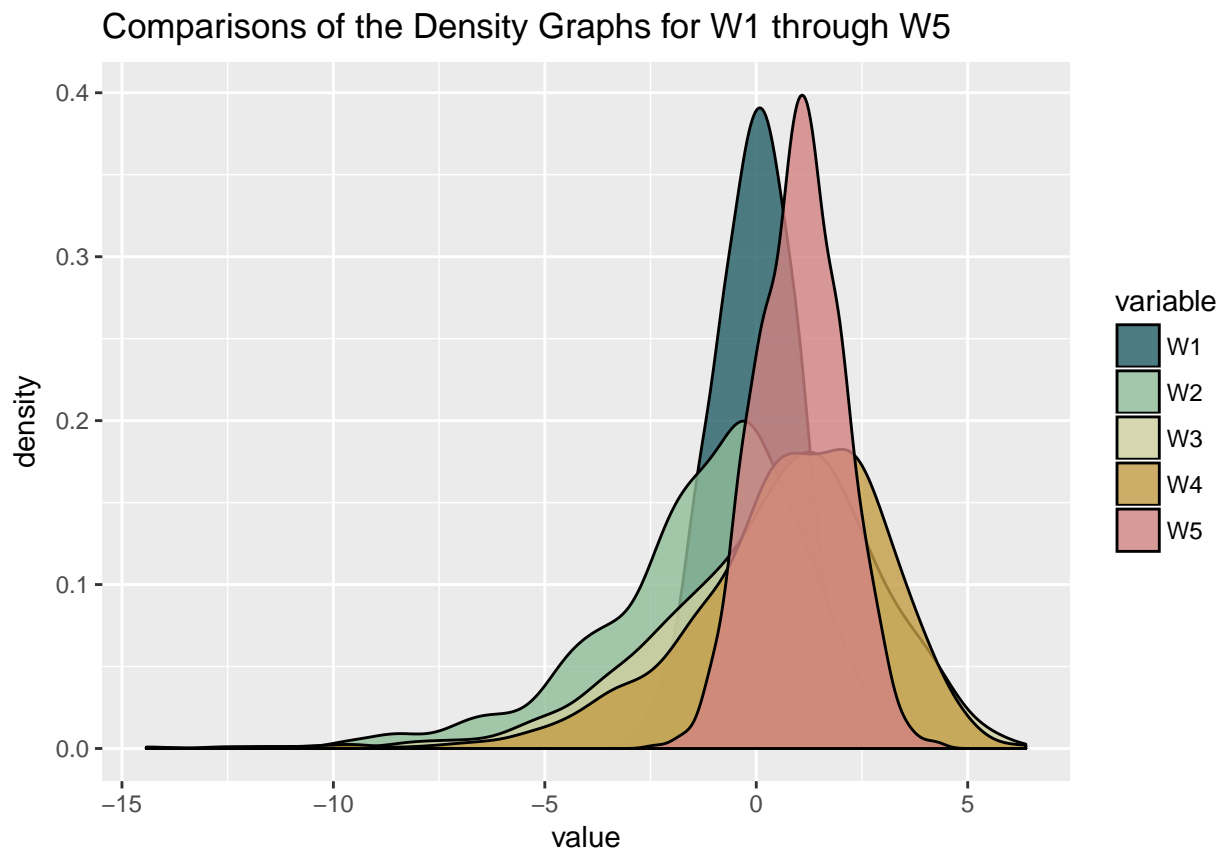


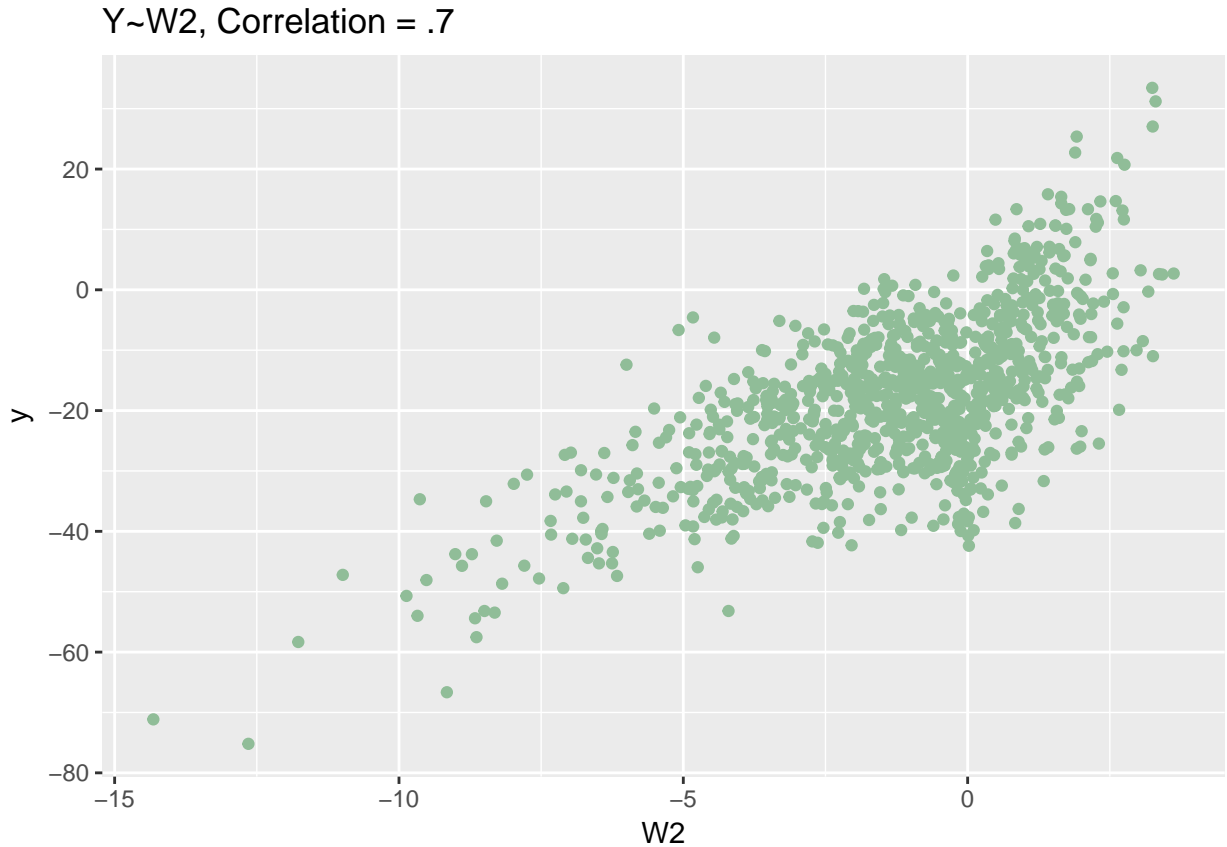
Figure 6 provides another way of visualizing some of the information given in Table 3. Here we can see the densities as well as the paired correlations of the first four variables in  $D_3$ .

**Figure 7:**



There is more variation between the densities of  $\omega_1, \dots, \omega_5$  than we have seen in the other data sets.  $\omega_2, \omega_3$ , and  $\omega_4$  have greater spread than their counterparts that are generated under the normal distribution.

### 2.1.1 Figure 8:



As the relationship between  $Y$  and  $\omega_2$  was so striking, it is nice to see a scatter plot that represents it.

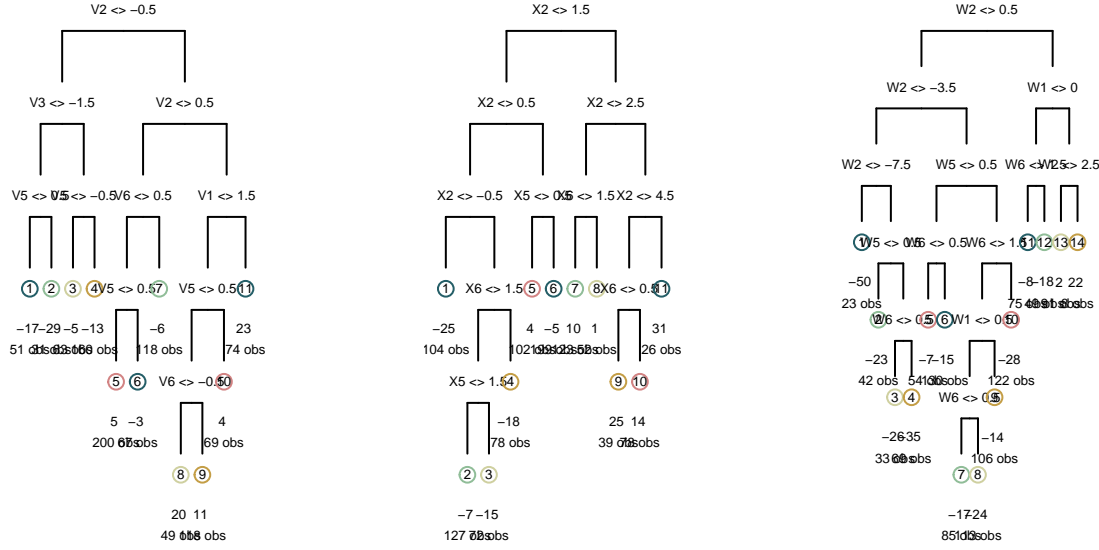
## 2.2 Models and Comparisons

### CART: Regression Trees

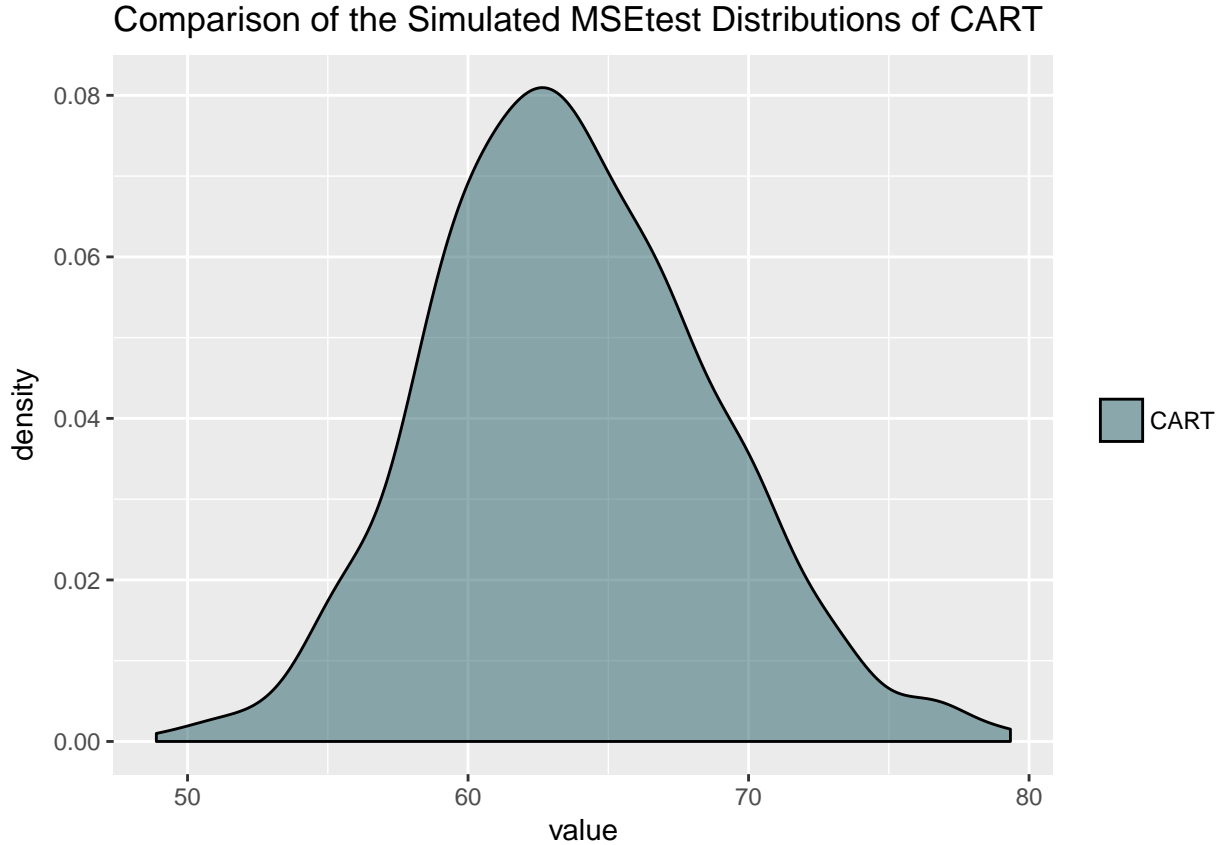
As outlined in the 1984 textbook, *Classification and Regression Trees*, Brieman, Friedman, Olshen, and Stone described their method for creating, pruning, and testing regression trees. There are essentially three steps: one, decide on a variable to split over, two, partition that variable space in two distinct partitions, and three, set our initial predictions for each partition to be mean value of the response according to the observed responses corresponding to the values in the partitions. Recursively, this process is repeated for each new partition until some stopping condition is reached. This is a top down, greedy algorithm that functions by creating as large a tree as possible and then is pruned down to prevent over fitting.

#### 2.2.1 Figure 9:

CART Representing  $Y$ ~, from datasets D1, D2, and D3







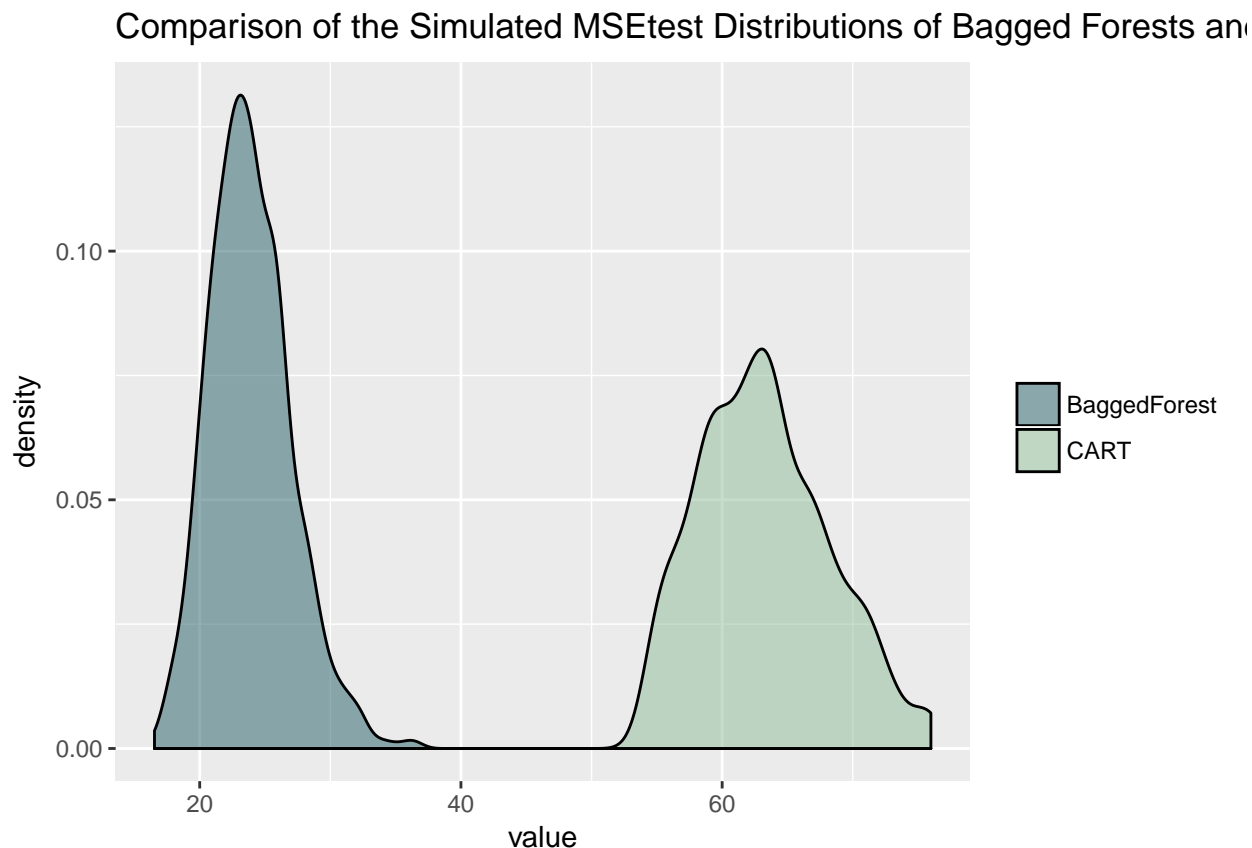
The distribution of 1000 CART trees'  $MSE_{test}$  is roughly normal with a variance of `var(testmseC)`.

## 2.3 Bagged Forests

As one can see in the Figure 10, there is a fair amount of variability in a single tree, they are heavily dependent on fluctuations in the starting data set. As mention briefly in the introduction, bagged forests present one solution to this problem. To create a bagged forest, as outlined in *An Introduction to Statistical Learning* by James, Witten, Hastie and Tibshirani, 2013, many bootstrapped samples are taken from the initial dataset and trees are fitted to them. The final predictions are, then, averaged over all of the trees. This ensures that while each tree has high variance, when they are aggregated the variance will decrease.

Let's put that to the test here using our dataset  $D3$  again. We'll build 100 forests of 100 trees each and compare the variability of the  $MSE$  distributions.

### 2.3.1 Figure 11:



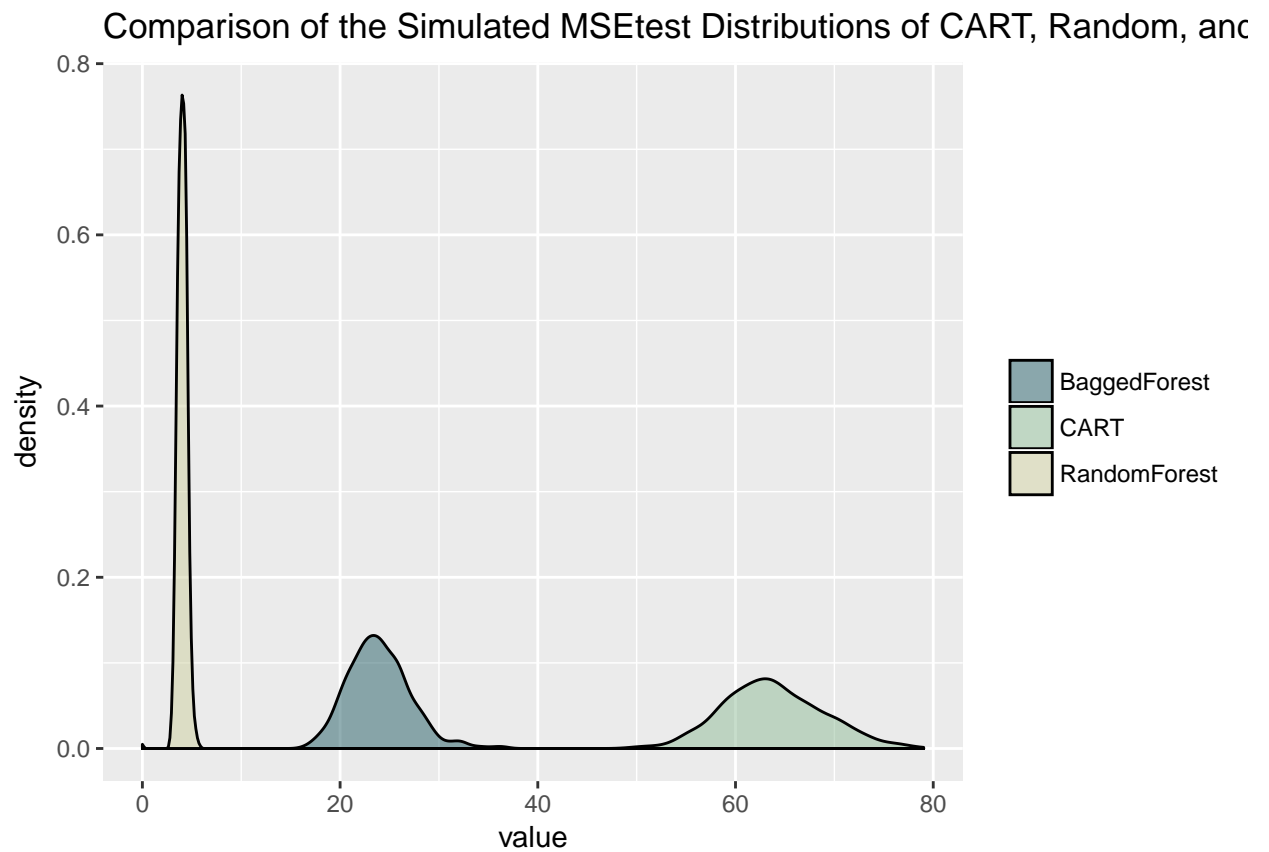
As one can see, the values of  $MSE_{test}$  for the bagged forest were entirely below the  $MSE_{test}$  for the trees and the variance was much smaller.

## 2.4 Random Forests

As the number of trees grown in each forest increases, the  $MSE_{test}$  decreases (cite). Still, this can become computationally intensive on larger data sets where we would like very accurate models. Random forests are often seen as a solution to this problem. In a bagged forest, every variable is considered when each split is made but in a random forest only  $mtry, mtry \leq p$  are considered. This allows us to assume that the trees have a level of independence not found in bagged forests, and that a small random forest will often out perform the bagged forest.

For an illustration, let's build a random forest on  $D3$  and compare the  $MSE$ .

### 2.4.1 Figure 12:



### 2.4.2 TO DO

- Fix the algorithms to work for regression
- Make the OOB references more clear
- Expand around viz created today 3/8



# Chapter 3

## Random Forest Variable Importance

### 3.1 Breiman et al Introduce Permuted Variable Importance (1984)

#### 3.1.1 Variable Importance on a Single Tree

Breiman et al in *Classification and Regression Trees* (1984) propose a method for variable importance for individual trees that stems from their definition of  $\tilde{s}$ , a surrogate split. Surrogate splits help Breiman et al deal with several common problems one may have: modeling with missing data, diagnosing masking, and variable importance. They are defined using logic that resembles that behind random forests.

##### Definitions

Assume the standard structure for tree models. Let  $D$  be the dataset composed of  $D = Y, X_1, \dots, X_p$ , where the model we would like to estimate is of the form  $T : Y \sim X_1, \dots, X_p$ . For any node  $t \in T(D)$ ,  $s^*$  is the best split of the node into daughters  $t_r$  and  $t_l$ . Take  $X_i \in D$  and let  $S_i$  be the set of all of the splits on  $X_i$  in  $T$ . Then set  $\bar{S}_i$  equal to the complement of  $S_i$ ,  $\bar{S}_i = S_i^c$ . For any possible split  $s_i \in S_i \cup \bar{S}_i$ ,  $s_i$  will split the node  $t$  into two daughters,  $t_{i,l}$  and  $t_{i,r}$ . Count the number of times that  $s^*$  and  $s_i$ , while splitting differently, generate the same left daughter  $t_l$  as  $N_{LL}$  and the number of times they generate the same right daughter as  $N_{RR}$ . Then the probability that a case falls within  $t_L \cap t'_L$  is  $P(t_L \cap t'_L) = \sum_j \frac{\pi(j)N_j(LL)}{N_j}$  and the probability that a case falls within  $t_R \cap t'_R$  is  $P(t_R \cap t'_R) = \sum_j \frac{\pi(j)N_j(RR)}{N_j}$ . Where  $\pi(j)$  is the prior assumed for the  $j$ th variable. Finally, the probability that a surrogate split predicts  $s^*$  is  $P(s^*, s_M) = (t_R \cap t'_R) + P(t_L \cap t'_L)$ . Then the surrogate split is the value of  $s^*$  that maximizes this probability. It is denoted  $\tilde{s}$ .

A surrogate split  $\tilde{s}$ , is one that estimates the best possible univariate split  $s^*$  on node  $t$ .

**Defintion: Variable Importance, Single Tree**

$$VI_{tree}(X_i, T) = \sum_{t \in T} \Delta RSS(\tilde{s}_i, t)$$

Or the decrease of RSS attributable to  $X_i$  across the tree  $T$ . In *Classification and Regression Trees*, Breiman et al, outline several potential problems with this method that they do not attempt to solve. First, that this is only one of a number of reasonable ways to define variable importance. Second, the variable importances for variables  $X_1, \dots, X_p$  can be effected by outliers or random fluctuations within the data. (Ch 5.3)

### 3.1.2 Variable Importance for a Random Forest

One way to define variable importance for a random forest follows directly from Breiman et al's definition for a single tree. Recall that each tree in a random forest is fit to a bootstrapped sample of the original observations. To estimate the test error, therefore, no cross validation is needed - each tree is simply tested against the test set of observations that were not in that tree's initial training set. To determine variable importance for a predictor  $X_j$ , we look at the RSS of the each tree's prediction that did not split on  $X_j$ . These values are then averaged over the subset forest that did not include  $X_j$ . A large value would imply that in trees that included  $X_j$ , the predictive capabilities were increased.

To formalize that idea, let's refer to the set of trees that did not consider  $X_j$ ,  $T_{x_j}^c$ . Now,  $T_{x_j}^c \subset R$ , the random forest. The subset of the original data that will be tested on each tree,  $t$ , is  $\bar{B}^t$ . The dimensions of  $\bar{B}^t$  are  $\nu_t \times p$ , where  $p$  is the number of predictors and  $\nu \leq n$ . The number of trees in  $T_{x_j}^c$  is  $\mu$  where  $\mu \leq ntree$

Now, base variable importance is:

$$VI_\alpha(X_j, R) = \sum_{t \in T_{x_j}^c} \frac{1}{\nu_t} RSS(t, \bar{B}_t)$$

However, this method poses some problems. Namely, while variable importance for random forests is more stable than for the variable importance values for CART, (this is because the model is less variable in general), it is lacking the traditional inferential capabilities of other regression models. In an effort to derive a p-value for variable importance values, Breiman 2001b, describes a *permuted variable importance* or  $VI_\beta$  that does not utilize  $T_{x_j}^c$ .

## 3.2 Strobl et al Respond (2008)

Strobl et al (2008) respond to Breiman's method with one main argument: the null hypothesis implied by the permutation distribution utilized in permuted variable importance is that  $X_i$  is independent of  $Y$  **and**  $X_j \notin X_1, \dots, X_p$  so the null hypothesis will be rejected in the case where  $X_j$  is independent of  $Y$  but not some subset of the other predictors. As correlation among the predictors is very common in data sets that are used for random forests, this is a large problem for Breiman's method.

**Algorithm 3** Permuted Variable Importance for Random Forests,  $VI_\beta$ 

- 
- 1: Fit a random forest,  $R$  on the dataset  $D$  fitting the model  $Y \sim X_1, \dots, X_p$ .
  - 2: **for** each  $X_i \in X_1, \dots, X_p$  **do**
  - 3:   **for** each  $t \in R$  **do**
  - 4:     Calculate:  $\Phi_o = \frac{1}{\nu_t} RSS(t, \bar{B}^t)$
  - 5:     Permute  $X_i$ . Now find  $\Phi^* = \frac{1}{\nu_t} RSS(t, \bar{B}_t^*)$
  - 6:     The difference between these values,  $\Phi^* - \Phi_o$ , is the variable importance for  $X_j$  on  $t$ ,
  - 7:   **end for**
  - 8:   Average over all  $t \in R$
- 

$$VI_\beta(X_j) = \frac{1}{ntree} \sum \Phi^* - \Phi_o$$

$$VI_\beta(X_j) = \frac{1}{ntree} \sum \frac{1}{\nu_t} RSS(t, \bar{B}_t^*) - \frac{1}{\nu_t} RSS(t, \bar{B}^t)$$

9: **end for**

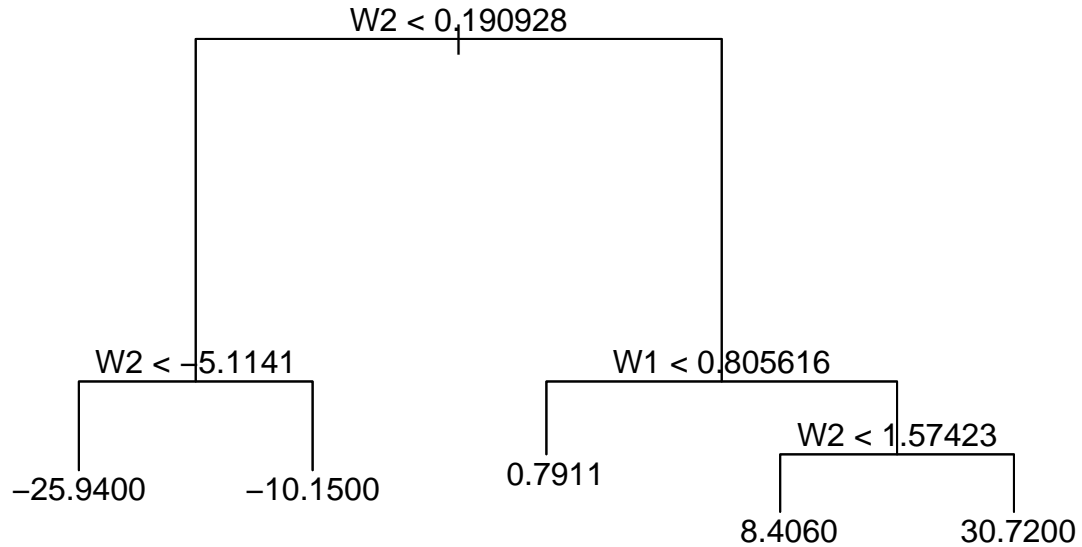
---

To alleviate this difficulty, Strobl et al propose a permutation scheme under the null hypothesis that  $X_j$  given it's relationship with the other predictors is independent of  $Y$ .

Finally, as above, the permuted variable importance for each variable  $X_i \in X_1, \dots, X_p$  is the sum of  $VI_{perm2}(X_i, T_i)$  over all  $T_i \in RF$ .

AN EXAMPLE:

Fit a simple tree:



**Algorithm 4** Conditional Variable Importance for Random Forests

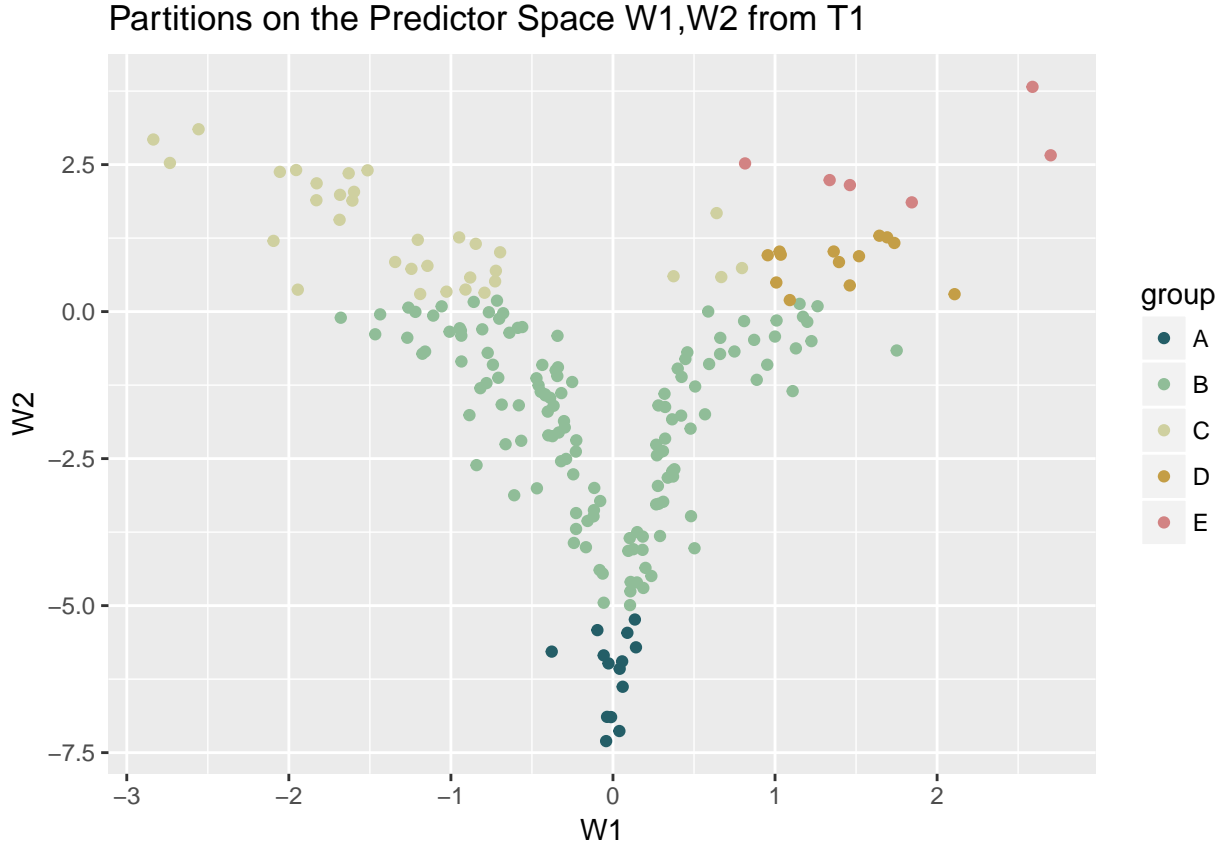
- 1: Fit a random forest,  $RF$  on the dataset  $D$  fitting the model  $Y \sim X_1, \dots, X_p$
- 2: **for** each  $T_i \in RF$  **do**
- 3:     Calculate the oob error before permutation:

$$\frac{\sum_{i \in \bar{B}^T} I(y = \hat{y})}{|\bar{B}^T|}$$

- 4:     **for** each  $X_i \in X_1, \dots, X_p$  **do**
- 5:         Permute the dataset,  $D$ , conditional on the grid created by the splits on  $X_i$  in  $T_i$
- 6:         Calculate the oob error again and compare with the estimate before permutation

$$VI_{perm2}(X_i, T_i) = \frac{\sum_{i \in \bar{B}^T} I(y = \hat{y})}{|\bar{B}^T|} - \frac{\sum_{i \in \bar{B}^T} I(y = \hat{y}^p)}{|\bar{B}^T|}$$

- 7:     **end for**
- 8: **end for**

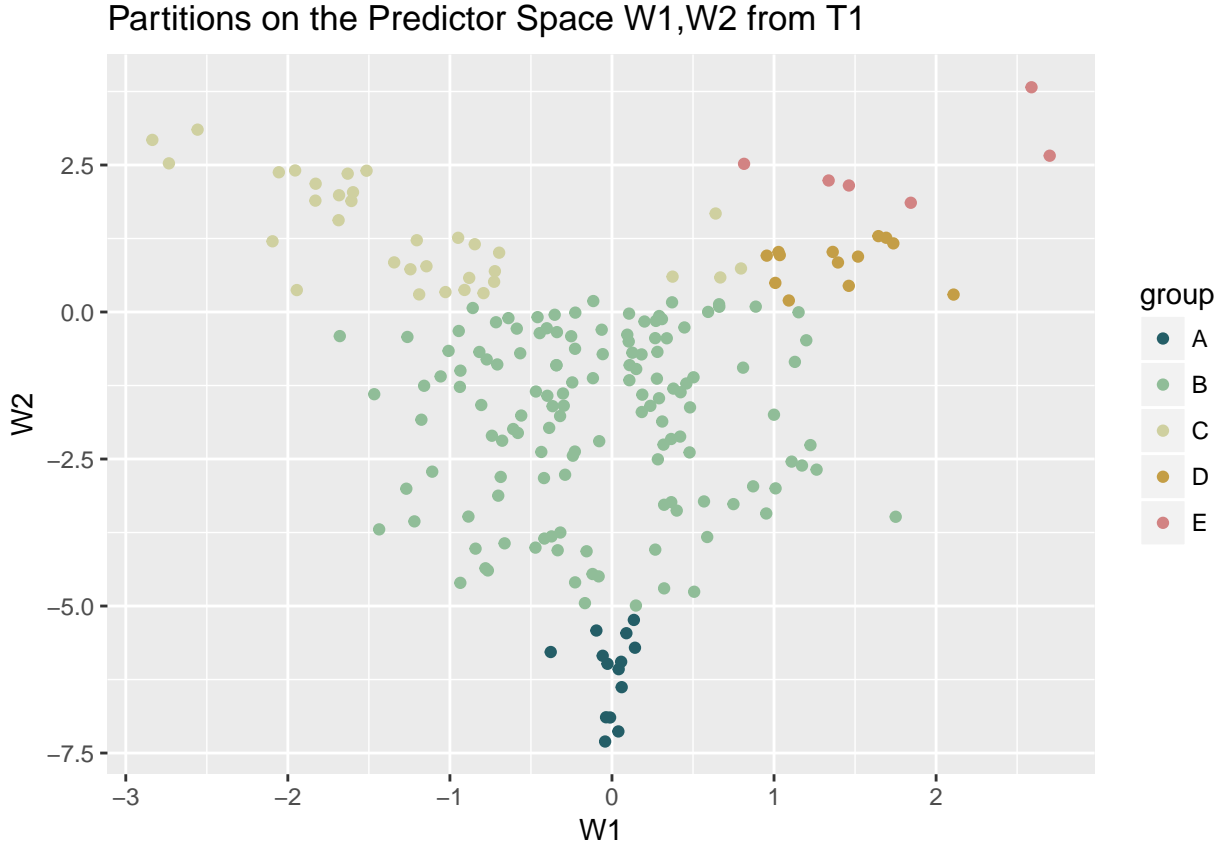


| Group | Predicted Value of Y | Min(W2) | Max(W2) | Min(W1) | Max(W1) |
|-------|----------------------|---------|---------|---------|---------|
| A     | -25.9369010231033    | -7.3    | -5.24   | -0.38   | 0.14    |
| B     | -10.1521754868715    | -4.99   | 0.19    | -1.68   | 1.75    |



| Group | Predicted Value of Y | Min(W2) | Max(W2) | Min(W1) | Max(W1) |
|-------|----------------------|---------|---------|---------|---------|
| C     | 0.791053130734821    | 0.3     | 3.1     | -2.84   | 0.8     |
| D     | 8.40580623141462     | 0.2     | 1.29    | 0.96    | 2.11    |
| E     | 30.7209296380884     | 1.86    | 3.82    | 0.82    | 2.7     |

If we permute the  $\omega_1$  values in group  $B$ , this is what that plot looks like:



### 3.3 Inferential Variable Importance

This thesis hopes to be a reponse to conditional variable importance as outlined by Strobl et al 2008. First is that the practice of permuting given the partitions from the model  $Y \sim X_1, \dots, X_p$  instead of  $X_j \sim X_1, \dots, X_p$



## Chapter 4

# INFTrees and INFforests Variable Importance

### 4.1 Theory

### 4.2 Implementation



# Conclusion

If we don't want Conclusion to have a chapter number next to it, we can add the `{.unnumbered}` attribute. This has an unintended consequence of the sections being labeled as 3.6 for example though instead of 4.1. The  $\text{\LaTeX}$  commands immediately following the Conclusion declaration get things back on track.

## More info

And here's some other random info: the first paragraph after a chapter title or section head *shouldn't be* indented, because indents are to tell the reader that you're starting a new paragraph. Since that's obvious after a chapter or section title, proper typesetting doesn't add an indent there.



# Appendix A

## The First Appendix

This first appendix includes all of the R chunks of code that were hidden throughout the document (using the `include = FALSE` chunk tag) to help with readability and/or setup.

In the main Rmd file:

```
# This chunk ensures that the reedtemplates package is  
# installed and loaded. This reedtemplates package includes  
# the template files for the thesis and also two functions  
# used for labeling and referencing  
if(!require(devtools))  
  install.packages("devtools", repos = "http://cran.rstudio.com")  
if(!require(reedtemplates)){  
  library(devtools)  
  devtools::install_github("ismayc/reedtemplates")  
}  
library(reedtemplates)
```

In :

```
# This chunk ensures that the reedtemplates package is  
# installed and loaded. This reedtemplates package includes  
# the template files for the thesis and also two functions  
# used for labeling and referencing  
if(!require(devtools))  
  install.packages("devtools", repos = "http://cran.rstudio.com")  
if(!require(dplyr))  
  install.packages("dplyr", repos = "http://cran.rstudio.com")  
if(!require(ggplot2))  
  install.packages("ggplot2", repos = "http://cran.rstudio.com")
```

```
if(!require(reedtemplates)){  
  library(devtools)  
  devtools::install_github("ismayc/reedtemplates")  
}  
library(reedtemplates)  
library(tree)  
library(reedtemplates)  
library(MASS)  
library(ggplot2)  
library(randomForest)  
library(maptree)  
library(knitr)  
library(GGally)  
library(reshape)  
#flights <- read.csv("data/flights.csv")  
thesis <- c("#245e67", "#90bd98", "#cfd0a0", "#c49e46", "#d28383")
```



# Appendix B

## The Second Appendix: CTree

### B.0.1 Conditional Inference Trees

As mentioned in the introduction, CART has the tendency to bias towards variables with the most possible splits and overfitting. There is little head paid to statistical significance or general statistical theory. *Conditional Inference Trees* are a method proposed by Hothorn et al, 2006, that utilizes permutation theory to create and algorithm that is sensitive to these issues. A crutial difference between CTree and CART is that while CART is a top down algorithm, CTree initially assumes each row of the dataset is a node and then gradually prunes them.

---

**Algorithm 5** Conditional Inference Trees

---

```
1: for  $w_i, i \in \{w_1, \dots, w_n\}$  do
2:   Test the global null hypothesis of independence between any of the  $m$  covariates
   and the response.
3:   if  $H_O$  cannot be rejected then
4:     Stop
5:   else
6:     Select predictor  $X_j$  with the strongest linear association to  $Y$ 
7:   end if
8:   Choose a set  $A \in X_j$  such that  $A \cup X_j \setminus A = X_j$ 
9:   The case weights,  $w_{left}$  and  $w_{right}$  are then defined as  $w_{left,i} = w_i I(x_j \in A)$ 
   and  $w_{right,i} = w_i I(x_j \notin A)$ 
10: end for
```

---

The case weights,  $w_i \in w_1, \dots, w_n$ , correspond to nodes are defined as:

$$w_i = I(x_i \in N_t)$$

Where  $x_i$  is a vector of observations and  $N_t$  is a node in the tree.

**CTree fitted to  $D_3$**

[illegible]

# References

- Angel, E. (2000). *Interactive computer graphics : A top-down approach with opengl*. Boston, MA: Addison Wesley Longman.
- Angel, E. (2001a). *Batch-file computer graphics : A bottom-up approach with quicktime*. Boston, MA: Wesley Addison Longman.
- Angel, E. (2001b). *Test second book by angel*. Boston, MA: Wesley Addison Longman.