

Team Checkpoint 1

Forest Cover Team B

Annie Condon, Matt Hayden, Matt Robertson, Yvette Gonzalez

Contents

| | | |
|----------|----------------------------------|-----------|
| 1 | Introduction | 3 |
| 2 | The Modeling Problem | 3 |
| 3 | The Data | 3 |
| 3.1 | Response Variable | 4 |
| 3.2 | Continuous Variables | 4 |
| 3.3 | Binary Variables | 5 |
| 4 | Exploratory Data Analysis | 7 |
| 5 | Conclusion | 14 |

1 Introduction

The U.S. Forest Service relies on an accurate understanding of its forests composition in order to best protect and manage the forest land. Conducting accurate inventory of forest composition by direct observation or remotely sensed data is often too expensive and time consuming to do at large-scale. Predictive analytics can be employed to use the results of a small-scale survey to create a model that can be applied across a large region, using descriptive features extracted from maps of the area.

In this paper, our objective is to predict the forest cover type given a set of cartographic features and a variety of multiclass classification models. Our models will be evaluated using predictive accuracy.

2 The Modeling Problem

Our modeling problem is to predict the forest cover type as a multiclass classification problem based on the associated features. A multiclass classification problem classifies instances into one of the more than two classes. Our forest cover type is defined as one of seven, mutually exclusive, forest cover type classes, shown in table 1 below.

Table 1: Cover Type Classes

| Cover Type | Number of Observations |
|-------------------|------------------------|
| Spruce/Fir | 211840 |
| Lodgepole Pine | 283301 |
| Ponderosa Pine | 35754 |
| Cottonwood/Willow | 2747 |
| Aspen | 9493 |
| Douglas/Fir | 17367 |
| Krummholz | 20510 |

Several algorithms have been developed to solve multiclass classification problems. The possible algorithms we will consider for our problem are: neural networks, k-nearest neighbors, random forest and support vector machines.

We will evaluate our models using classification accuracy.

3 The Data

Our data set consists of 581,012 observations of the 30 x 30 meter cells of forest and 54 features associated with each cell. The features are derived from 12 attributes, with area and soil type binarized so that there are 4 binary area designators and 40 binary soil type designators. There is no missing data. The feature descriptions are listed in the table below:

Table 2: Features

| Feature | Descriptions |
|-----------------|--|
| Elevation | Elevation in meters |
| Aspect | Aspect in degrees azimuth |
| Slope | Slope in degrees |
| HDist.Hydrology | Horizontal distance to nearest surface water feature in meters |
| VDist.Hydrology | Vertical distance to nearest surface water feature in meters |

| Feature | Descriptions |
|------------------|---|
| HDist.Roadway | Horizontal distance to nearest roadway in meters |
| Hillshade.9am | Hillshade index at 9am, summer solstice |
| Hillshade.12pm | Hillshade index at noon, summer solstice |
| Hillshade.3pm | Hillshade index at 3pm, summer solstice |
| HHDist.FirePoint | Horizontal Distance to nearest wildfire ignition points |
| Area | Wilderness area designation - 4 binary areas |
| SoilType | Soil Type designation - 40 binary values |

3.1 Response Variable

Figure one shows the frequency of each class of cover type in our data set. 85 percent of the observations fall into classes 1 and 2, Spruce/Fir and Lodgepole Pine. Class 4, Cottonwood/Willow has the least amount of observations at 2,747.

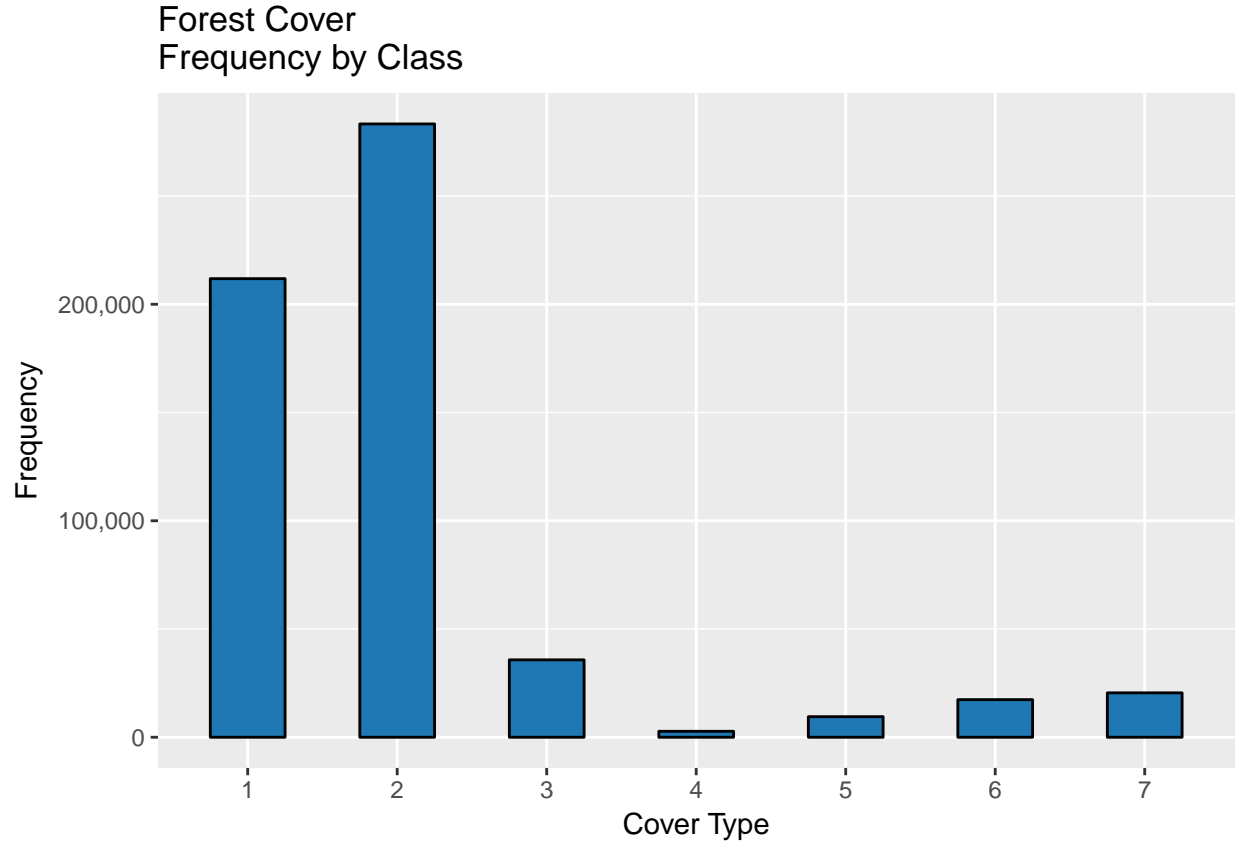


Figure 1: Frequency of each Cover Type Class

3.2 Continuous Variables

Table 3 includes some standard statistical measures of central tendency and variation for our continuous variables.

Table 3: Summary Statistics for Continuous Data

| | mean | sd | median | min | max | type |
|-----------------|-----------|-----------|--------|------|------|------------|
| Elevation | 2959.3653 | 279.9847 | 2996 | 1859 | 3858 | Continuous |
| Aspect | 155.6568 | 111.9137 | 127 | 0 | 360 | Continuous |
| Slope | 14.1037 | 7.4882 | 13 | 0 | 66 | Continuous |
| HDist.Hydrology | 269.4282 | 212.5494 | 218 | 0 | 1397 | Continuous |
| VDist.Hydrology | 46.4189 | 58.2952 | 30 | -173 | 601 | Continuous |
| HDist.Roadway | 2350.1466 | 1559.2549 | 1997 | 0 | 7117 | Continuous |
| Hillshade.9am | 212.146 | 26.7699 | 218 | 0 | 254 | Continuous |
| Hillshade.12pm | 223.3187 | 19.7687 | 226 | 0 | 254 | Continuous |
| Hillshade.3pm | 142.5283 | 38.2745 | 143 | 0 | 254 | Continuous |
| HDist.FirePoint | 1980.2912 | 1324.1952 | 1710 | 0 | 7173 | Continuous |

3.3 Binary Variables

Our data's binary variables consist of 4 different wilderness area designators and 40 different soil type designators. Table 4 and 5 below represent the frequency of each of the areas and soil types, in descending order.

Table 4: Area Type Counts

| Name | Count |
|-------|--------|
| Area1 | 260796 |
| Area3 | 253364 |
| Area4 | 36968 |
| Area2 | 29884 |

Table 5: Soil Type Counts

| Name | Count |
|------------|--------|
| SoilType29 | 115247 |
| SoilType23 | 57752 |
| SoilType32 | 52519 |
| SoilType33 | 45154 |
| SoilType22 | 33373 |
| SoilType10 | 32634 |
| SoilType30 | 30170 |
| SoilType12 | 29971 |
| SoilType31 | 25666 |
| SoilType24 | 21278 |
| SoilType13 | 17431 |
| SoilType38 | 15573 |
| SoilType39 | 13806 |
| SoilType11 | 12410 |
| SoilType4 | 12396 |
| SoilType20 | 9259 |
| SoilType40 | 8750 |
| SoilType2 | 7525 |
| SoilType6 | 6575 |

| Name | Count |
|------------|-------|
| SoilType3 | 4823 |
| SoilType19 | 4021 |
| SoilType17 | 3422 |
| SoilType1 | 3031 |
| SoilType16 | 2845 |
| SoilType26 | 2589 |
| SoilType18 | 1899 |
| SoilType35 | 1891 |
| SoilType34 | 1611 |
| SoilType5 | 1597 |
| SoilType9 | 1147 |
| SoilType27 | 1086 |
| SoilType28 | 946 |
| SoilType21 | 838 |
| SoilType14 | 599 |
| SoilType25 | 474 |
| SoilType37 | 298 |
| SoilType8 | 179 |
| SoilType36 | 119 |
| SoilType7 | 105 |
| SoilType15 | 3 |

4 Exploratory Data Analysis

Our next step is to explore the relationships in our data. We will begin by looking at boxplots of our scaled continuous variables in order to understand their relative distribution. We note that `Hillshade.12pm` and `VDist.Hydrology` have more pronounced skews than the other variables. We will need to investigate transformations if we use modeling techniques that can be negatively effected by outliers.

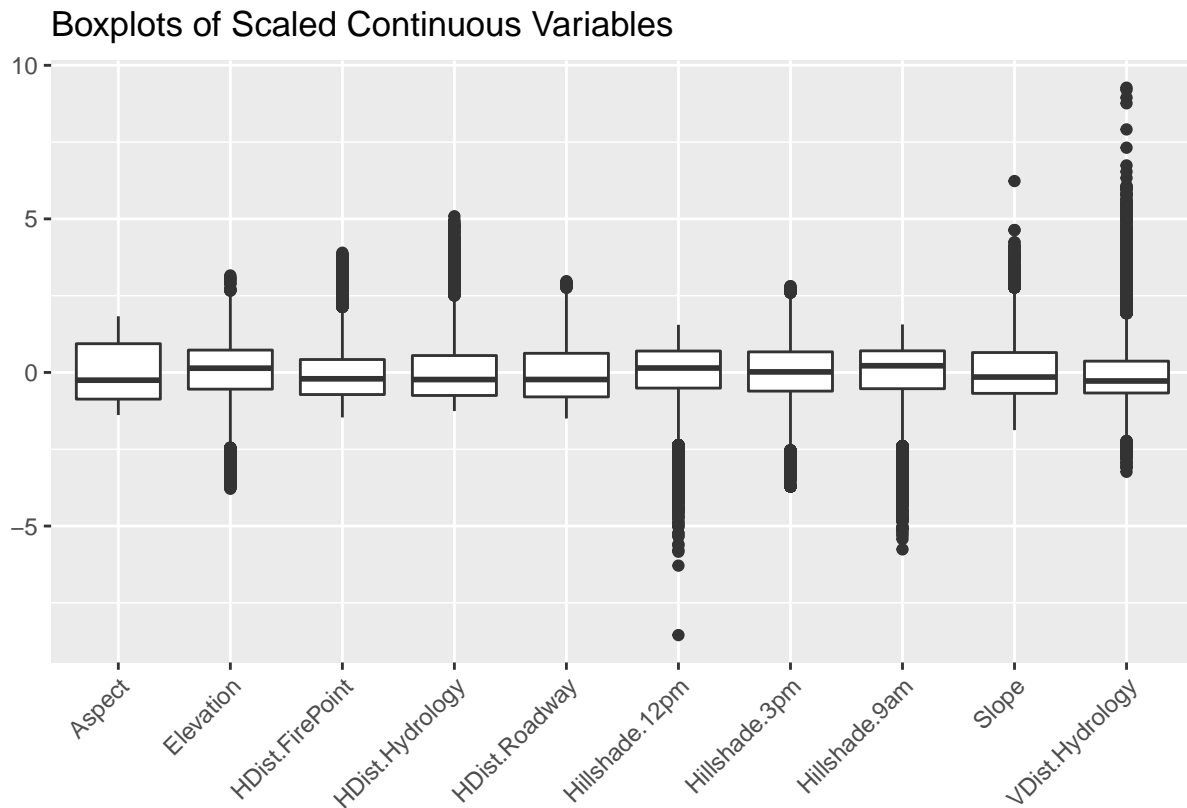


Figure 2: Boxplots of Scaled Continuous Variables

Our data exploration is informed by our statistical problem, which is one of multiclass classification. We will therefore be interested in exploring our predictors by each of our classes. The density plots below superimpose the density estimates for each variable by class, 1-7. We note that our variable Aspect shows multimodal distributions, indicating that it may have more than one grouping included. We will also note that several of our variables show distinct distributions by class, indicating that they would serve as a good predictor.

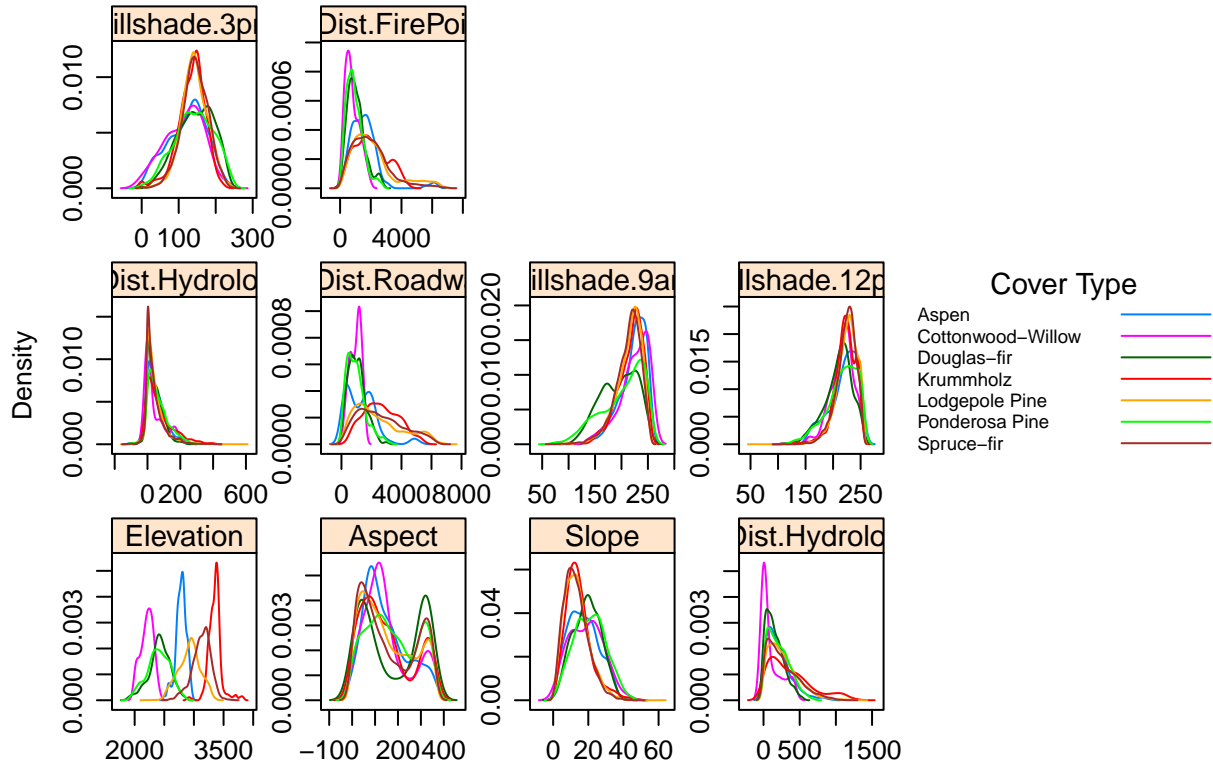


Figure 3: Density Plots of Continuous Variables by Class

Figure 4 shows density plots by each of the four wilderness areas.

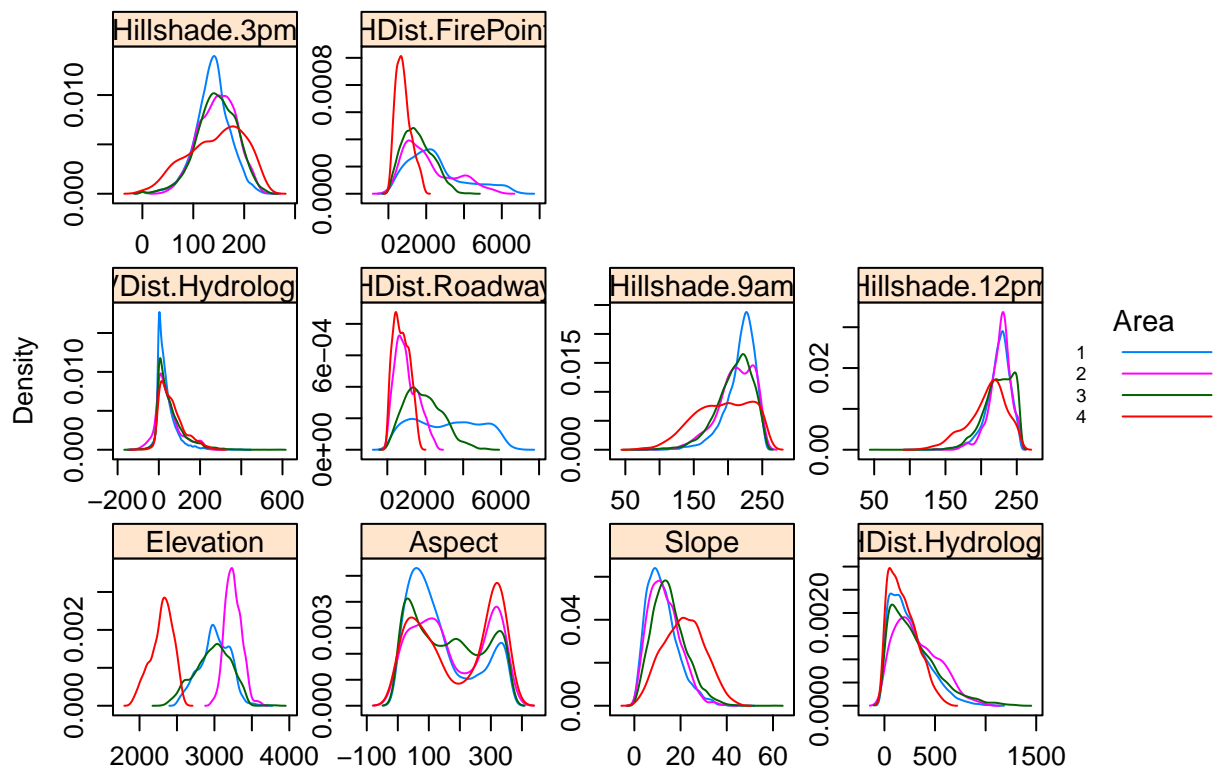


Figure 4: Density Plots of Continuous Variables by Area

Next, we will look at the correlations between the our predictor variables. Highly correlated predictors can have a negative effect on some of our modeling algorithms. Below is a correlation matrix with dark blue colors indicating a strong positive correlation and dark red colors indicating a strong negative correlation. We see the following variables with higher correlations that should be investigated further:

High Positive Collinearity

- VDist.Hydrology and HDist.Hydrology
- Aspect and Hillshade.3pm
- Hillshade.3pm and Hillshade.12pm

High Negative Collinearity

- Aspect and Hillshade.9am
- Slope and Hillshade.12pm
- Hillshade.3pm and Hillshade.9am

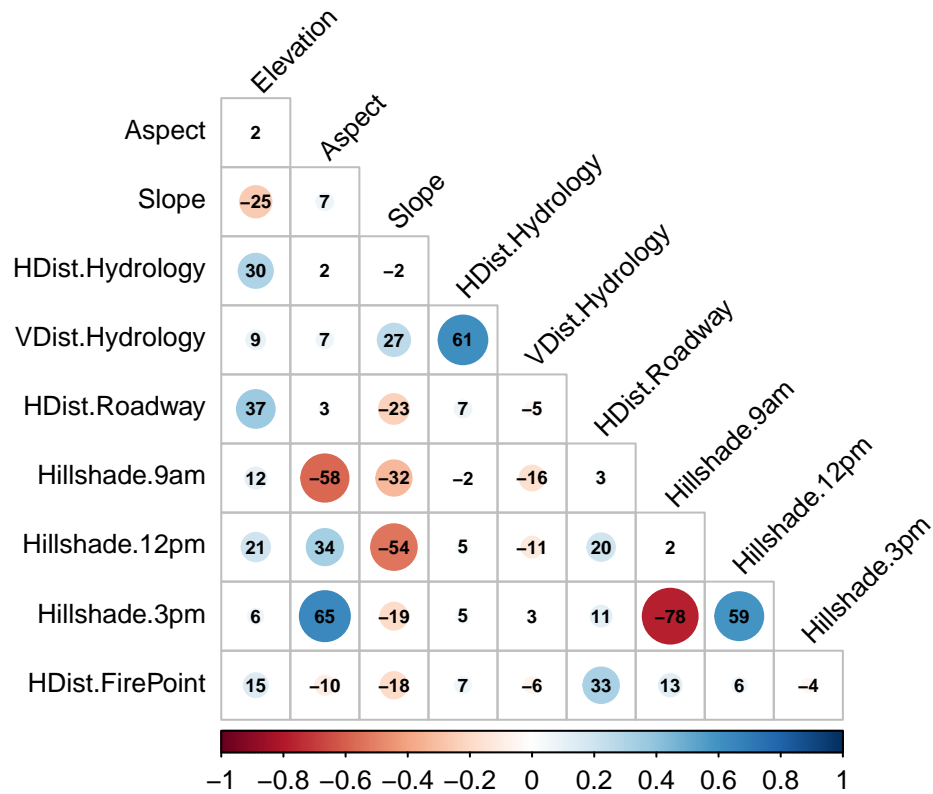


Figure 5: Correlation Matrix

Figure 6 shows a barchart of our 40 different soil types which shows the proportion of frequency of each class. We can see that soil types 1-7 have similar proportions, as do 19-33 and 38-40. A few soil types have only one cover type class, making them especially good predictors.

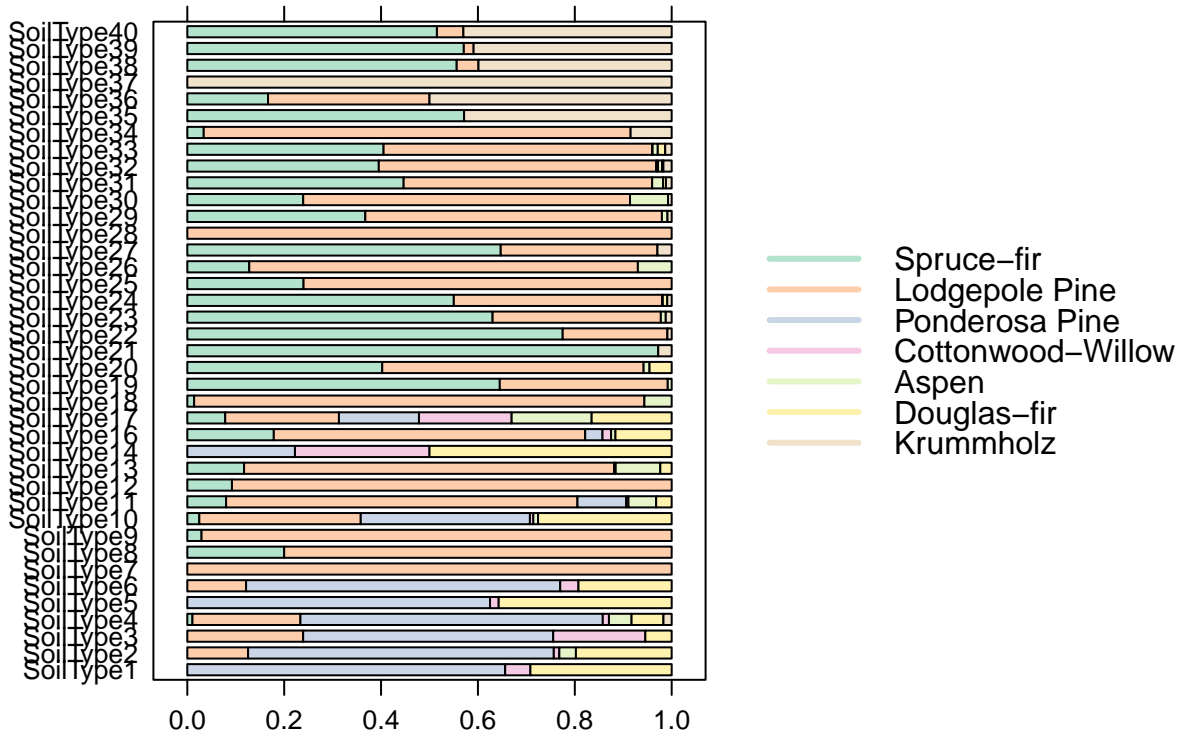


Figure 6: Soil Barchart

Figure 7 shows a barchart of our 4 different wilderness areas with the proportion of frequency of each class. We can see that the class composition in area 4 is distinguishable from the other areas, making it a good predictor.

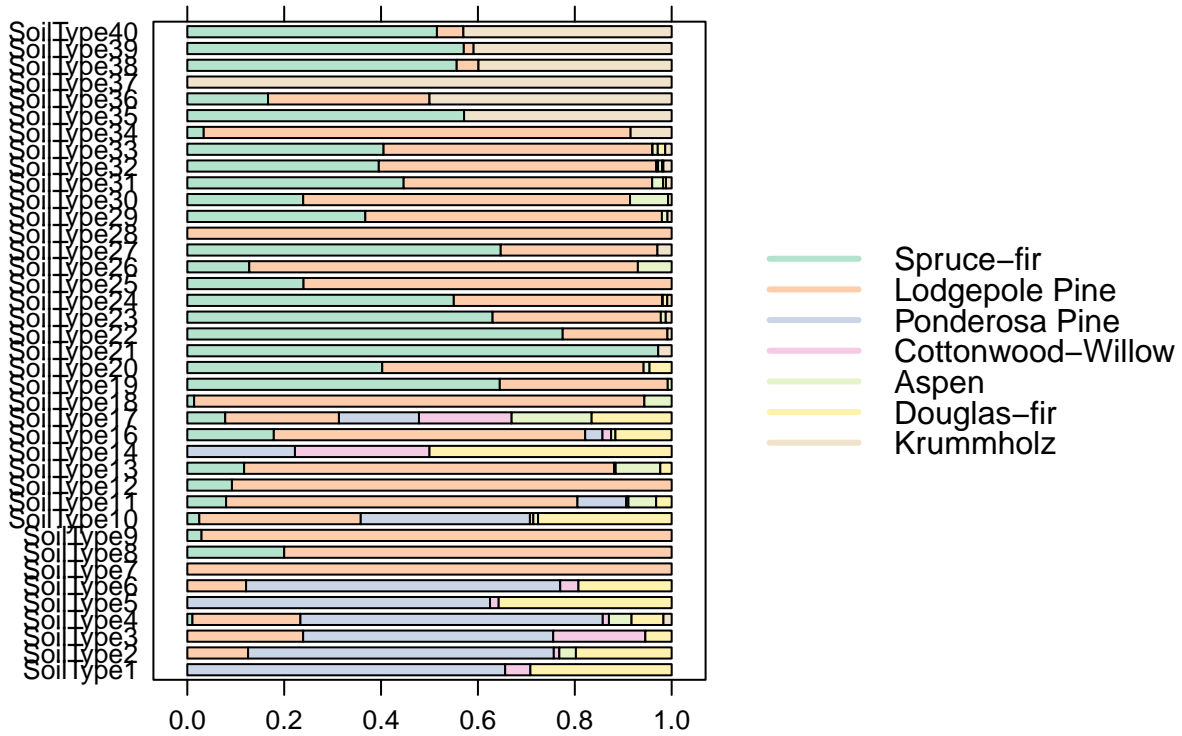


Figure 7: Area Barchart

Next, we look at the horizontal distance to the roadway, by class. We can see that most observations in classes 5-7 are closest to the roadways."

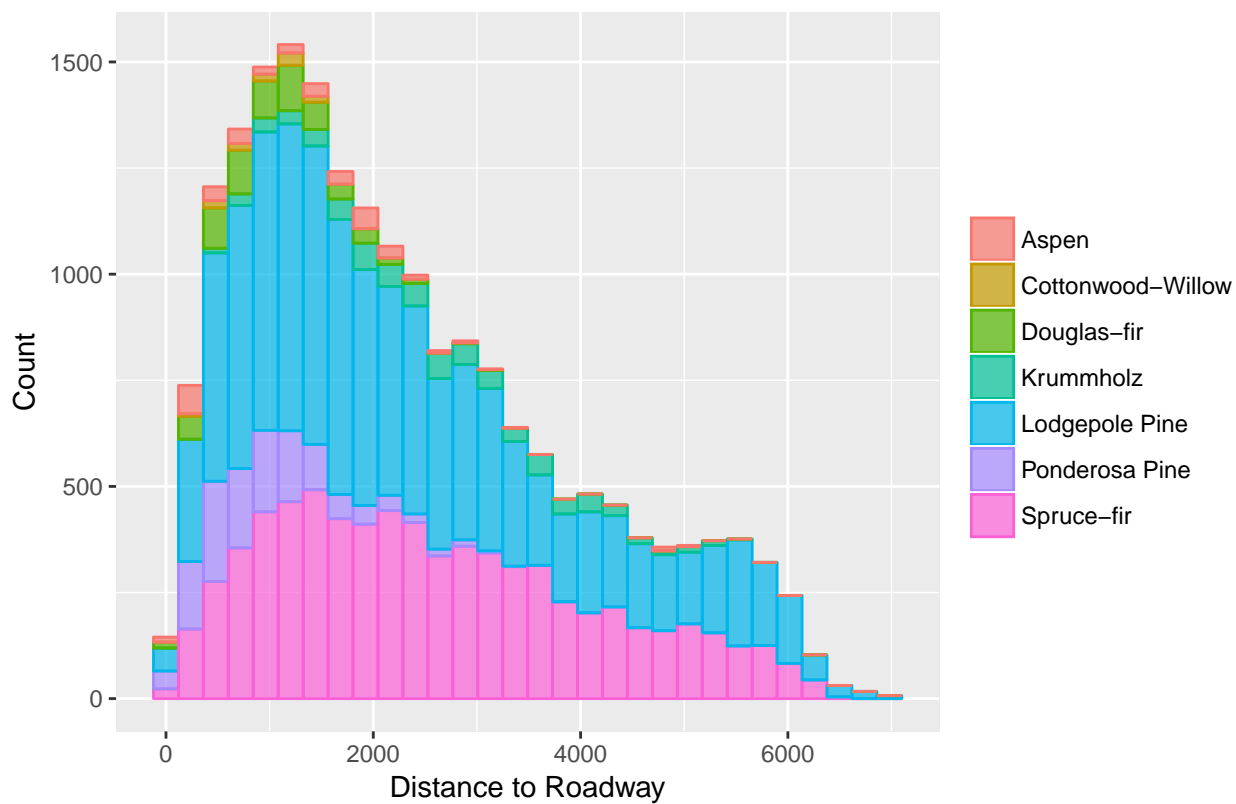
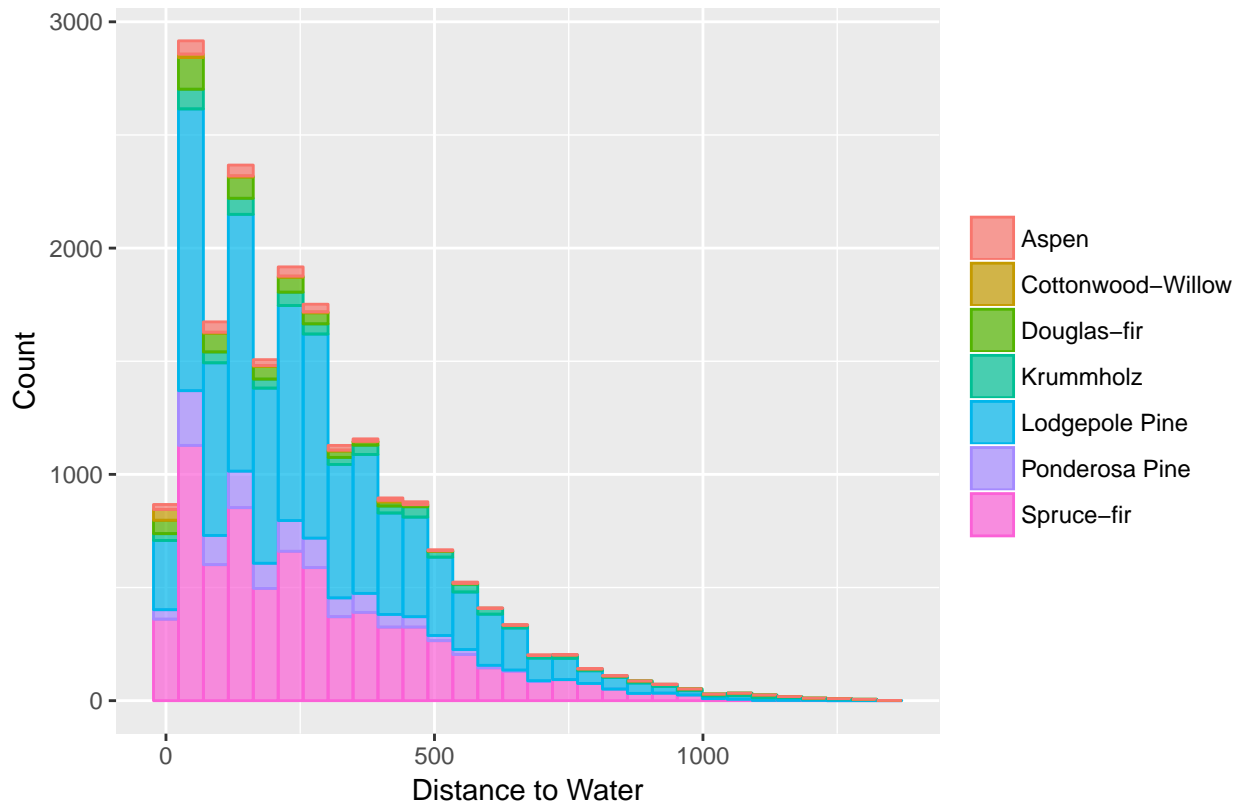


Figure 9 shows horizontal distance to the water, by class. Again, we see that most observations in classes 5-7 are closest to a water source.



5 Conclusion

Our initial analysis for our forest cover type prediction problem included defining our modeling problem, doing a data quality and inventory check and performing a preliminary exploratory data analysis. The results of our data quality check showed that we had no missing data and gave us a cursory understanding of our data set.

Our preliminary exploratory analysis revealed some potentially useful predictors and helped us to understand the relationships in our data. Our next step will be to build a few exploratory models in order to gain further insights into predictors we can use in our modeling.