# Determining Forest Cover Type Using Classification

Francesca Chiappetta

# Intro

## Background

- Growing demand of carbon credits is leading to a burgeoning carbon market.

- One way to generate carbon credits is through sequestering carbon in forest ecosystems.

## Context

- Determining the amount of carbon sequestered by a tree or many tree is complicated.

- It heavily depends on type of tree (and DBH)

- If this data is available, it can be combined with ecological models or remote sensing & ML to project carbon sequestration.

## Problem

- Machine learning has become an increasingly utilized method of classifying land cover.

- Although many methods use spatial data from remote sensing images, here only cartographic variables are used to predict tree species.

# Data

**Features Variables:**
Elevation, Aspect, Slope,
Horizontal_Distance_To_Hydrology,
Vertical_Distance_To_Hydrology,
Horizontal_Distance_To_Roadways,
Hillshade_9am, Hillshade_Noon,
Hillshade_3pm,
Horizontal_Distance_To_Fire_Points,
Soil_Type (40 binary columns)

**Target Variables:**
Cover_Type (7 types of trees)
  - 1 -- Spruce/Fir
  - 2 -- Lodgepole Pine
  - 3 -- Ponderosa Pine
  - 4 -- Cottonwood/Willow
  - 5 -- Aspen
  - 6 -- Douglas-fir
  - 7 -- Krummholz

# Expected Accuracy

"The overall mean absolute classification accuracy for the neural network method was 70.52%, with a 95% confidence interval of 70.26% to 70.80%." -- Blackard, 1998
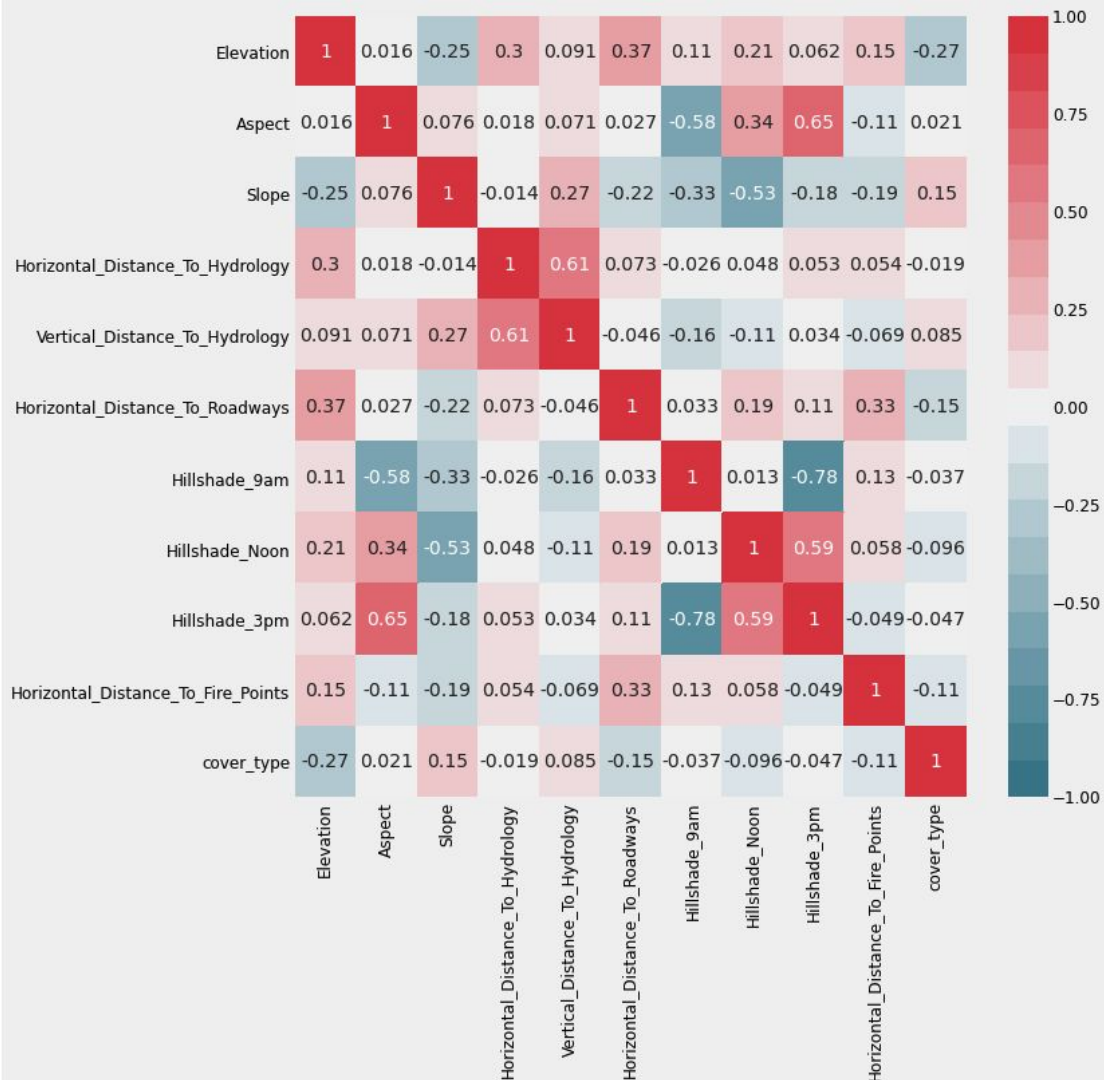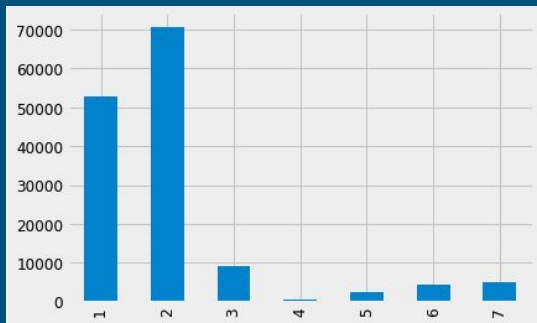
Blackard, Jock A. 1998. "Comparison of Neural Networks and Discriminant Analysis in Predicting Forest Cover Types." Ph.D. dissertation. Department of Forest Sciences. Colorado State University. Fort Collins, Colorado. 165 pages.

**Limitations:**

- This may be an oversimplification of the ecological processes underlying tree species distribution and growth
- This data is specific to Colorado & only 7 tree species. The models here wouldn't be applicable to tropical forests

# EDA

- 581,012 rows and 51 columns
  - Subsetted ¼ of the data
  - 145,253 rows and 51 columns
- Numeric features were not normally distributed
- Target feature had significantly more observations in class 1 & 2

# Results: Logistic Regression

Simple Logistic Regression

GridSearchCV and K-Fold Cross Validation

L2 penalty (Lasso Regression)

**Accuracy:**

**Train : 0.70**

**Test: 0.70**

** Null model would predict an accuracy of 0.49 based on the most frequent cover type

**Accuracy:**

**Train : 0.70**

**Test: 0.70**

*precision & recall scores were significantly lower for less frequent observations

**Accuracy:**

**Train : 0.70**

**Test: 0.70**

# Results: Random Forest Classifier



Feature importance →

| Random Forest | Hyper Parameter Tuning |
|---|---|

**Accuracy:**

**Train : 1.0**

**Test: 0.91**

*Overfitting the model (high variance)

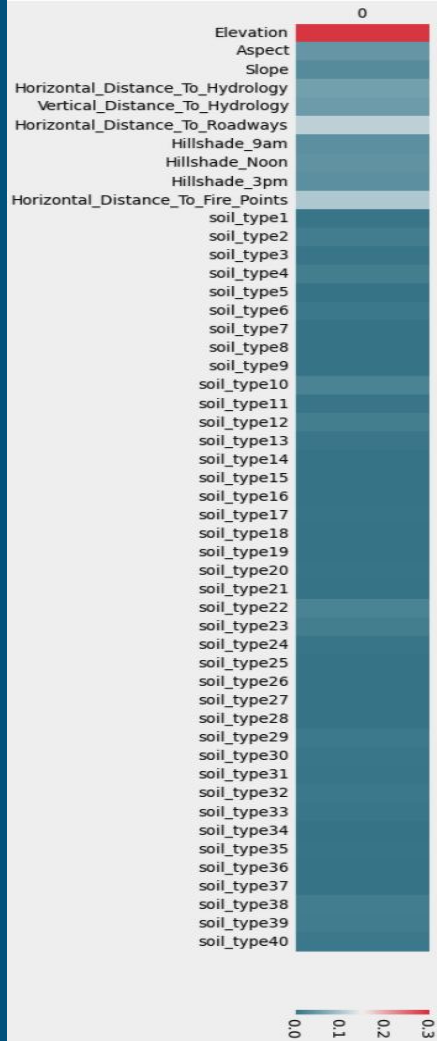**Accuracy:**

**Train : 0.99**

**Test: 0.92**

*Overfitting the model (high variance)

# Discussion

**Tuning more or other models:**

- Neural Networks

**Recommendations:**

- Using remote sensing data can likely predict cover type better than cartographic data
- It could be that a combination of remote sensing data & cartographic data fed into ML algorithms could produce accurate classification of tree cover type