

# Forest Cover Type Prediction

W207 Final Project

April, 2021

Authors: Aidan Jackson, Andi Morey, Naga Chandrasekaran, and Scott Gatzemeier



# Agenda

1. Exploratory Data Analysis
2. Clean/Format Data and Feature Engineering
3. Initial Machine Learning Models
4. Hyperparameter Tuning
5. Evaluation of the Best Model
6. Interpret Model Results
7. Summary & Conclusions

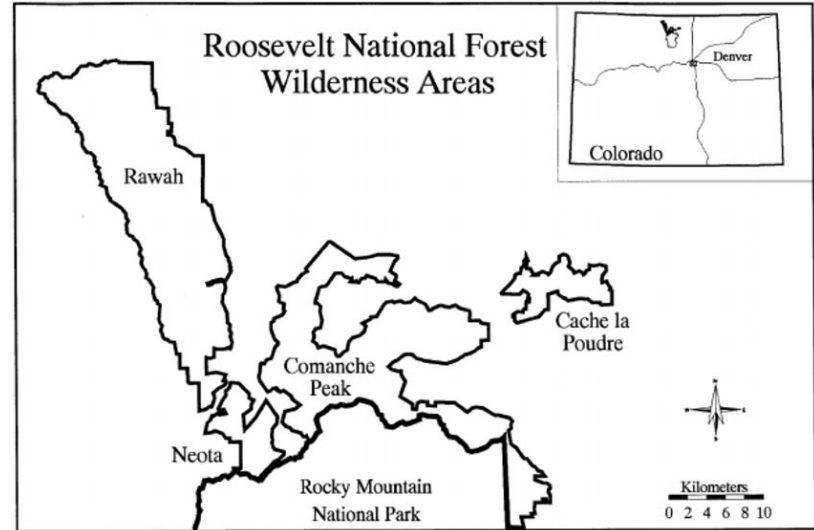


# Problem Statement

## Cover Types

Spruce/Fir  
Lodgepole Pine  
Ponderosa Pine  
Cottonwood/Willow  
Aspen  
Douglas-fir  
Krummholz

## Areas



Classify the cover type for areas

# Dataset Information

## Data



- 581,012 instances
- 565,892 test set
- 15,120 training set

## Attributes



- Numeric, Categorical
- 12 categories
- 54 total columns
- 44 data columns binary

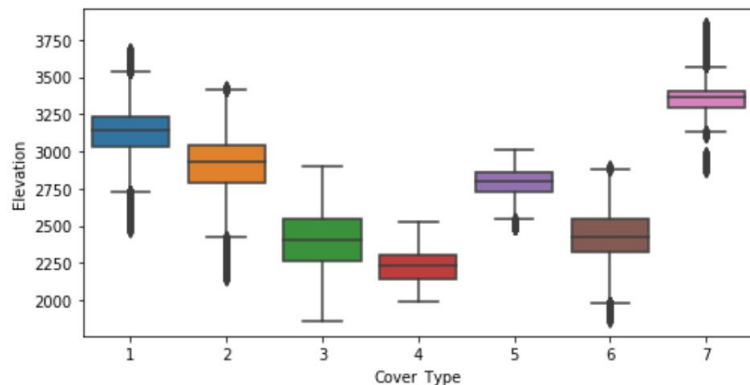
## Attribute Details

1. Elevation
2. Aspect
3. Slope
4. Horizontal distance to hydrology
5. Vertical distance to hydrology
6. Horizontal distance to roadways
7. Horizontal distance to firepoints
8. Hill shade 9am: RGB representation
9. Hill shade noon: RGB representation
10. Hill shade 3pm: RGB representation
11. Wilderness areas: 4 binary data
12. Soil type: 40 binary data



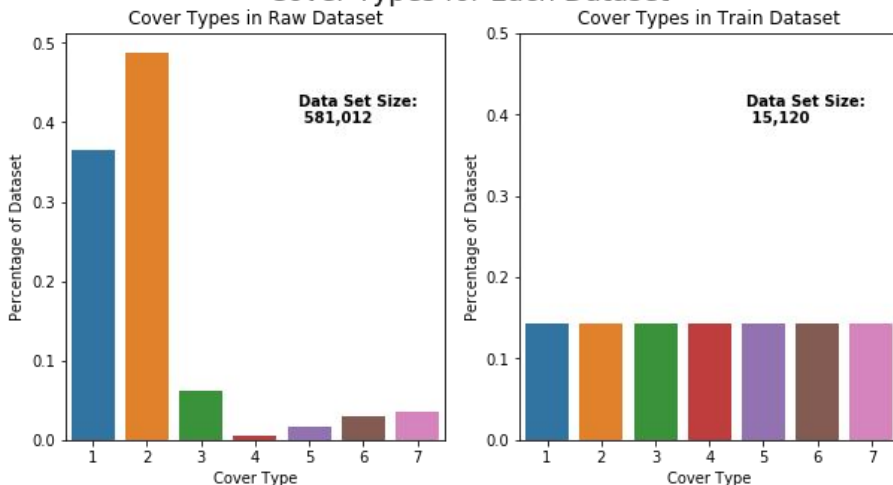
# Exploratory Data Analysis

## Cover Type Key Important Feature



## Cover Type Distributions

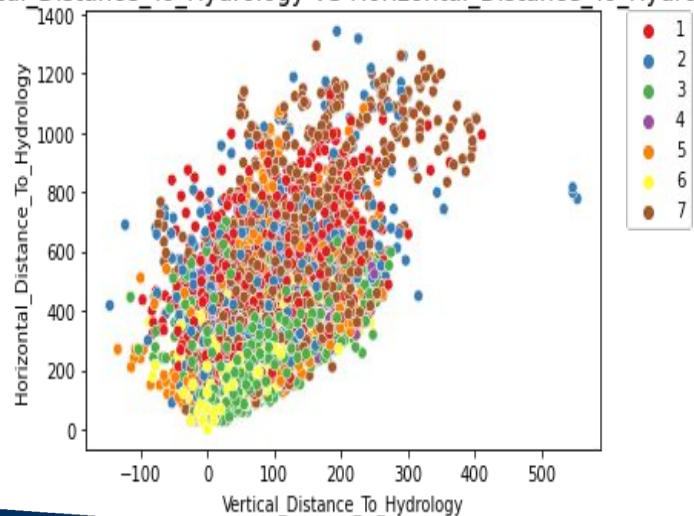
Cover Types for Each Dataset



# Exploratory Data Analysis

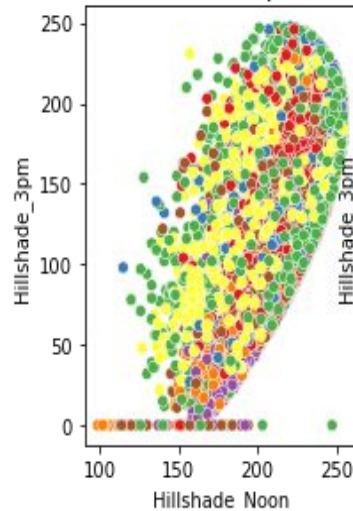
## Training Data Frame

Vertical\_Distance\_To\_Hydrology VS Horizontal\_Distance\_To\_Hydrology

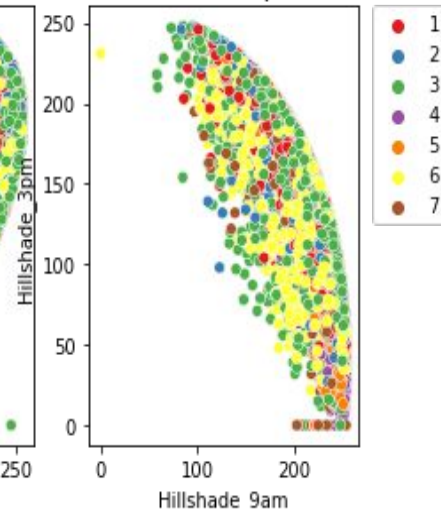


## Hillshade Correlation

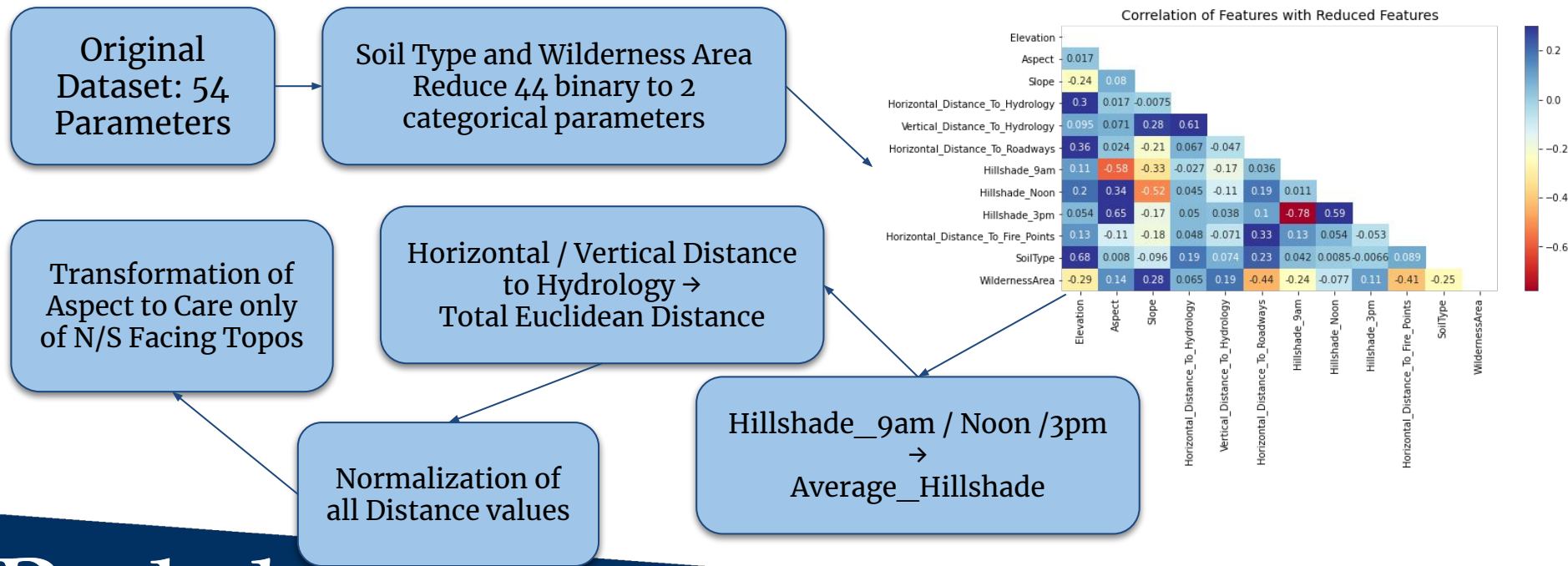
Noon VS 3pm



9am VS 3pm



# Feature Engineering





# Machine Learning Model Progression

Tuned & Featured Engineered Model Results

Model	Kaggle Accuracy, Before (%)	Kaggle Accuracy, After (%)
K-Nearest Neighbor	63	71
Naive Bayes	42	42
Logistic Regression	40	59
Decision Tree	66	77
Neural Network	35	72
Tie Breaker	-	72

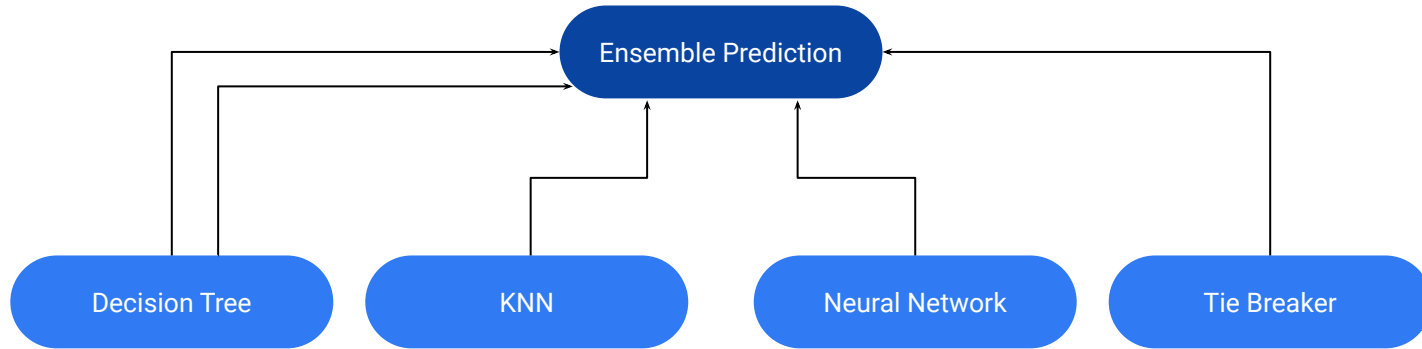


# Hyperparameter Tuning

- Naive Bayes and Logistic Regression discarded due to low accuracy
- Random Forest had best individual performance
- KNN:  $K = 1$ , Euclidean distance
- NN: 9 hidden layers, 100 nodes each

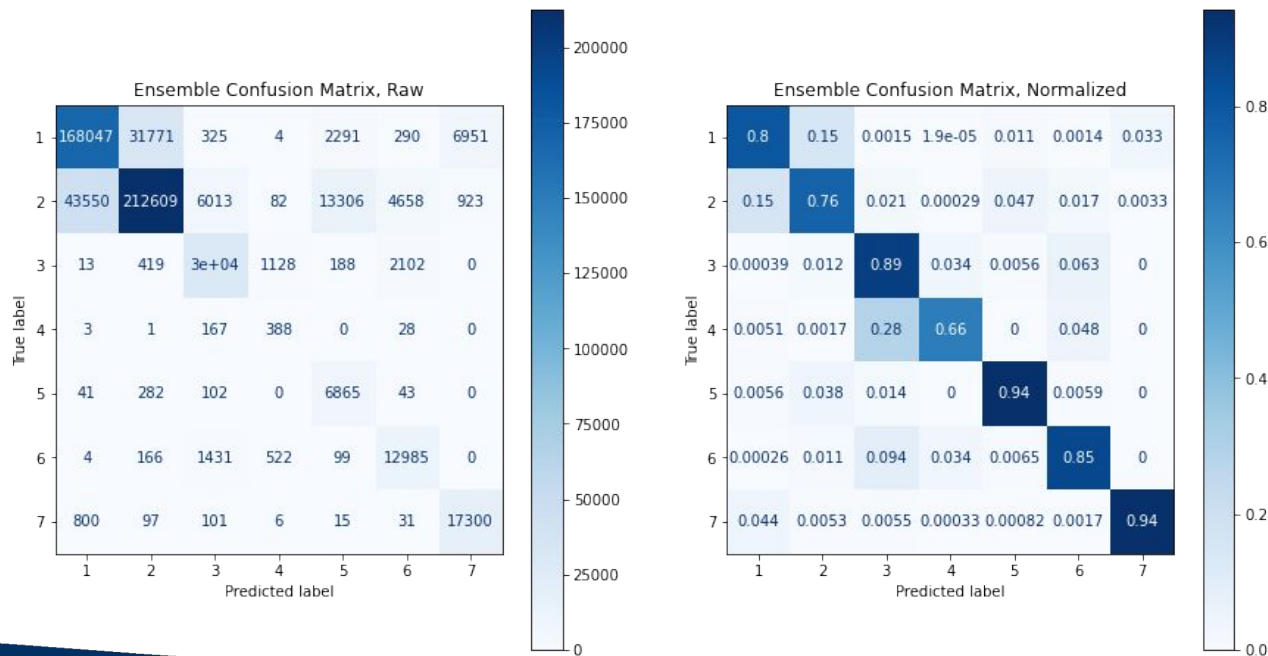
# Evaluation of the Best Model

- Best accuracy: 79.579%
- Position 197 / 1693 on the Leaderboard





# Interpret Model Results

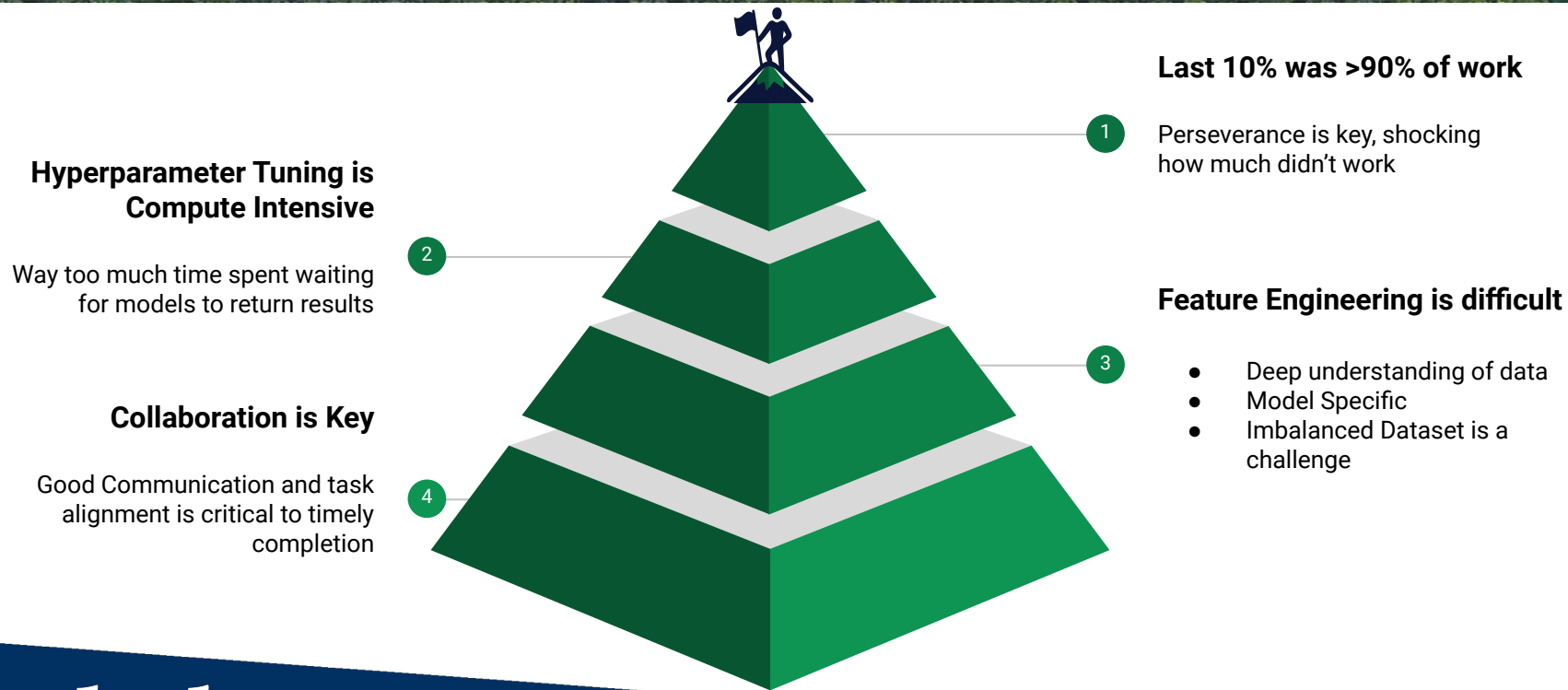


# Summary & Conclusions

- **Ensemble Model Achieved Best Results**
  - Power of Group Consensus
- **Cover Type 1 and 2 Dominate in Test Data**
  - Make up 85% of Data
- **Future Use Case Challenges**
  - Specific to Wilderness areas in Colorado
  - Data Collection Method Unknown



# Lessons Learned





# Questions?