

# Forest Cover Type EDA

We preliminary explored the data to uncover patterns and test hypotheses

Topic	Finding
Data Background	Each data observation represents a 30m x 30m patch of land in Roosevelt National Forest of northern Colorado. Our model will predict a classification for the forest cover type. There are 7 forest cover types.
Data Shape	<ul style="list-style-type: none"><li>● 581,012 Rows</li><li>● 55 Columns</li><li>● 0 Null Values</li></ul>
Column Types	<ul style="list-style-type: none"><li>● 10 variables are Continuous</li><li>● 44 variables are Binary</li></ul>
Test Variable Distribution	<ul style="list-style-type: none"><li>● Class imbalance is present</li><li>● Lodgepole Pine: accounts for 49% of data observations</li><li>● Spruce/Fir: 36%</li><li>● Ponderosa Pine: 6%</li><li>● Douglas Fir, Aspen, Cottonwood/Willow: 5%</li></ul>
Data Challenges	<ul style="list-style-type: none"><li>● 40 different types of Soil variables (all binary)</li><li>● 4 different types of Wilderness variables (all binary)</li></ul>
Variables Correlation	<ul style="list-style-type: none"><li>● Elevation, Aspect, Horizontal Distance to Fires, and Horizontal Distance to Roads variables have high correlation with Forest Cover Types variable</li><li>● The 4 Wilderness variables have low correlation with Forest Cover Types variable</li></ul>