

# Analysis of Cartographic Features and Forest Covertypes

Mikkel Bjornson

10/9/2021

## Introduction

Dominant species in tree cover can affect the structure of an ecosystem. The topographic features of that environment may effect which species dominate. Initial analysis focuses on the effects of elevation and other topographic features on the presence of spruce and fir trees. By quantifying these effects, insights into the niche of these trees can be better understood. The effects of elevation on tree species are of particular interest. In addition, Mapping tree species over large areas can be difficult and costly. Associations between tree species and topography may help estimate the forest cover. The topographic features are explored in relation to predicting forest cover type.

## Data Description

The analysis uses the cover type data (Bache and Lichman 2013). Four wilderness areas within the Roosevelt National Forest are sampled using remote mapping techniques. Cartographic features are obtained for thirty by thirty meter plots. Cover type is obtained from the USFS resource information center. The data contains twelve explanatory variables and one response. Variables include:

Variable Name	Description
Elevation	Meters above sea level that the plot is located
Aspect	Degress from north that the slope of the plot faces
Slope	Steepness of a hill the plot is on, 0 indicating flat
Horz_dist_hydro	Horizontal distance in meters to the nearest body of water
Vert_dist_hydro	Vertical distance in meters to the nearest body of water, negative values indicating water at lower elevations
horz+dist_road	Horizontal distance in meters to the nearest road
Hillshade9	Calculated illumination of the plot at 9am on the summer solistice, based of topographical features, 0 indicating shaded plot
Hillshade12	Calculated illumination of the plot at noon on the summer solistice, based of topographical features, 0 indicating shaded plot
Hillshade3	Calculated illumination of the plot at 3pm on the summer solistice, based of topographical features, 0 indicating shaded plot
Horz_dist_fire	Meters to the nearest forest fire starting point
wild_area	Indicates which of the four wilderness areas the sample comes from

Variable Name	Description
soil_type	Indicates which of the 40 different soil types is present in the plot
cover_type	Indicates the dominant tree species on the plot

Analysis begins with an inferential study of the effects of elevation on spruce/fir cover when controlling for other cartographic variables. The cover type variable is reduced to a binary response. Logistic Regression methods are selected for the inferential analysis (Kutner et al 2004, James et al 2021). A full model, a reduced model removing all distance measures, and a reduced model with only elevation are fit. Drop in deviance testing is used to determine the best model.

Further analysis focuses on the predictive power of cartographic features. Both a multiclass logistic model using a one vs. rest method (Kutner et al 2004, James et al 2021), and a multiclass random forest (James et al 2021) are explored. The multiclass logistic regression model split the cover type variable into 7 binary response variables. Logistic regression models are built for each cover type, using a classification tree to select explanatory variables. Predictions are given by the model with the largest response value. The several random forest models are built with varying numbers of trees and features to split. The model with the lowest validation error is selected.

Initial exploration, revealed apparent differences in elevation for each the seven cover types. There is also apparent differences between soil and cover type. However, there are no recorded instances of spruce cover in Montane alluvial soil (type elu51). It is unknown whether spruce cannot grow in this soil type or it is just not observed. The perfect separation prevents the use of maximum likelihood methods, and thus has to be excluded from the logistical regression analyses. There was no strong graphical evidence of associations with the other variables.

## Statistical Modeling

### Inferential model

Three models are fit to the data set for the purpose of estimating the effect of elevation on Spruce/Fir tree cover. An initial rich model is fit using all available variables.

$$\text{logit}(\text{spruce}/\text{fir}) = \beta_0 + \beta_1(\text{Elevation}) + \beta_2(\text{Aspect}) + \beta_3(\text{Slope}) + \beta_4(\text{Horz\_dist\_hydro}) + \beta_5(\text{Vert\_dist\_hydro}) + \dots \\ \beta_6(\text{horz\_dist\_road}) + \beta_7(\text{Hillshade9}) + \beta_8(\text{Hillshade12}) + \beta_9(\text{Hillshade3}) + \beta_{10}(\text{Horz\_dist\_fire}) + \epsilon$$

Two reduced models are also evaluated:

- reduced model 1:

$$\text{logit}(\text{spruce}/\text{fir}) = \beta_0 + \beta_1(\text{Elevation}) + \epsilon$$

- reduced model 2:

$$\text{logit}(\text{spruce}/\text{fir}) = \beta_0 + \beta_1(\text{Elevation}) + \beta_2(\text{Aspect}) + \beta_3(\text{Slope}) + \dots \\ \beta_7(\text{Hillshade9}) + \beta_8(\text{Hillshade12}) + \beta_9(\text{Hillshade3}) + \epsilon$$

### Predictive model

The first predictive model is a multiclass logistic predictive model fit the following logistic models:

$$\text{logit}(\text{spruce}/\text{fir}) = \beta_0 + \beta_1(\text{Elevation}) + \epsilon$$

$$\text{logit}(\text{Lodge pole pine}) = \beta_0 + \beta_1(\text{Elevation}) + \epsilon$$

$$\begin{aligned} \text{logit}(\text{spruce}/\text{fir}) = \beta_0 + \beta_1(\text{Elevation}) + \beta_2(\text{Vert\_dist\_hydro}) + \beta_3(\text{horz\_dist\_road}) + \dots \\ \beta_4(\text{Hillshade9}) + \beta_5(\text{Hillshade3}) + \epsilon \end{aligned}$$

$$\begin{aligned} \text{logit}(\text{cottonwood}) = \beta_0 + \beta_1(\text{Elevation}) + \beta_2(\text{Horz\_dist\_hydro}) + \beta_3(\text{horz\_dist\_road}) + \dots \\ \beta_4(\text{Hillshade9}) + \beta_5(\text{Horz\_dist\_fire}) + \epsilon \end{aligned}$$

$$\text{logit}(\text{Douglas fir}) = \beta_0 + \beta_1(\text{Elevation}) + \epsilon$$

$$\text{logit}(\text{Krummholz}) = \beta_0 + \beta_1(\text{Elevation}) + \epsilon$$

$$\text{logit}(\text{PonderosaPine}) = \beta_0 + \beta_1(\text{Elevation}) + \beta_2(\text{Hillshae12}) + \epsilon$$

Using a one vs. rest strategy, the model producing the largest response indicates the predicted value.

The second predictive model is a Random Forest model with 200 trees.

## Results

### Inferential study

The full model is chosen over both reduced models using the drop in deviance test (reduced model 1: p-value<0.00001, reduced model 2: p-value<0.00001). The table below indicates all explanatory variables contributed to the model. Isolating elevation, there is evidence to suggest that elevation is associated with the presence of spruce/fir as the dominant tree cover (p-value < 0.00001). With 95% confidence, the logistic model estimates the presence of spruce/fir cover increases by between 87.64% and 89.08% for every 100 meter increase in elevation.

	Estimate	Std. Error	z value	Pr(> z )
(Intercept)	-15.03020	0.24166	-62.19549	0.00000
Elevation	0.00633	0.00002	324.22677	0.00000
Aspect	-0.00046	0.00004	-11.78927	0.00000
Slope	-0.01956	0.00122	-16.00905	0.00000
Horz_dist_hydro	-0.00170	0.00002	-84.84012	0.00000
Vert_dist_hydro	-0.00264	0.00008	-34.36787	0.00000
Horz_dist_road	-0.00004	0.00000	-18.70981	0.00000
Hillshade9	0.00897	0.00143	6.26550	0.00000
Hillshade12	-0.03468	0.00117	-29.69117	0.00000
Hillshade3	0.01594	0.00117	13.60595	0.00000
Horz_dist_fire	-0.00001	0.00000	-3.82884	0.00013

## Predictive Model

The multiclass logistic regression model achieved a test error rate of 0.3549. Examining the per class accuracy rates in the table below, the model accurately estimated Lodge pole pine (accuracy $\approx$  81%) and Spruce/fir (accuracy $\approx$  64%). It suffered from very poor accuracy on the less common cover types.

The Random Forest model appears to fit the data better with a test error rate of 0.0358. The per class accuracy rate is above 90% for most cover types. Cottonwood/willow (accuracy $\approx$  87%) and Aspen (accuracy $\approx$  86%) are the least accurately predicted.

Cover	Logistic	Random_Forest
Spruce/Fir	0.6400	0.9602
Lodgepole Pine	0.8082	0.9765
Ponderosa Pine	0.2774	0.9623
Krummholz	0.0296	0.9561
Cottonwood/Willow	0.0000	0.8648
Aspen	0.0000	0.8401
Douglas-fir	0.0000	0.9096

## Conclusion

Elevation and other cartographic features appear to be associated with the dominant tree species in forests of Roosevelt National Forest. There is an estimated increase in the chance of finding spruce/fir between 87.64% and 89.08% per 100 meter increase in elevation. This suggests that spruce/fir forests perhaps favor higher elevations. In addition, the model finds significant effects on spruce/fir for all other variables. Although not the main focus on this study, this does suggest that the dominant species in forest cover is a complex mix of different environmental factors.

When attempting to predict forest cover, the logistic model failed to capture the complexity of the system with an accuracy rate of about 65%. The Random Forest model managed about a 96% success rate. The ability of random forest models to manage non-monotonic changes is likely a factor in the higher success rates. Further model building using xgboost and neural networks could possibly increase the predictive accuracy even further.

## Citations

Bache, K. & Lichman, M. (2013). UCI Machine Learning Repository <https://archive.ics.uci.edu/ml/datasets/covertypes>. Irvine, CA: University of California, School of Information and Computer Science

Kutner, M. H., Nachtsheim, C. J., & Neter, J. (2004). Applied Linear Regression Models (fourth). McGraw-Hill Irwin.

James, G., Witten, D., Hastie, T., & Tibshirani, R. (2021). An introduction to statistical learning with applications in R. Springer.