

# Machine Learning Competition: Report on Forest Cover Type Prediction

Fernandez, María, McIver, Jordan, Besanson, Gaston

March 27, 2015

## Introduction

THIS REPORT'S purpose is to explain our experience on the Forest Cover Type Prediction's Competition.<sup>1</sup> Our objective in this project was to *predict as best as possible the type of tree in a region with the purpose of minimizing the fires*. Even though that we used the same loss for each type of misclassification -in other words, all trees are equally important-, we decided to create new features<sup>2</sup>. **This feature is explained in detailed in the attached PowerPoint presentation.**

<sup>1</sup> <https://inclass.kaggle.com/c/prediction-of-a-forest-covertime>

<sup>2</sup> The recommended approach would be consider penalizing more the misclassification of very flammable trees

## Data Exploration and Feature Creation

THE DATASET used in this report consists of 54 attributes (or features) and 50.000 observations for the training (plus one class attribute) and 100.000 observations for the competition. The attributes are:

- 10 quantitative data

Feature	Range
Elevation	1859–3858 (m)
Aspect	0–360 (azimuth)
Slope	0–66
Horizontal Distance to Hydrology	0–1397 (m)
Vertical Distance to Hydrology	173–601 (m)
Horizontal Distance to Roadways	0–7117 (m)
Hillshade 9 a.m.	0–255 (index)
Hillshade Noon	0–255 (index)
Hillshade 3 p.m.	0–255 (index)
Horizontal Distance to Fire Points	0–7173 (m)

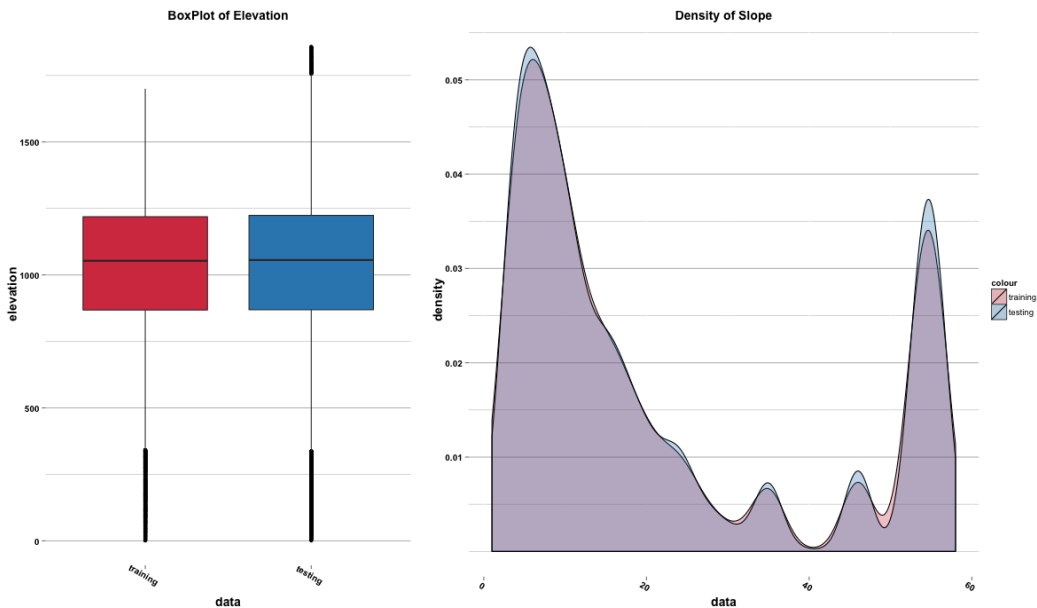
- 4 binary wilderness areas and 40 binary soil type variables

- *Cover Type:*

Class	Tree Type
1	Spruce/fir
2	Lodgepole pine
3	Ponderosa pine
4	Cottonwood/willow
5	Aspen
6	Douglas-fir
7	Krummholz

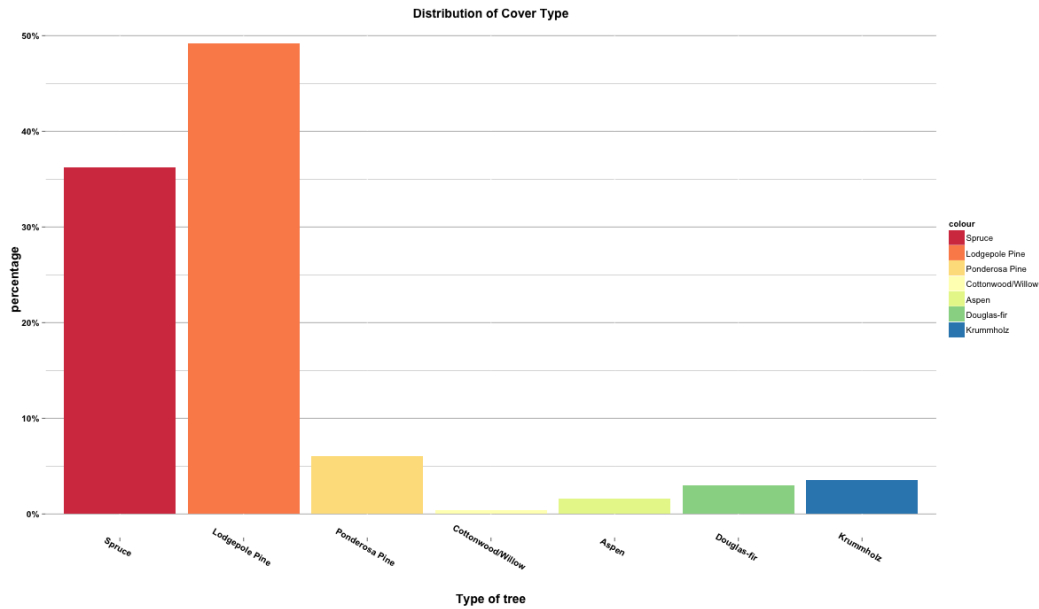
FIRST, we quickly checked that our training set was representative of the testing set data. The means and quantiles of each feature were very similar between the two datasets. That gave us a hint, that the model training error would be a good proxy for the test error; *dismissing our initial concerns of overfitting*. The following box-plot shows the similarity of Elevation of the two datasets. Besides, a density plot of the Slope shows that the two datasets are equally distributed.<sup>3</sup>

<sup>3</sup> Please for additional plots see the presentation.



SECOND, we noticed that the training set is **imbalanced**<sup>4</sup>; meaning that we have two type of trees *Spruce* (Class 1) and *Lodgepole Pine* (Class 2) that represent over 80% of the observations.

<sup>4</sup> This is also present in the original dataset <http://kdd.ics.uci.edu/databases/covertypes/covertypes.html>.



WE CREATED six new variables to try to identify features important to fire risk.<sup>5</sup> Finally, we applied a *normalization* on both the training and the test sets to the 60 features.<sup>6</sup>

New Feature	Type of Variable
No Sun at 3pm	Binary variable
Zero Horizontal Distance to Hydrology	Binary variable
Crow-fly Distance	Continuous variable
Fire Risk	Continuous variable
Shift-Vertical Distance to Hydrology	Continuous variable
Shift-Horizontal Distance to Hydrology	Continuous variable

*For further detail see attached PowerPoint Presentation.*

## Machine Learning Methods

We are going to study the performance of following methods on the dataset.

### Support Vector Machine (SVM)<sup>7</sup>

SVM tries to find a hyperplane that splits the data with the maximum margin. A clear shortcoming for plain-vanilla SVM is that this

<sup>5</sup> we try, but did not continue, to reduce dimensionality through Principal Component Analysis, but the loss of variance decreased the models performance (there were loss of information). We tried the linear and polynomial transformation of the training set data and run it with the Support Vector Machine method. Please see /Report/Graphs/old folder to see Plots and PCA code. Besides, we decided not to collapse the binary variables. Please see: Nagel et al. *Exploring Non-Linear Data Relationships in VR using the 3D Visual Data Mining System*.

<sup>6</sup> `preProcess` function from the `Caret` package.

<sup>7</sup> For this method, we used the `e1071` package

dataset is not linearly separable. To address this, we introduce the kernel transformations and try to classify the points. But, there is a second problem that the classes are imbalanced and this can produce suboptimal models which are biased towards the majority class and have low performance on the minority class, like most of the other classification paradigms.<sup>8</sup> We include weights for the classes to try to overcome this problem.

#### PARAMETERS TO TRAIN:

1. *kernel*: nonlinear kernel function (such as a Gaussian radial basis function<sup>9</sup> or Polynomial kernel).
2. *cost*: is the parameter for the soft margin cost function, which controls the influence of each individual support vector; this process involves trading error penalty for stability.
3. *gamma*: is a kernel parameter,

#### *K* Nearest Neighbor<sup>10</sup>

KNN is a lazy learning algorithms<sup>11</sup>, that stores all available cases and classifies new cases based on a similarity measure. We considered because it doesn't need for the classes to be linearly separable and because it is very simple to implement with only two parameter to consider: the distance metric<sup>12</sup> and the neighbors being considered. On the down side, this algorithm is also sensitive to unbalanced dataset like this one; where infrequent classes are therefore often dominated in most neighborhoods.<sup>13</sup>

#### PARAMETERS TO TRAIN:

1. *k*: is the number of neighbors considered.

#### *Random Forest*<sup>14</sup>

RANDOM FORESTS are an ensemble learning method that construct a number of decision trees at training time and outputting the class that is the mode of the classes output by individual trees. Random Forests are a combination of tree predictors where each tree depends on the values of a random vector sampled independently with the same distribution for all trees in the forest. The basic principle is that a group of "weak learners" can come together to form a "strong learner". *Random Forests are a wonderful tool for making predictions considering they do not overfit because of the law of large numbers.*<sup>15</sup> An important characteristic of this method is that provides methods for balancing error in class population unbalanced datasets.

<sup>8</sup> Further readings, please see: Trebar, Mira, Steele, Nigel, *Application of distributed SVM architectures in classifying forest data cover types*. Batuwita, Rukshan, Palade, Vasile, *Class Imbalance Learning Methods for Support Vector Machines*

<sup>9</sup>  $K(x_i, x_j) = e^{-\frac{\|x_i - x_j\|^2}{2\sigma^2}}$ , where  $\sigma$  is the spread or standard deviation.

<sup>10</sup> For this method, we used the `class` package

<sup>11</sup> While a training dataset is required, it is used solely to populate a sample of the search space with instances whose class is known. No actual model or learning is performed during this phase.

<sup>12</sup> Given that we used the build-in function of a package we were stuck with the Euclidean distance that works well with continuous variables but for categorical data other metrics can be considered.

<sup>13</sup> This can be alleviated through balanced sampling of the more popular classes, which we didn't do for this method

<sup>14</sup> For this method, we used the `caret` package

<sup>15</sup> Further readings, please see: Sug, Hyontai, *Better Induction Models for Classification of Forest Cover*

## PARAMETERS TO TRAIN:

1. *Randomly Selected Predictors*: is the number of attributes to pick randomly to generate each subtree in a tree in the forest. The recommended default parameter value is the square root of the number of parameters.
2. *Number of trees*: The recommended default parameter value is among the hundreds.

*Gradient Boosting Machine (GBM)*<sup>16</sup>

ON THE CONTRARY to Random Forest, that rely on simple averaging of models in the ensemble; the main idea of boosting is to add new models to the ensemble sequentially. At each particular iteration, a new weak, base-learner model is trained with respect to the error of the whole ensemble learnt so far. To establish a connection with the statistical framework, a gradient-descent based formulation of boosting methods was derived. This formulation was called the gradient boosting machines.

IN GBMs, the learning procedure consecutively fits new models to provide a more accurate estimate of the response variable. The principle idea behind this algorithm is to construct the new base-learners to be maximally correlated with the negative gradient of the loss function, associated with the whole ensemble.<sup>17</sup> We picked this method because it introduces a lot of freedom into the model design.

<sup>16</sup> For this method, we used the H2o package. We try first with the gbm package, but was unstable.

<sup>17</sup> Further readings, please see: Natekin, Alexey, Knoll, Alois, *Gradient boosting machines, a tutorial*

## PARAMETERS TO TRAIN:

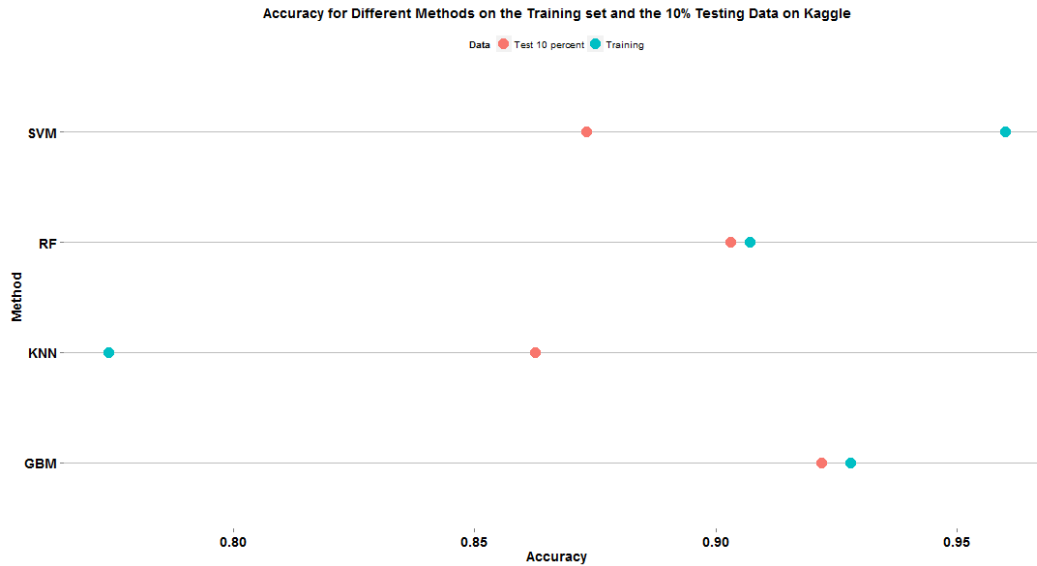
1. *Number of trees*.
2. *Maximum Depth*: The maximum number of edges to generate between the first node and the terminal node.
3. *Minimum Rows*: The minimum number of observations to include in a terminal leaf.
4. *Shrinkage*: The rate at which the algorithm should converge.

*Results*

FOR EACH METHOD we try different parameters, especially with the ensemble methods that where our top picks for this competition. To pick among these parameters, we did 10-fold cross validation. Given the amount of parameters that we try for Random Forest and GBM, the training of the models were done in batches to reduce the computational burden of not being able to use the Virtual Machines.

Please see in the Appendix the detail information on the training of the models for the different methods.

THE FOLLOWING plot shows the best model for each method and its accuracy in the training set and in the 10% leaderboard from the Kaggle Competition.



## Conclusion

CONSIDERING our results on the training set and acknowledging the results on the 10% of the Kaggle Competition, our final submission for this competition is the GBM model with the following parameters:

1. *Number of trees*: 250
2. *Maximum Depth*: 18
3. *Minimum Rows*: 10
4. *Shrinkage*: 0.1

FINALLY, we want to show the Confusion matrix for this model:

Predicted								
Actual	1	2	3	4	5	6	7	Error
1	16523	1501	1	0	11	0	65	0.08718
2	1016	23382	59	0	71	47	16	0.04916
3	0	113	2811	15	4	83	0	0.07105

4	0	1	52	131	0	16	0	0.34500
5	19	250	16	0	520	5	0	0.35802
6	4	125	198	10	2	1172	0	0.22435
7	177	18	0	0	0	0	1566	0.11073
Totals	17739	25390	3137	156	608	1323	1647	0.07790

*The larger misclassification happens between Class 1 and Class 2 as we mention in the start of this Report, it accounts for 66% of the misclassification between Classes.*

## Appendix

### Random Forest

AT FIRST we used 3 and 5 folds Cross Validation to discard some of parameters.

mtry	Accuracy	CV	Public Score (10%)
8	0.8548800	3 Folds	
16	0.8802400	3 Folds	
24	0.8843601	3 Folds	
32	0.8862000	3 Folds	
40	0.8865000	3 Folds	
48	0.8868199	3 Folds	
56	0.8872399	3 Folds	
59	0.8867800	3 Folds	
8	0.8620399	5 Folds	
16	0.8894999	5 Folds	
24	0.8926399	5 Folds	
32	0.8953000	5 Folds	
40	0.8958948	5 Folds	
48	0.8992497	5 Folds	
56	0.8970799	5 Folds	
59	0.8956799	5 Folds	
48	0.9030799	10 Folds	0.90770

mtry	Accuracy	CV	Public Score (10%)
56	0.9017999	10 Folds	

### GBM

INITIALLY we left untouched the following parameters shrinkage (equal to 0.1) and the Minimum Rows (equal to 10). But when we found the best GBM model we try a grid for different Minimum Rows: 1,5,7.

Depth	Trees	Accuracy	CV	Public Score (10%)
8	100	0.87672	3 Fold	
10	100	0.89244	3 Fold	
12	100	0.90042	3 Fold	
14	100	0.90382	3 Fold	
10	250	0.89106	3 Fold	
12	250	0.90648	3 Fold	
14	250	0.90692	3 Fold	
16	250	0.90904	3 Fold	
14	500	0.90506	3 Fold	
16	500	0.90652	3 Fold	
14	250	0.9145	5 Fold	
16	250	0.91694	5 Fold	
14	500	0.9124	5 Fold	
16	500	0.91258	5 Fold	
16	250	0.91872	10 Fold	
<b>18</b>	<b>250</b>	<b>0.9221</b>	<b>10 Fold</b>	<b>0.92790</b>
22	250	0.92094	10 Fold	

### KNN

K	Accuracy	Public Score (10%)
<b>1</b>	<b>0.77418</b>	<b>0.8626</b>



K	Accuracy	Public Score (10%)
2	0.74394	
3	0.76902	

## SVM

Cost	Gamma	Accuracy	Public Score (10%)
0.1	0.5	0.77244	
1	0.5	0.89456	
10	0.5	0.96026	
0.1	1	0.71196	
1	1	0.92776	
10	1	0.98742	
0.1	2	0.56352	
<b>1</b>	<b>2</b>	<b>0.96114</b>	<b>0.87320</b>
10	2	0.998	