

# Forest Fires in Portugal

João Teixeira - up200705307, Nuno Peixoto - up200804621

25-01-2015



# Introduction

- Monitoring and forecasting forest fires in Portugal;
- The several variables may influence the burnt area;
- In 2003, Portugal faced the worst forest fire losing 8.6% of of the total area;
- Elevation, slope or density are some of the specifications of the data set;

**Objective** : Explore and predicte the data of the forest.

# Exploratory analysis of the data

- Global Summary
- Main Variables
- Target Variable

# Global Summary

- Number of Columns: 81.
- Number of Rows: 990.
- Number of Data: 80190.
- Target Value: 1 (TotalBurntArea) - Numeric variable.
- Number of Unknown Values: 0.

## Global Summary (cont.)

**Climate Variables** - The climatic conditions may affect the probability of a fire to occur;

**Landscape Variables** - The landscape has been extensively associated with fire occurrence;

**Socio-economic Variables** - Human have impact in historical fire patterns;

**Topographic Variables** - The topographic features may influence the fire ignitions;

# Main Variables

In the following table we have the **TOP5** main variables:

attr_importance	attribute
0.2037	ELEV_MAX
0.1962	bio1
0.1926	ELEV_MEAN
0.1898	bio7
0.1844	DensPop01

## Main Variables (Number of outliers)

- ELEV\_MAX: 8 (0.81%)
- Bio1: 21 (2.12%)
- ELEV\_MEAN: 9 (0.91%)
- Bio7: 1 (0.1%)
- DensPop01: 132 (13.33%)



## Main Variables (Standard Deviation)

- ELEV\_MAX: 339.100654
- Bio1: 14.710837
- ELEV\_MEAN: 251.9971412
- Bio7: 30.9059137
- DensPop01: 1222.3683295

# Target Variable

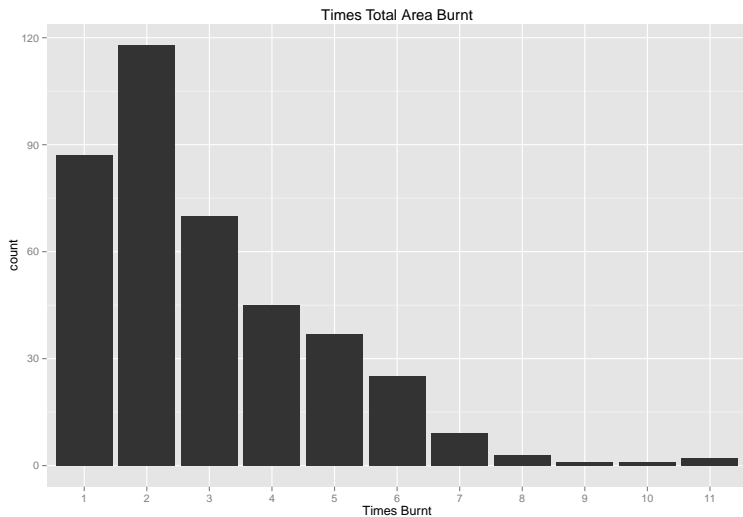
Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
0	0	609	2550	2752	68981

## Target Variable (Number of outliers)

- TotalBurntArea: 106 (10.71%)

We can see that more than 10% of the total burnt area values are considered outliers.

# Target Variable (Total Area vs. Total Burnt Area)



# Data Pre-Processing

- Remove None importance Variables
- Normalizing Value

## Remove None importance Variables

attr_importance	attribute
NaN	TCI_STD
NaN	LPI
NaN	ED
NaN	FRAC_SD
NaN	IJI
NaN	ENN_AM
NaN	eucalipto_AREA_perc
NaN	outfolhosas_AREA_perc

# Normalizing Value

- Data normalization pre-processing we will use for the analysis;

##	ELEV_MAX	ELEV_MEAN	ELEV_STD	SLOPE_MAX	SLOPE_MEAN	SLOPE
## 1	396	168.5080	76.7385	44.9590	18.87750	11.9
## 2	706	604.4890	42.7725	39.1152	8.99396	6.0
## 3	88	34.2032	23.7021	14.3287	2.32026	1.7

##	ELEV_MAX	ELEV_MEAN	ELEV_STD	SLOPE_MAX	SLOPE_MEAN
## 1	-0.2493703	-0.569479	0.2047242	0.5829093	1.0775585
## 2	0.6648126	1.160624	-0.4090371	0.2295771	-0.3401967
## 3	-1.1576551	-1.102441	-0.7536369	-1.2690826	-1.2975129

# Predictive Models

In order to find the best regression model that can predict the target variable of the test data set with less error, we analysed the following forecasting models:

- Multiple Linear Regression
- Regression Trees
- K-Nearest Neighbors (KNNs)
- Support Vector Machines (SVMs)
- Artificial Neural Networks (ANNs)
- Random Forest (Ensembles)



# Model comparison

The metric evaluation to be considered for these forecasts is the MAE - Mean Absolute Error. We use a cross validation method with 2 repetitions of 3 folds.

To evaluate the models we will use the performanceEstimation package that provides a set of functions and arguments that allow us to change the values of parameters in order to check the best fit for an specific model.

# Clustering

The following clustering methods were used to try to find different groups of observations present in the data set:

- Clustering Large Applications (CLARA)
- Partitioning Around Medoids (PAM)
- Hierarchical Clustering
- K-Means Clustering

## Clustering results

Using a R script with the help of the silhouette function we could find the best number of clusters for each used method

CLARA		PAM		hclust	
nClusters	SilhCo	nClusters	SilhCo	nClusters	SilhCo
2	0.7144933	2	0.7336859	2	0.55027091
4	0.6361579	3	0.6187632	3	0.32391734
3	0.6243904	4	0.4812423	5	0.17051099
5	0.3102661	5	0.4318594	4	0.16893921
10	0.2789919	6	0.3595906	6	0.09489308
9	0.2749504	7	0.3234947	7	0.09160570
6	0.2645142	9	0.2750479	8	0.08740959

# Silhouette Plot

Silhouette plot of  $(x = k2\$cluster, \text{dist} = \text{dist}(\text{train\_data[, -73]})$

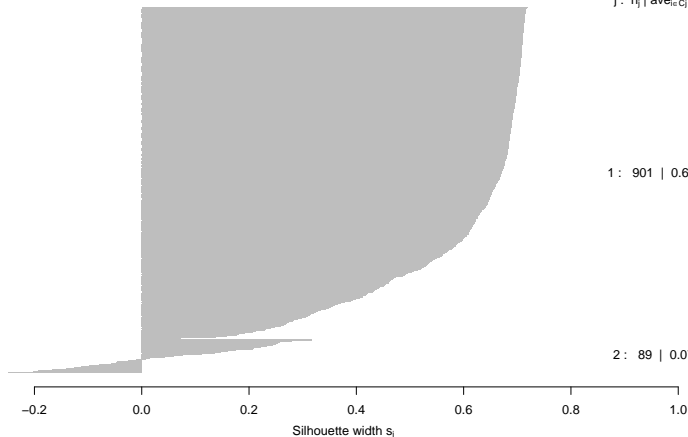
$n = 990$

2 clusters  $C_j$

$j : n_j \mid \text{ave}_{i \in C_j} s_i$

1 : 901 | 0.60

2 : 89 | 0.07



Average silhouette width : 0.55