

# Data Processes Project

*Alejandro Basco Plaza, Miriam Zaragoza Pastor, Jonatan Ruedas Mora, Raúl Cruz Benita  
and Laura Sánchez de Rojas Huerta*

*20/12/2019*

## Contents

Abstract . . . . .	1
Introduction and related work . . . . .	1
Dataset description . . . . .	2
Data wrangling . . . . .	3
Descriptive statistics analysis . . . . .	4
PCA analysis . . . . .	5
Strength of relationships . . . . .	5
Prediction . . . . .	5
Results . . . . .	5
Discussion and future work . . . . .	5

## Abstract

This is an R Markdown document. Markdown is a simple formatting syntax for authoring HTML, PDF, and MS Word documents. For more details on using R Markdown see <http://rmarkdown.rstudio.com>.

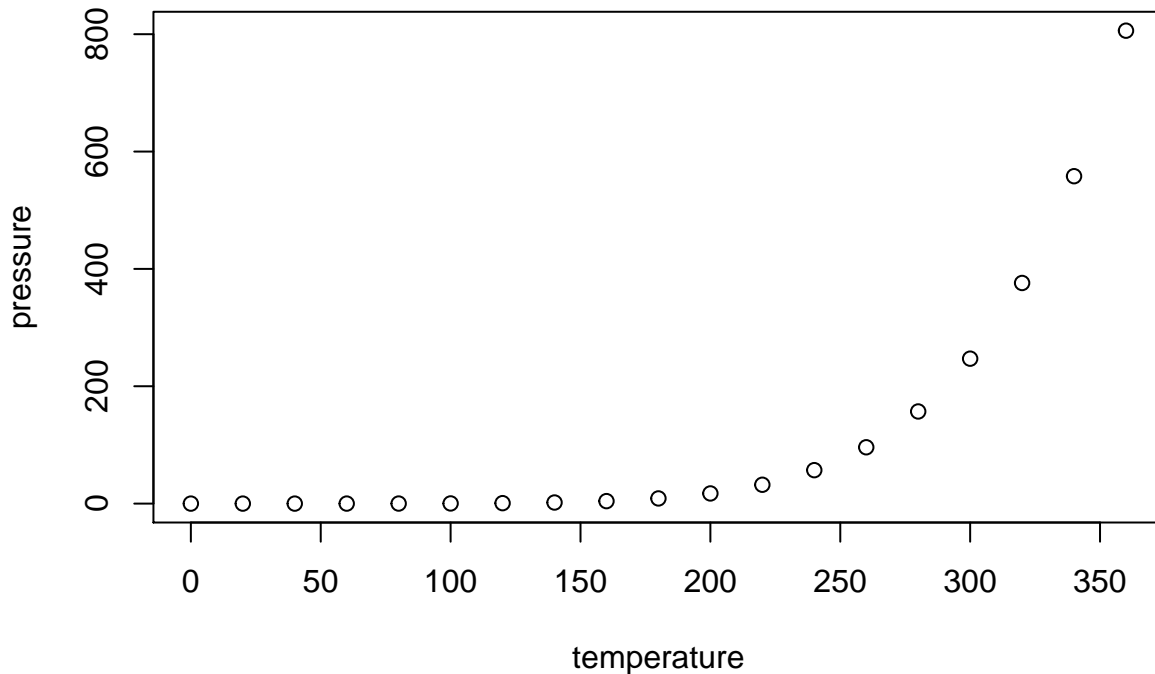
When you click the **Knit** button a document will be generated that includes both content as well as the output of any embedded R code chunks within the document. You can embed an R code chunk like this:

```
summary(cars)
```

```
##      speed      dist
##  Min.   : 4.0    Min.   :  2.00
##  1st Qu.:12.0    1st Qu.: 26.00
##  Median :15.0    Median : 36.00
##  Mean   :15.4    Mean   : 42.98
##  3rd Qu.:19.0    3rd Qu.: 56.00
##  Max.   :25.0    Max.   :120.00
```

## Introduction and related work

You can also embed plots, for example:



Note that the `echo = FALSE` parameter was added to the code chunk to prevent printing of the R code that generated the plot. `## Exploratory data analysis`

### Dataset description

Before exposing the results of the exploratory data analysis carried the data, important information and background about the dataset is introduced below:

The dataset used for this work was obtained from Kaggle's website. In this page, we can find a lot of information and datasets about the air quality in Madrid.

The data in this data set has been collected from the original files provided by Madrid Open Data. Decide soluciones organization processed them and uploaded to the Kaggle's website.

Originally, Madrid Open Data get the data from 24 automatic stations around Madrid. Those stations, which are set around all the districts of the city, record information about pollution in the area. The data was generated from those stations whose pollution sensors measure air quality in Madrid city.

Regarding the observations, all csv files include twelve month data, except the last one that only has data until May. According to this, we have the followings rows:

CSV	Observations
Stations.csv	24 rows
Madrid2001.csv	217872 rows
Madrid2002.csv	217296 rows
Madrid2003.csv	243984 rows
Madrid2004.csv	245496 rows
Madrid2005.csv	237000 rows
Madrid2006.csv	230568 rows
Madrid2007.csv	225120 rows
Madrid2008.csv	226392 rows
Madrid2009.csv	215688 rows
Madrid2010.csv	209448 rows

CSV	Observations
Madrid2011.csv	209928 rows
Madrid2012.csv	210720 rows
Madrid2013.csv	209880 rows
Madrid2014.csv	210024 rows
Madrid2015.csv	210096 rows
Madrid2016.csv	209496 rows
Madrid2017.csv	210120 rows
Madrid2018.csv	69096 rows

We can also differentiate two domains whose columns represent different data:

On the one hand, we have Stations.csv with six columns, which contains information about the stations previously mentioned, in which the data about the pollution levels of the air is collected. On the other hand, we have Madrid20xx.csv, which contain the data about Madrid’s air pollution itself. Some of these csv files have 14 columns, others 16 and others 17. We have the following columns in the csv files:

CSV	Features
Stations.csv	6 columns
Madrid2001.csv	16 columns
Madrid2002.csv	16 columns
Madrid2003.csv	16 columns
Madrid2004.csv	17 columns
Madrid2005.csv	17 columns
Madrid2006.csv	17 columns
Madrid2007.csv	17 columns
Madrid2008.csv	17 columns
Madrid2009.csv	17 columns
Madrid2010.csv	17 columns
Madrid2011.csv	14 columns
Madrid2012.csv	14 columns
Madrid2013.csv	14 columns
Madrid2014.csv	14 columns
Madrid2015.csv	14 columns
Madrid2016.csv	14 columns
Madrid2017.csv	16 columns
Madrid2018.csv	16 columns

## Data wrangling

In order to work with the previously described data, some steps were performed in order to make the data more manageable:

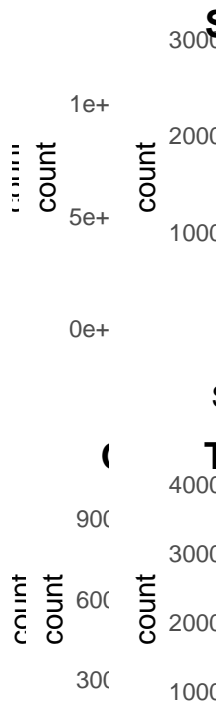
- All the yearly datasets were unified in a single one, adding an additional feature to the existing date column that all of the datasets present: the year. Thus, after this transformation, a single dataset contained all the information about Madrid’s air pollution in the last 17 years, instead of 17 different datasets (one per year).
- In the resulting dataset, there were several missing values: as it has been shown in the table above, some datasets have a different number of columns, so when the previous step was performed and the datasets were unified in a single one, the number of missing values of the resulting dataset was very high. In order to solve this, the rows containing missing values have been removed from the dataset.

## Descriptive statistics analysis

To start the analysis, we are going to analyse the distribution of the variables in our dataset in terms of center, dispersion, skewness and kurtosis descriptive measures.

At first, to observe the distribution of the 10 quantitative variables in our airquality dataset, the histogram for each variable is provided below:

```
##           date  BEN   CO  EBE NMHC  NO_2   O_3   PM10 SO_2  TCH
## 1 2001-08-01 01:00:00   NA 0.37   NA   NA 58.40 34.53 105.00 6.34   NA
## 2 2001-08-01 01:00:00 1.50 0.34 1.49 0.07 56.25 42.16 100.60 8.11 1.24
## 3 2001-08-01 01:00:00   NA 0.28   NA   NA 50.66 46.31 100.10 7.85   NA
## 4 2001-08-01 01:00:00   NA 0.47   NA   NA 69.79 40.65  69.78 6.46   NA
## 5 2001-08-01 01:00:00   NA 0.39   NA   NA 22.83 66.31  75.18 8.80   NA
## 6 2001-08-01 01:00:00 2.11 0.63 2.48 0.05 66.26 33.50 122.70 6.36 1.23
##           TOL station
## 1         NA 28079001
## 2 10.82 28079035
## 3         NA 28079003
## 4         NA 28079004
## 5         NA 28079039
## 6 13.28 28079006
```



PCA analysis

Strength of relationships

Prediction

Results

Discussion and future work