



Asociación Española para la Inteligencia Artificial (AEPIA)

UIMP

Universidad Internacional
Menéndez Pelayo

Practical: Data Preparation

Máster Universitario en Investigación en
Inteligencia Artificial

Ciencia de Datos y Aprendizaje Automático

José Hernández-Orallo

Translated from original material by M.José Ramírez-Quintana
ETSINF

Universitat Politècnica de València

September 20, 2016

- Exercise 1: Inspection of data.

The “titanic.csv” file (available on the platform) contains data on the sinking of the Titanic. Copy the file in your working directory. Then, go to R and use the command

```
titanic <- read.csv(file.choose(),header=TRUE, sep=',')
```

and choose the relevant csv file. You may write the file name (including the path to the file) instead of `file.choose()`, i.e. “titanic.csv”.

Notice how you can change the field separator character according to what is used in the csv file, so that the file is interpreted in the correct way. Show the names of the columns. Observe that the first column (whose name is "X") is redundant (it denotes the identifier of each instance) so it could be removed. To do this, use the `subset` command as follows (use help if needed):

```
titanic<-subset(titanic,select=-X)
```

Now try the following commands:

```
> titanic
> head(titanic)
> summary(titanic)
> plot(titanic)
```

Which variables are quantitative and which variables are categorical? How can we know it?

- Exercise 2: Working with basic graphics.

Download the file “cars.csv” from the platform. This file contains information about the speed and stopping distances of cars.

- 2.1 Make a plot of the distance field in terms of the speed field (use the \$ syntax).
- 2.2 Make a histogram of the distance variable.
- 2.3 Make a histogram of the speed variable.
- 2.4 Modify the previous plots to show the name of the variables (“speed” or “distance”) as the title of the axis. Change the title of the three graphics, and also use colours for the histograms and titles. Save the new graphics as pdf files.

- Exercise 3: Transformations of variables and datasets.

Remove the first column of the cars data frame. Now, assume that data from two more cars are made available:

speed	dist
21	47
34	87

- 3.1 Construct a new data frame with the above data.
- 3.2 Add the constructed data frame to the `cars` data frame.

3.3 Sort the data in the resulting dataset by column `speed` (ascending). There are two ways to do this: using the `order()` command or combining the `with` and the `order()` commands. (Suggestion: search on the internet “how to sort a data frame by columns”).

- Exercise 4: Data manipulation. Download the file “airquality.csv” from the platform. This dataset contains some New York air quality measurements. Solve the following questions:

1. Extract the first 2 rows of the data frame and print them to the console. What does the output look like?
2. How many observations (i.e. rows) are in this data frame?
3. What is the value of `Ozone` in the 40th row?
4. How many missing values are in the `Ozone` column of this data frame?
5. What is the mean of the `Ozone` column in this dataset? Exclude missing values (coded as NA) from this calculation.
6. Extract the subset of rows of the data frame where `Ozone` values are above 31 and `Temp` values are above 90. What is the mean of `Solar.R` in this subset?

- Exercise 5: Data transformation (2).

With the data frame “airquality.csv” solve the following exercises:

1. Discretise the `Ozone` column into five bins (`‘bin1’`, `‘bin2’`, ...) of equal width and a sixth bin (`‘binNA’`) for NA.
2. Discretise the `Solar` column into four bins of equal size and a fifth bin for NA.
3. Create a new column `AbsDay` from the columns `Month` and `Day` such that counts the number of days passed from `Month=5` and `Day=1`.

- Exercise 6: Data transformation (3).

With the data frame “titanic” solve the following exercises:

1. Numerise the `class` column, where `Crew=4`, `1st=3`, `2nd=2` and `3rd=1`.
2. Transform the `titanic` data frame into a new data frame (`titanic2`) with as many examples as passengers using the `Freq` column. In other words, there should be no rows for those for which `Freq=0` and there should be 35 replicated rows for those with `Freq=35`.

3. Compare the plots of the original titanic data frame with the new one.

- Exercise 7: Data selection.

1. Calculate a correlation matrix for the air dataset. Do you see a pair of attributes that are redundant?
2. Calculate a correlation matrix for the cars dataset. Do you see a pair of attributes that are redundant?
3. Using the data frame 'air', perform a simple random sampling of 50 examples.
4. Using the data frame 'air', perform a stratified random sampling of 5 examples of each month.