

# Advanced Statistical Modelling: Logistic Regression

*Ricard Meyerhofer & Joel Cantero*

*4/11/2019*

## Exploratory data analysis

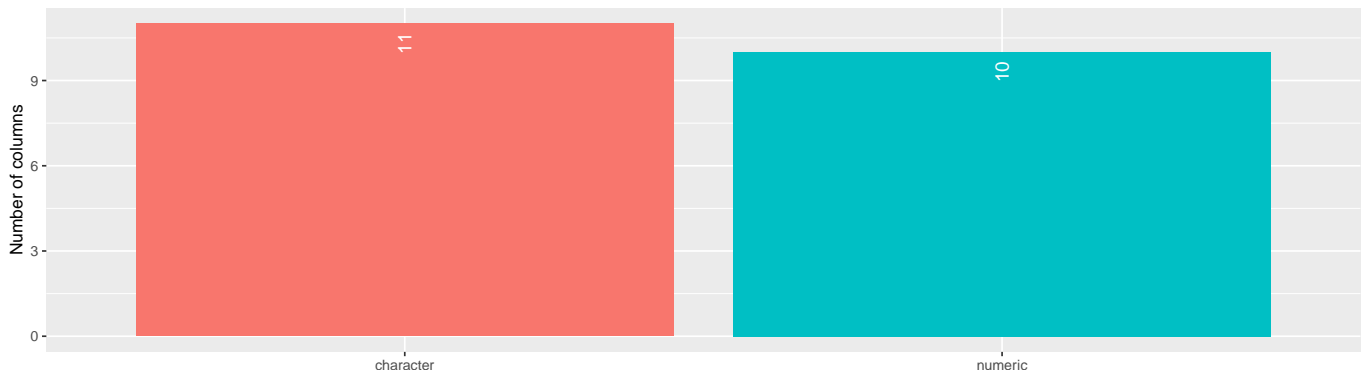
As explained in the problem statement, our dataset is composed by 28645 calls from JYB. JYB has the purpose of reducing the telemarketing costs by decreasing the number of calls to clients not likely to buy the product. This is the list of the available variables:

Variable	Description	Attribute type
id	Customer ID	Client
age	age in years	Client
job	(admin., blue-collar, entrepreneur, housemaid, management, retired, self-employed, services, student, technician, unemployed, unknown)	Client
marital	Marital status (Divorced, married, single, unknown)	Client
education	Education level (basic.4y, basic.6y, basic.9y, high.school, illiterate, professional.course, university.degree, unknown)	Client
default	is he/she a defaulter? (No, yes, unknown)	Client
housing	does he/she has a mortgage? (No, yes, unknown)	Client
loan	does he/she has a personal loan? (No, yes, unknown)	Client
contact	phone type (cellular, telephone)	Call
month	month of the call	Call
day_of_week	day of the call (mon, tue, wed, thu, fri)	Call
campaign	Number of contacts made this campaign for this client (including the current one)	Campaign
pdays	number of days that have passed since the customer was contacted for the last time for a previous campaign (999 means that it was not previously contacted)	Campaign
previous	number of calls made to this client before this campaign	Campaign
poutcome	previous campaign result (failure, nonexistent, success)	Campaign
emp.var.rate	employment variation rate (quarterly)	Indicators
cons.price.idx	Consumer Price Index (monthly)	Indicators
cons.conf.idx	Consumer confidence index (monthly)	Indicators
euribor3m	euribor a 3 meses (daily)	Indicators
nr.employed	number of employed (quarterly)	Indicators
Y	The customer subscribed the deposit? (yes,no)	Response

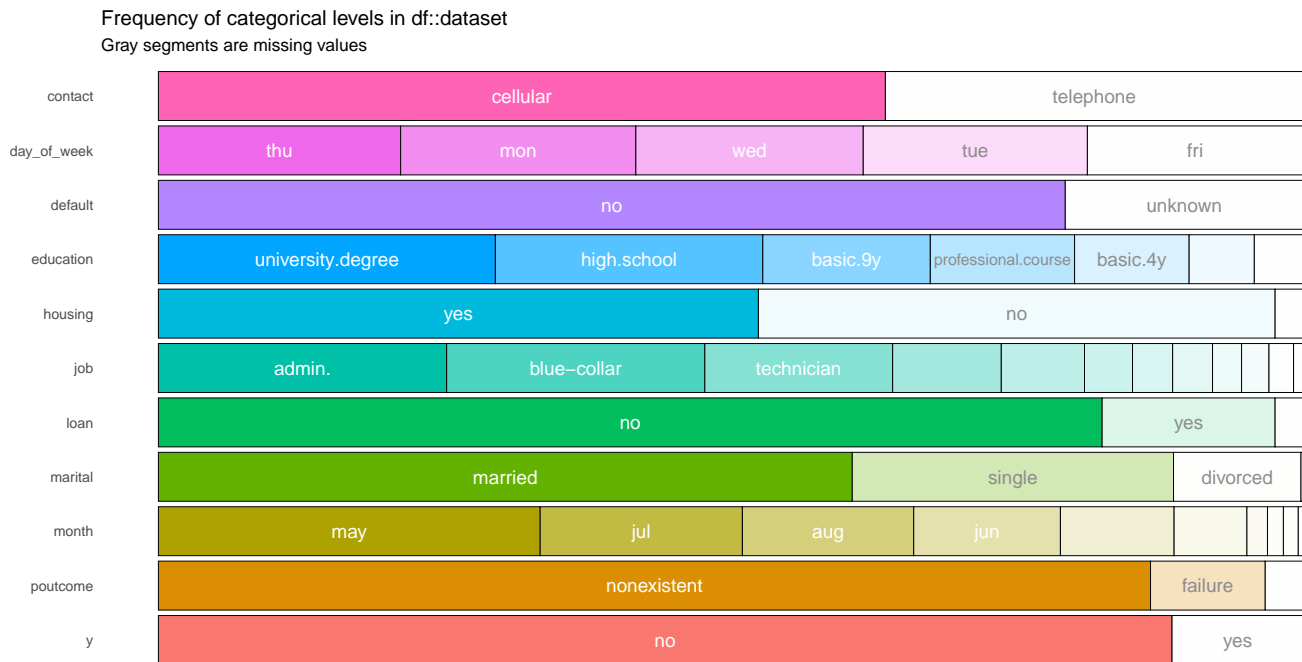
As we can see in the plot, we have 11 factors and 10 numeric values. Furthermore, we have seen that there our dataset is complete which means that it has no missing values. However, this does not imply that there are no outliers.

df::dataset column types

df::dataset has 21 columns.



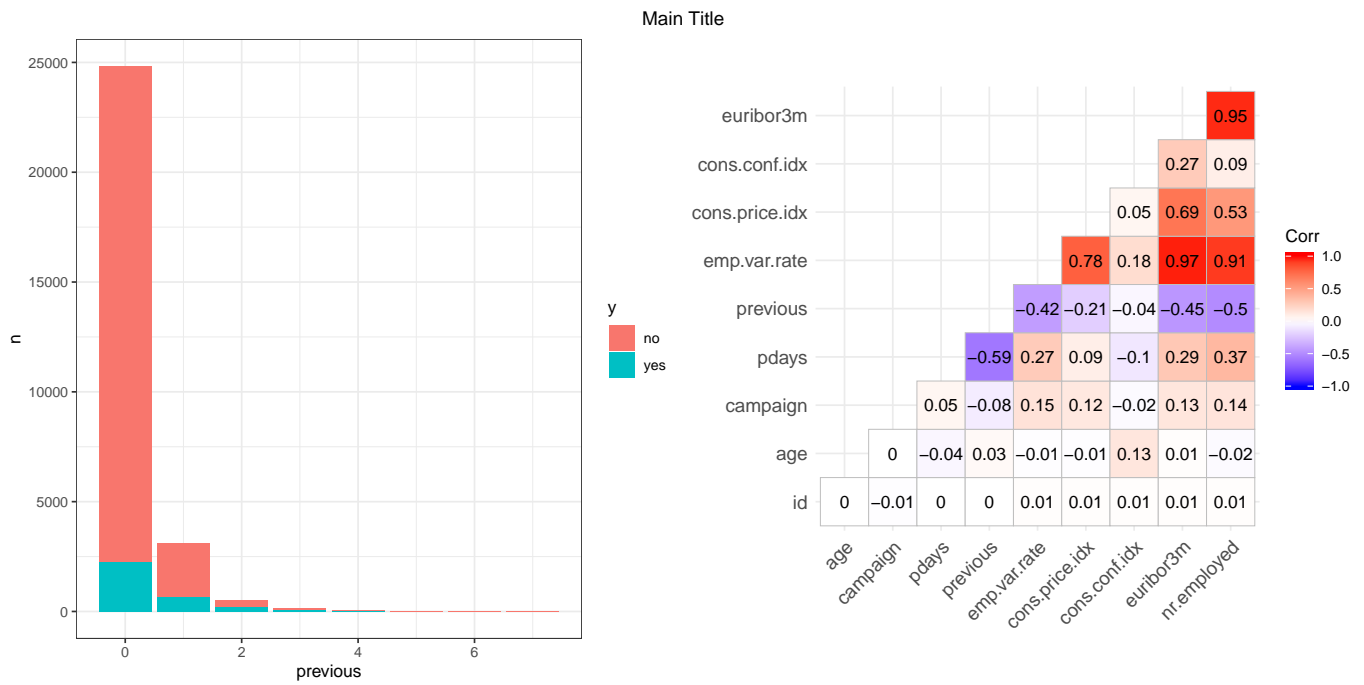
Exploration of the categorical variables where we can see a map with all their possible values and their representation. We can see that there is a clear unbalance in some variables such as *poutcome*, *load*, *contact* or the response variable *y*. So we might need to resample our dataset.



If we now explore of the numerical variables, we see that most of variables are not following a normal. This might be because they are ratios and also are taken most of them in a monthly basis.



We can also see that the correlation between the variable *previous* and *pdays* is of 0.50 which means that they are moderately correlated. But if we pay attention to number of 0's and 999 and their description, we can see that they are describing the same thing. So we are going to see how these variables are related to the response variable and as we can see below, *pdays* is not necessary. Furthermore, if we take a look at the correlation matrix, we can see that



Also we can do a general plot to have a first idea of which variables help us to distinguish between the two groups. As we can see below, it looks like *housing*, *contact* and *potcome* distinguish between the two groups.

