

Basic Operations in R

<code>abs (x)</code>	Returns the absolute value of x
<code>as.data.frame (x)</code>	Converts x (generally a matrix) into a data frame
<code>as.Date (V, "FM")</code>	Tells to R that the variable <i>V</i> contents dates in the format <i>FM</i> .
<code>as.character (V)</code>	Converts the variable <i>V</i> into characters
<code>as.factor (V)</code>	Converts the variable <i>V</i> in a factor variable
<code>as.matrix (x)</code>	Converts x (usually a table or vector) into a matrix
<code>as.numeric (L)</code>	Converts the logical expression <i>L</i> into numbers (0 or 1)
<code>as.numeric (V)</code>	Converts the variable <i>V</i> in a numerical variable
<code>as.vector (x)</code>	Converts x into a vector (x is usually a matrix)
<code>c (a, b, ... , n)</code>	Creates a vector containing the elements (numbers) a, b, ..., n
<code>c ("a", "b", ... , "n")</code>	Creates a vector containing the elements (characters) "a", "b", ..., "n"
<code>cbind (F₁, F₂, ... F_n)</code>	Combines the data frames <i>F₁</i> , <i>F₂</i> , ... <i>F_n</i> into only one data frame, taking into account that <i>F₁</i> , <i>F₂</i> , ... <i>F_n</i> have the same observations in different variables
<code>colMeans (M)</code>	Gets the mean from all the elements contained in every column from matrix <i>M</i>
<code>colnames (F)</code>	Gets the names of the columns contained in the data frame <i>F</i> (<i>F</i> can also be a matrix)
<code>colSums (M)</code>	Adds up all the elements contained in every column from matrix <i>M</i>
<code>complete (mice(F))</code>	Completes the missing values of the data frame <i>F</i> based on the values in the other variables. Requires loading the "mice" package. A split seed could also be useful.
<code>cor (F)</code>	Calculates the correlation between all the variables contained in the data frame <i>F</i>
<code>cor (F[c("V₁", " V₂", ..., " V_n ")])</code>	Calculates the correlation between the variables <i>V₁</i> , <i>V₂</i> , ..., <i>V_n</i> , all contained in the data frame <i>F</i> .
<code>cor (V₁, V₂)</code>	Calculates the correlation between the variables <i>V₁</i> and <i>V₂</i>
<code>coredata (F)</code>	Extracts the core data from the data frame <i>F</i>
<code>coredata (V)</code>	Extracts the core data from the variable <i>V</i>
<code>data.frame (V₁, V₂, ..., V_n)</code>	Creates a data frame combining the vectors <i>V₁</i> , <i>V₂</i> , ..., <i>V_n</i>
<code>diag (M)</code>	Gets the diagonal elements from a table / matrix
<code>dim (x)</code>	Gets the dimension from the object x
<code>dim (x) = c (nr, nc)</code>	Changes the dimension of the object x, resulting in an object with nr rows and nc columns
<code>F [, !(names(F) %in% X)]</code>	Removes all the variables whose names are contained in the vector <i>X</i> from frame <i>F</i>
<code>F [c("V₁", "V₂", ..., "V_n")]</code>	Shows only the selected variables <i>V₁</i> , <i>V₂</i> , ..., <i>V_n</i> of the data frame <i>F</i>
<code>F [n,]</code>	Gets the element n from the data frame <i>F</i>
<code>F [,n]</code>	Gets the column n from the data frame <i>F</i>
<code>F [sample(nrow(F), n),]</code>	Takes a random sample of n elements from the data frame <i>F</i>
<code>F\$V_r</code>	Accesses to the variable <i>V_r</i> contained in the data frame <i>F</i>
<code>getwd ()</code>	Gets the current working directory
<code>grepl("W", V, fixed=TRUE)</code>	If the word <i>W</i> appears in the variable <i>V</i> , returns "TRUE". If not, returns "FALSE"
<code>head (x)</code>	Shows the six first elements of x
<code>image (M, axes=FALSE)</code>	Plots an image from the matrix <i>M</i> , using default colours
<code>is.na (V)</code>	Returns a vector of TRUE/FALSE for if the variable <i>V</i> is missing.
<code>jitter (x)</code>	Adds or subtracts a small amount of random noise to x
<code>length (x)</code>	Provides the length of the element x
<code>ls ()</code>	Shows a list with all the variables created in the current R session
<code>M\$residuals</code>	Gets the residuals from the model <i>M</i>
<code>M₁ * M₂</code>	Multiplicates <i>M₁</i> and <i>M₂</i> , element by element (escalar product)
<code>match ("Name", V)</code>	Searches where the word "Name" is in the variable <i>V</i>

<code>matrix (V, byrow=TRUE, nrow=n)</code>	Creates a matrix from the vector V , dividing its elements into n rows
<code>mean (V)</code>	Calculates the mean of the variable V
<code>min (V)</code>	Returns the minimum value of the variable V
<code>months (V)</code>	Extracts the months from the variable V , which has a date format
<code>names (F)</code>	Gets the names of all the variables contained in the data frame F
<code>na.omit (F)</code>	Removes observations with any missing value (or "NA" value) from the data frame F
<code>na.rm = TRUE</code>	Additional argument that removes the "NA" values from V . It can be used with functions like mean, sd, max, min, tapply, etc.
<code>nchar (x)</code>	Counts the number of characters contained in x
<code>ncol (F)</code>	Gets the number of variables (or columns) contained in the data frame F
<code>nrow (F)</code>	Gets the number of observations (or rows) contained in the data frame F
<code>paste0 ("A", V₁, V₂, "B", ...)</code>	Pastes the text A with the text contained in the variables V_1 and V_2 , and with the text B, without leaving spaces between any items (Output: AV ₁ V ₂ B)
<code>paste ("A", V₁, V₂, "B", ...)</code>	Does the same as paste0, but leaving one space between the items (Output: A V ₁ V ₂ B).
<code>rbind (F₁, F₂, ... F_n)</code>	Combines the data frames $F_1, F_2, \dots F_n$ into only one data frame, taking into account that $F_1, F_2, \dots F_n$ have different observations with the same variables
<code>read.csv ("File.csv")</code>	Reads the .csv file "File" into R
<code>relevel (V, "Name")</code>	Sets the value "Name" as the reference of the unordered factor variable V
<code>rm (V)</code>	Removes the variable V
<code>rowMeans (M)</code>	Gets the mean from all the elements contained in every row from matrix M
<code>rowSums (M)</code>	Adds up all the elements contained in every row from matrix M
<code>sd (V)</code>	Calculates the standard deviation of the variable V
<code>seq (a, b, c)</code>	Creates a vector with a sequence of numbers from a to b with increments of c
<code>sign (n)</code>	Returns 1 if n is a positive number, -1 if n is a negative number, and 0 if n is zero.
<code>sort (table(V))</code>	Sorts a table which counts the number of observations in the variable V by the number of observations in that variable
<code>sort (V)[k]</code>	Gets the value of the element k contained in the variable V , sorted from the element with the least value (k = 1) to the element with the highest value (k = n), where n is the total number of observations contained in the data set.
<code>sqrt (x)</code>	Returns the square root of x
<code>str (F)</code>	Gets the structure from the data frame F
<code>subset (F, C)</code>	Gets a subset from the data frame F with the elements satisfying condition C
<code>sum(diag(table(T)))/nrow(F)</code>	If T is a classification matrix from a logistic / classification tree model, and F is a data frame corresponding, this command gets directly the accuracy
<code>sum (x)</code>	Gets the sum from all the elements contained in x (x could be a vector or a table/matrix)
<code>summary (F)</code>	Gets the summary from the data frame F
<code>summary (M)</code>	Gets the summary from the model M
<code>summary (V)</code>	Gets the summary from the variable V
<code>table (V)</code>	Counts the number of observations in each category in the variable V
<code>table (V₁, V₂, ..., V_n)</code>	Counts the number of observations in each category in the variables V_1, V_2, \dots, V_n
<code>tail (x)</code>	Shows the six last elements of x
<code>tapply (V₁, V₂, O)</code>	Splits the data by the variable V_2 , and applies the operation O to the variable V_1
<code>V = NULL</code>	Removes the variable V
<code>V[n]</code>	Gets the element n from the vector V
<code>weekdays (V)</code>	Extracts the weekdays from the variable V , which has a date format
<code>which (V == "Text")</code>	Finds the element in a data frame whose variable V is exactly "Text"
<code>which.max (V)</code>	Finds the element in a data frame whose variable V has the maximum value
<code>which.min (V)</code>	Finds the element in a data frame whose variable V has the minimum value

<code>write.csv (F, "File.csv")</code>	Saves the data frame <i>F</i> in the .csv file called "File.csv"
<code>%in%</code>	In operator
<code>&</code>	And operator
<code>?function</code>	Opens the R help and shows the features of the function
<code> </code>	Or operator

Linear Regressions in R

Creating a one variable linear regression from the data frame F , with VD as the dependent variable and VI as the independent variable

```
lm (VD ~ VI, data=F)
```

Creating a multivariable linear regression from the data frame F , with VD as the dependent variable and VI_1, VI_2, \dots, VI_n as independent variables

```
lm (VD ~ VI1 + VI2 + ... + VIn, data=F)
```

Building a linear model from the data frame F to predict the dependent variable, VD , using all the other variables in the data frame.

```
lm (VD ~ ., data=F)
```

Getting the summary from the model M

```
summary (M)
```

Computing the SSE of the model M

```
sum ((M$residuals)^2)
```

Making predictions using the model M in the training set F

```
predict (M)
```

Making predictions using the model M in the new data frame or testing set $TestF$

```
predict (M, newdata=TestF)
```

Computing the SST, where F is the training set, $TestF$ is the name of the testing set and V is the dependent variable we are trying to predict

```
sum ((TestF$V - mean(F$V))^2)
```

Computing the out-of-sample SSE of the model M , where $TestF$ is the name of the testing set, V is the dependent variable and P is the vector where the predictions are stored

```
sum ((TestF$V - P)^2)
```

Computing the out-of-sample R^2 of the model M , taking into account that SSE refers to out-of-sample SSE

```
1 - SSE/SST
```

Creating automatically a model with a good compromise of model simplicity and R^2 from the model M , which has been built previously

```
step (M)
```

Splitting Randomly a dataset in R (Continuous problems)

Sets a seed with the number n . If we split randomly the same data set two or more with the same seed n , the same splitting will be got all the times

```
set.seed (n)
```

Splits the data frame F in two data sets, leaving a proportion of p data in the first data set, and a proportion of $(1-p)$ data in the second set.

```
sample (1:nrow(F), size=p*nrow(F))
```

Creates the first dataset (from the data set F), given the split is called S

```
F[S,]
```

Creates the second dataset (from the data set F), given the split is called S

```
F[-S,]
```

Splitting Randomly a dataset in R (All problems)

Loads the library `caTools`

```
library (caTools)
```

Sets a seed with the number n . If we split randomly the same data set two or more with the same seed n , the same splitting will be got all the times

```
set.seed (n)
```

Splits the outcome variable V in two data sets, leaving a proportion of p data in the first data set, and a proportion of $(1-p)$ data in the second set, making sure that the outcome variable is well-balanced in each set

```
sample.split (V, SplitRatio = p)
```

Creates the first dataset (from the data set F), given the split is called S

```
subset (F, S == TRUE)
```

Creates the second dataset (from the data set F), given the split is called S

```
subset (F, S == FALSE)
```

Creating Logistic Regressions in R

Creating a one variable logistic regression from the data frame F , with VD as the dependent variable and VI as the independent variable

Creating a multivariable logistic regression from the data frame F , with VD as the dependent variable and VI_1, VI_2, \dots, VI_n as independent variables

Creating a logistic model from the data frame F to predict the dependent variable, VD , using all the other variables in the data frame except V_1, V_2, \dots, V_n .

Creating a logistic model from the data frame F to predict the dependent variable, VD , using all the other variables in the data frame.

Getting the summary from the model M

Making predictions of probabilities using the model M in the training set F

Making predictions using the model M in the new data frame or testing set $TestF$

```
glm (VD ~ VI, data=F, family=binomial)
```

```
glm (VD ~ VI1 + VI2 + ... + VIn, data=F, family=binomial)
```

```
glm (VD ~ ., data=F, family=binomial)
```

```
glm (VD ~ . - V1 - V2 - ... - Vn, data=F, family=binomial)
```

```
summary (M)
```

```
predict (M, type="response")
```

```
predict (M, type="response", newdata=TestF)
```

Computing ROC and AUC in R

Loads the library "ROCR"

Creates a prediction taking into account the predictions (P) made by the model and the actual values of the variable V

Does the same than the previous command, but with classification trees. The `[,2]` means that the prediction function has to take the second column

Plots the ROC curve. PR is the output of the prediction function. The other arguments are optional, and are used to colour the plot and to add values of thresholds to it.

Computes the AUC from the prediction made in the previous step, called $P2$

```
library (ROCR)
```

```
prediction (P, V)
```

```
prediction (P[,2], V)
```

```
plot (performance(PR, "tpr", "fpr"), colorize = TRUE,  
      print.cutoffs.at = seq(0,1,0.1), text.adj = c(-0.2,1.7))
```

```
AUC = as.numeric(performance(P2, "auc")@y.values)
```

Cross Validation Method in R (CART)

Loads the library "caret"

Loads the library "e1071"

Defines the number of folds that we want, in this case, n

Defines the cp parameters from a to b in increments of c

Gets the optimal value of the cp parameter. F is the name of the training set, C is the output of the `trainControl` function and G is the output of the `expand.grid` function

```
library (caret)
```

```
library (e1071)
```

```
trainControl (method="cv", number = n)
```

```
expand.grid (.cp=seq(a, b, c))
```

```
train (VD ~ VI1 + VI2 + ... + VIn, data=F,  
      method="rpart", trControl=C, tuneGrid=G)
```

CART Models in R

Loads the library "rpart"

Loads the library "rpart.plot"

Creates a classification tree from the data frame F , with VD as the dependent variable and VI_1, VI_2, \dots, VI_n as independent variables. If no additional arguments are added, R will built the tree by default. Some additional arguments that can be added are the following:

Sets that in each split must be a minimum of n points

Sets the cp (complexity parameter) value equal to k

Sets the matrix M as a penalty matrix

Creates a regression tree from the data frame F , with VD as the dependent variable and VI_1, VI_2, \dots, VI_n as independent variables, with a complexity parameter of n . If no additional arguments are added, R will built the tree by default. Some additional arguments that can be added are the following:

Sets the cp (complexity parameter) value equal to k

Plots the tree T , with n significant digits in each split

Loads the library "randomForest"

Converts the variable we are trying to predict, V , into a factor.

Creates a random forest from the data frame F , with VD as the dependent variable and VI_1, VI_2, \dots, VI_n as independent variables, with a minimum of n points in every split and a number m of trees.

Makes predictions in unseen data $TestF$, in classification problems, based on the tree model T and using a threshold of 0.5. If the argument "newdata" is removed, the predictions will be made in the training set

In classification problems, gives the probabilities of every of the possible outputs happening. If the argument "newdata" is removed, the predictions will be made in the training set

Makes predictions in unseen data $TestF$, in regression problems, based on the tree model T . If the argument "newdata" is removed, the predictions will be made in the training set

Makes predictions in unseen data $TestF$, in classification problems, based on the random forest model R . If the argument "newdata" is removed, the predictions will be made in the training set. Adding the argument $type="prob"$ is also optional. If it's not added, the output will be a prediction; if it's added, the output will be the probability of something happening.

```
library (rpart)
```

```
library (rpart.plot)
```

```
rpart ( $VD \sim VI_1 + VI_2 + \dots + VI_n$ , data= $F$ , method="class")
```

```
minbucket = n
```

```
cp = k
```

```
parms=list(loss= $M$ )
```

```
rpart ( $VD \sim VI_1 + VI_2 + \dots + VI_n$ , data= $F$ )
```

```
cp = k
```

```
prp ( $T$ , digits = n)
```

```
library (randomForest)
```

```
as.factor ( $V$ )
```

```
randomForest ( $VD \sim VI_1 + VI_2 + \dots + VI_n$ , data= $F$ ,  
method="class", nodesize=n, ntrees=m)
```

```
predict ( $T$ , newdata= $TestF$ , type="class")
```

```
predict ( $T$ , newdata= $TestF$ )
```

```
predict ( $T$ , newdata= $TestF$ )
```

```
predict ( $R$ , newdata= $TestF$ , type="prob")
```

Interpreting Random Forests

Plots a chart for each variable measuring the number of times that variable was selected for splitting. M is the name of the model, VU is the output from the first line, and VUS is the output from the second line

Plots the reduction in impurity of each variable in the model M

```
varUsed (M, count = TRUE)
sort (VU, decreasing = FALSE, indexreturn = TRUE)
dotchart (VUS$x, names(M$forest$levels[VUS$ix]))
```

```
varImpPlot (M)
```

Text Analytics in R

Sets the language to default (English)

Reads the .csv file "File" into R, keeping the text in a way that can be analyzed

Loads the library "tm"

Loads the library "SnowballC"

Converts the text contained in the variable V into a corpus (a corpus is a collection of documents)

Gets the text contained in the element n of the corpus C

Changes all the text contained in the corpus C to lowercase.

Removes all the punctuation signs contained in the corpus C

Gets a list with the english stopwords.

Removes all english stopwords from the corpus C

Removes all english stopwords and the words " W_1 ", " W_2 ", ..., " W_n " from the corpus C

Stems all the words from the corpus C (stemming a word is getting the root from it)

Creates a document term matrix counting the number of times that every word contained in the corpus C appears on it

Shows the matrix created with the previous command, called M

Finds the terms appearing at least n times in the matrix M (generated from the corpus C)

Removes the terms appearing in the matrix M in a proportion lower than $(1 - p)$

Converts the matrix M into a standard matrix, and then, into a data frame

If the data frame F got from the previous step has any variable with its name starting with a number, fixes it

```
Sys.setlocale ("LC_ALL", "C")
read.csv ("File.csv", stringsAsFactors=FALSE)
```

```
library (tm)
library (SnowballC)
Corpus (VectorSource (V))
```

```
C[[n]]$content
```

```
tm_map (C, content_transformer(tolower))
```

```
tm_map (C, removePunctuation)
```

```
stopwords ("english")
tm_map (C, removeWords, stopwords("english"))
tm_map (C, removeWords, c("W1", "W2", ..., "Wn",
stopwords("english")))
tm_map (C, stemDocument)
```

```
DocumentTermMatrix (C)
```

```
inspect (M)
```

```
findFreqTerms (M, lowfreq=n)
```

```
removeSparseTerms (M, p)
```

```
as.data.frame (as.matrix(M))
```

```
make.names (colnames(F))
```

Hierarchical clustering in R

Gets the euclidean distance in every variable contained in the data frame (or vector, but in this case there is only one variable) F . (Be careful! In large data sets this command could collapse the computer)

`dist (F, method = "euclidean")`

Computes the hierarchical clustering inputting the distances D , where D is the output got from the previous command

`hclust (D, method = "ward.D")`

Plots the dendrogram of the clustering C , where C is the output of the hclust function

`plot (C)`

In the graph plotted with the previous command, shows in "col" colour the result of dividing the data set by n clusters

`rect.hclust (C, k=n, border="col")`

Divides the observations into n different clusters.

`cutree (C, k=n)`

Runs the second argument (colMeans) on each element of the first argument. The result is finding the centroids of every cluster contained in the group of clusters G (G is the output of the cutree function) from every variable from variable a to variable b (a and b are numbers) contained in the data frame F

`lapply (split (F, G), colMeans)`

Gets in which cluster of all the clusters contained in the group G is the observation n

`G [n]`

Gets all the observations of the variable V contained in the cluster c

`c$V`

K-means clustering in R

Loads the library "flexclust"

`library (flexclust)`

Computes the K-means clustering inputting the vector or data frame *V*, dividing the data into *k* clusters and doing no more than *n* iterations. The output of this command is *K*

`kmeans (V, centers = k, iter.max = n)`

Gets in what cluster is placed every observation from the vector/data frame *V*

K\$clusters

Gets the centroids of the clusters

K\$centers

Gets the size of each cluster

K\$size

Converts the information got from the clustering information, *K* to an object of the class KCCA. *V* is the vector or data frame used in the kmeans function. This step is needed before using the predict function, and may take a long time. The output of this step will be *K.kcca*

`as.kcca (K, V)`

Using *K.kcca*, predicts in which cluster will be located every observation from the new data frame/vector, called *NV*

`predict (K.kcca, newdata = NV)`

Normalising data

Loads the library "caret"

`library (caret)`

Gets the mean and the standard deviation from the object *x*. Then, normalises the object *y* with the features of the object *x*. If we want to normalise an object with its own features, then the object *y* would be the same as the object *x*

`predict (preProcess (x), y)`

Plots & Data Visualization in R

<code>abline (h=a)</code>	In a plot, draws an horizontal line in $y=a$.
<code>abline (v=a)</code>	In a plot, draws a vertical line in $x=a$.
<code>abline (v=as.Date(c("Dt")))</code>	In a plot whose x-axis are dates, draws a vertical line in the point where date = "Dt"
<code>boxplot ($V_1 \sim V_2$)</code>	Creates a boxplot with the variable V_1 sorted by the variable V_2 . V_2 is not always needed
<code>hist (V)</code>	Creates an histogram of the variable V
<code>lines (x, y)</code>	Adds a line to the line plot created previously with the variables x and y
<code>points(x[C], y[C], col="col", pch=19)</code>	In an existing plot, colours with the colour "col" the subset of points fulfilling the condition C
<code>plot (x, y)</code>	Creates a basic plot with x as the independent variable, and y as the dependent one
<code>plot (x[a:b], y[a:b])</code>	Creates a basic plot, but only with the observations ranged between a and b

The following section refers to additional arguments that can be added in most of the types of plots listed above:

<code>breaks = N</code>	In an histogram, divides the entire range of it in N pieces. The width of every piece is calculated by $Width = \frac{\max(V) - \min(V)}{N}$
<code>col = "col"</code>	Changes the colour of the plot to the colour "col". Examples: red, blue, green, yellow...
<code>lwd = w</code>	Changes the width of a line from default ($w = 1$) into w.
<code>main = "Title of the plot"</code>	Names the plot with the name "Title of the plot"
<code>type = "l"</code>	In a x-y plot, draws the output as a line
<code>xlab = "Axis x name"</code>	Labels the x-axis with the name "Axis x name"
<code>xlab = ""</code>	Labels the x-axis as default
<code>xlim = c (a, b)</code>	Limits the values on the x-axis from a to b.
<code>ylab = "Axis y name"</code>	Labels the y-axis with the name "Axis y name"
<code>ylab = ""</code>	Labels the y-axis as default

Other functions in R

`image (M, axes = FALSE, col = grey(seq(0, 1, length = 256)))`

`image (M, axes=FALSE, col = rainbow(k))`

`lag (zoo(V, -n, na.pad = TRUE)`

`merge (F1, F2, by.x="NameV1", by.y="NameV2", all.x=TRUE)`

`read.csv ("File.csv", header = FALSE)`
the

data frame hasn't got a header (or variable name)

`read.table ("File.csv", header = FALSE, sep = "S", quote = "\\")`

Plots an image from the matrix *M*, using the gray colour convention

Plots an image from the matrix *M*, using *k* different colours from the rainbow palette

Lags the variable *V* a number of *n* observations, adding NA values to the first *n* observations (needs loading the library "zoo" previously.

Joins two data frames, *F1* and *F2*, matching the variable called *NameV1* from *F1* data frame with the variable called *Name2* from *F2* data frame, keeping all rows from the data frame *F1*, even if some of the rows from *V1* doesn't match any row in *F2*.

Reads the .csv file "File" into R, but telling R that

row

Reads the file "File.txt", but telling R that the data frame hasn't got a header (or variable name) row and making sure that the text is read in properly