

Intro to Programming and Basic Analysis in R

Dr. Iulia Cioroianu

ESS 1B 2021

TOPIC 1. Introduction, installing and setting up R and RStudio

<https://cran.r-project.org/>

<https://www.rstudio.com/products/rstudio/download/>

<https://rstudio.cloud/>

In this class:

We will be using RStudio Cloud, which allows you to work on your own copy of each lab project.

For each lab, the Project will contain: * the .R file with the code that we'll go over together, as well as short exercises that you are encouraged to do together * the .Rmd file which contained the same code, with explanations ready to be knitted into.. * the .pdf file with code, explanations and results * any required data

Beyond this class, if you don't have it already, downloading and installing R onto your computer is easy:

Depending on your operating system, follow the instruction below.

a. Windows

Visit <https://cran.r-project.org/bin/windows/base/> and click on "Download R 3.5.1 for Windows" in the upper left corner. Save the R-3.5.1.exe installer in your Downloads directory (or any other location of your choice on your computer). Once the download has finished, double click the file. If asked if you want to run the software and/or if you are allowing the app to make changes to your computer click "Run" and/or "Yes". Select the language to be used in the installation (English is the default). Read the information provided and keep clicking "Next" on the following screens to keep all the default settings. Wait for the installation to finish. Proceed to the "Install RStudio" section below.

b. Mac OS X

Visit <https://cran.r-project.org/bin/macosx/> and click on "R-3.5.1.pkg" on the left side of the page. Wait for the installer to download, then click on it and click the "Continue" button until prompted to agree to the terms and conditions. Read the terms and conditions and click the "Agree" button. Then click the "Install" button. If prompted to log in to complete the installation, insert your computer username and password, and click "Install Software". When installation finishes, proceed to the "Install RStudio" section below.

c. Linux

Use the following command to install R-Base:

```
sudo apt-get install r-base
```

typing "y" if prompted along the way. Wait for the installer to finish and proceed to the Install RStudio section below.

Starting an R session

Windows: Rgui.exe.

MAC: R.app.

UNIX: you enter the command R (the path must include the R binaries).

Writing scripts

You can write anything in the console, but it is much more convenient to write a script that you can then run (in chunks or in full).

> Script/program - collection of expressions to be evaluated sequentially.

You can write scripts in the R GUI, but it is much more convenient to use a text editor such as Emacs, VIM, Eclipse, OR...

The amazing IDE (integrated development environment) for R: **RStudio**.

Downloading and installing RStudio:

Visit <https://www.rstudio.com/products/rstudio/download/#download> and click the installer that corresponds to your operating system to download RStudio Desktop. Wait for the installer to download, then click on it to proceed with the installation, clicking “Next”/“Keep” and “Install” along the way if prompted. Wait for the installer to finish.

Testing

You should now have both R and RStudio installed among your programs/applications. To test if the installation was successful, start RStudio by either clicking on the shortcut created or searching for it in your programs/applications list/search box.

In the “Console” window, type:

```
1+sqrt(4)
```

```
## [1] 3
```

And click enter. R will compute and display the answer for you. In this case, 3.

Setting the working directory

What is the current working directory?

```
getwd()
```

```
## [1] "/cloud/project"
```

Change the working directory to the folder we created at the beginning:

```
#setwd("Newdir")
```

Want to set back to one level up?

```
setwd('..')
```

What is in the working directory?

```
list.files()
```

```
## [1] "Data_analysis_R.R" "example_data.csv" "Intro_to_r.pdf"
## [4] "Intro_to_r.Rmd"    "Newdir"           "project.Rproj"
```

Very useful when you're not in RStudio Cloud - list the folders and sub-folders:

```
# list.dirs()
```

Getting help

In RStudio: the Help panel in the lower-right window. Try it out! In general, in R, find out more about an R command or function `x`, `help(x)` or `?x`.

```
help(median)
```

```
?lm
```

Search the titles, names, aliases, and keyword entries of the available help files for a keyword or phrase:

```
help.search("regular expression")
```

```
help.start() #HTML help interface
```

```
## starting httpd help server ... done
```

```
## If the browser launched by 'xdg-open' is already running, it is *not*
```

```
## restarted, and you must switch to its window.
```

```
## Otherwise, be patient ...
```

TOPIC 2: Basic computing and programming in R

This section is based on “Introduction to Scientific Programming and Simulation Using R” by Jones, Millardet and Robinson, which is a great resource if you are serious about using R in your work. Available at: <https://www.crcpress.com/Introduction-to-Scientific-Programming-and-Simulation-Using-R-Second-Edition/Jones-Maillardet-Robinson/9781466569997>

Arithmetic operations

```
((1 + 2/3)*3)^2
```

```
## [1] 25
```

```
sqrt(25) # square root
```

```
## [1] 5
```

```
exp(1) #exponential
```

```
## [1] 2.718282
```

```
options(digits = 20) #set number of digits to display
```

```
exp(1)
```

```
## [1] 2.7182818284590450908
```

```
options(digits = 6)
```

```
log(2.7183) #logarithm
```

```
## [1] 1.00001
```

Functions

Multiple built-in functions were used above. A function takes one or more arguments (or inputs) and produces one or more outputs (or return values). You write the name of the function followed by the argument in round brackets.

```
pi    #itself a function in R, without an argument
```

```
## [1] 3.14159
```

```
sin(pi/6)    #sine function
```

```
## [1] 0.5
```

```
round(pi)    #round to nearest integer
```

```
## [1] 3
```

```
seq(from = 1, to = 7, by = 2)    #produces sequences of numbers from 1 to 9, every 2
```

```
## [1] 1 3 5 7
```

```
seq(from = 1, to = 7)    #the default value for the third argument is 1
```

```
## [1] 1 2 3 4 5 6 7
```

If you do not specify the names of the arguments (from, to, by) then the order you list them in in the function is assumed to be the default one.

```
seq(1, 9, 3)
```

```
## [1] 1 4 7
```

```
?seq
```

Variables

"Expression" - a phrase of code that can be executed. "Assignment" - saving the evaluation of an expression, assigning it to a variable.

```
a=1    #variables are created the first time you assign a value to them.
```

```
a    #display the value of variable a
```

```
## [1] 1
```

```
A <- 2    #names are case sensitive.
```

```
print(A)    #display the value of variable A
```

```
## [1] 2
```

```
a+A/10    #evaluate value of this expression
```

```
## [1] 1.2
```

```
a <- a+1    #you can have a variable both on the left and the right
```

```
a
```

```
## [1] 2
```

Using <- or = ?

Can use either = or <- for assignment. They do the same thing. Some prefer one, some prefer the other. In this tutorial I will use <- because it is more widely used. I prefer = because it is easier to use and similar to other programming languages that I use. To get <- in RStudio, use Alt- as a shortcut. Note that it includes spaces around it. It is good practice to have these spaces, although they are not (always) necessary. It is also good practice to give informative names to your variables.

Vectors

A vector (or array) is an indexed list of variables. Created the first time a value is assigned to them, just like variables.

A variable is a vector with length 1.

Vectors are created by using functions such as `c()`-combine, `seq(from, to, by)`-sequence, `rep(x, times)`-repetition. Let's create some vectors. Check the "Environment" panel in the left to see them.

```
v1 <- c(1, 2, 3, 4)
```

```
v2 <- seq(1,7,2)
```

```
v3 <- rep(2,3)
```

```
v4 <- c(v2,v3) #we can combine vectors as well
```

To refer to element *i* of vector *x*, we use *x*[*i*].

```
v1[1]
```

```
## [1] 1
```

```
v2[2:3] #Get element f "from:to"
```

```
## [1] 3 5
```

```
length(v3) #the number of elements of vector v3
```

```
## [1] 3
```

Algebraic operations on each element separately (elementwise):

```
v1+v2
```

```
## [1] 2 5 8 11
```

```
v1*v2
```

```
## [1] 1 6 15 28
```

```
v1*v3
```

```
## Warning in v1 * v3: longer object length is not a multiple of shorter object  
## length
```

```
## [1] 2 4 6 8
```

Note the warning message. The two vectors need to either be the same length or one be a multiple of the other. If not, R duplicates the shorter one, but produces the warning.

```
2+v3 #So if one of them is of length 1, this works.
```

```
## [1] 4 4 4
```

Applying functions to vectors

Note that some functions applied to vectors act on the whole vector:

```
sum(v1)
```

```
## [1] 10
```

```
mean(v1)
```

```
## [1] 2.5
```

```
var(v1)
```

```
## [1] 1.66667
```

```
prod(v1)
```

```
## [1] 24
```

```
min(v1)
```

```
## [1] 1
```

```
max(v1)
```

```
## [1] 4
```

While other functions act elementwise:

```
sqrt(v1)
```

```
## [1] 1.00000 1.41421 1.73205 2.00000
```

```
sort(v1, decreasing=TRUE)
```

```
## [1] 4 3 2 1
```

Logical expressions: True or False

Comparison operators <, >, <=, >=, == (equal to), and != (not equal to). Logical operators & (and), | (or), and ! (not).

```
v1==v2 #0 and 1 can be used to mean true or false
```

```
## [1] TRUE FALSE FALSE FALSE
```

```
c(0,1,1) | c(1,1,0) #True if either one of them true
```

```
## [1] TRUE TRUE TRUE
```

```
xor(c(0,1,1), c(1,1,0)) #exclusive disjunction -
```

```
## [1] TRUE FALSE TRUE
```

```
#either one or the other but not both
```

Selecting a subset of a vector for which a certain condition is true:

```
v1 > 2
```

```
## [1] FALSE FALSE TRUE TRUE
```

```
v1[v1 > 2]
```

```
## [1] 3 4
```

Finding the position of elements which meet a certain condition:

```
which(v4>=5)
```

```
## [1] 3 4
```

Missing data

```
m <- NA # assign NA to variable a
```

```
is.na(m)    #is it missing?
```

```
## [1] TRUE
```

```
m <- c(1,NA,3)    #NA in a vector
```

```
is.na(m)    #is it missing?
```

```
## [1] FALSE TRUE FALSE
```

```
mean(m)     # NAs can propagate
```

```
## [1] NA
```

```
mean(m, na.rm = TRUE)    # NAs can be removed
```

```
## [1] 2
```

```
length(m[!is.na(m)])    #how many non-missing elements are in m?
```

```
## [1] 2
```

Matrices

Created from a vector using the function `matrix(data, nrow = 1, ncol = 1, byrow = FALSE)`. If `length(data)<nrow*ncol`, then data is reused as many times as is needed.

```
matrix(0, 2, 2)
```

```
##      [,1] [,2]
```

```
## [1,]    0    0
```

```
## [2,]    0    0
```

```
matrix(c(1,2), 2, 2)
```

```
##      [,1] [,2]
```

```
## [1,]    1    1
```

```
## [2,]    2    2
```

```
matrix(c(1:6), 2, 3)
```

```
##      [,1] [,2] [,3]
```

```
## [1,]    1    3    5
```

```
## [2,]    2    4    6
```

```
matrix(1:6, 2, 3, TRUE)
```

```
##      [,1] [,2] [,3]
```

```
## [1,]    1    2    3
```

```
## [2,]    4    5    6
```

To create a diagonal matrix:

```
diag(1:3)
```

```
##      [,1] [,2] [,3]
```

```
## [1,]    1    0    0
```

```
## [2,]    0    2    0
```

```
## [3,]    0    0    3
```

We refer to the elements of a matrix using two indices: the first one for row, the second for column

```
A <- matrix(1:6, 2, 3) #create a new matrix
```

```
A
```

```
##      [,1] [,2] [,3]
## [1,]    1    3    5
## [2,]    2    4    6
```

```
A[1, 3] <- 0 #assign the value 0 to the element on row 1, column 3
```

```
A
```

```
##      [,1] [,2] [,3]
## [1,]    1    3    0
## [2,]    2    4    6
```

```
A[, 2:3] #display all rows, but only columns from 2 to 3.
```

```
##      [,1] [,2]
## [1,]    3    0
## [2,]    4    6
```

```
A[1,]
```

```
## [1] 1 3 0
```

Get the number of rows and columns of the matrix:

```
nrow(A)
```

```
## [1] 2
```

```
ncol(A)
```

```
## [1] 3
```

To join matrices with rows of the same length (stacking vertically):

```
rbind(A, c(1:3))
```

```
##      [,1] [,2] [,3]
## [1,]    1    3    0
## [2,]    2    4    6
## [3,]    1    2    3
```

To join matrices with columns of the same length (stacking horizontally) use cbind(...).

```
cbind(A, c(1:2))
```

```
##      [,1] [,2] [,3] [,4]
## [1,]    1    3    0    1
## [2,]    2    4    6    2
```

Algebraic operations act elementwise on matrices:

```
A+A # addition
```

```
##      [,1] [,2] [,3]
## [1,]    2    6    0
## [2,]    4    8   12
```

```
A*A #element-wise product
```

```
##      [,1] [,2] [,3]
```



```
## [1,] 1 9 0
## [2,] 4 16 36
```

For matrix multiplication, use %*%:

```
A%*%t(A) # multiply A by its transpose
```

```
##      [,1] [,2]
## [1,] 10 14
## [2,] 14 56
```

```
solve(A%*%t(A)) #get the inverse of the above square matrix
```

```
##      [,1] [,2]
## [1,] 0.1538462 -0.0384615
## [2,] -0.0384615 0.0274725
```

From vector to matrix and the other way around

```
is.vector(A) #test if A is a vector
```

```
## [1] FALSE
```

```
is.matrix(A) #test if A is a matrix
```

```
## [1] TRUE
```

```
B <- as.matrix(v1) # create a matrix of 1 column out of a vector
```

```
B
```

```
##      [,1]
## [1,] 1
## [2,] 2
## [3,] 3
## [4,] 4
```

```
is.matrix(B)
```

```
## [1] TRUE
```

```
x <- as.vector(A) #create a vector from the columns of a matrix (stored columnwise)
```

```
x
```

```
## [1] 1 2 3 4 0 6
```

```
?is.vector
```

Control structures

- *if*

```
if (A[1,1]==1) {
  print("First element of first column is 1")}
```

```
## [1] "First element of first column is 1"
```

- if, else statement

```
if (A[1,1]==2) {
  print("First element of first column is 2")
} else {
```

```
print("First element of first column is not 2")
}
```

```
## [1] "First element of first column is not 2"
```

- *ifelse*

```
v <- c(1,2,3,4,5,6)
ifelse(v %% 2 == 0, "even", "odd")
```

```
## [1] "odd" "even" "odd" "even" "odd" "even"
```

- *for* loop

```
for(i in v2) {
  print(i)
}
```

```
## [1] 1
## [1] 3
## [1] 5
## [1] 7
```

Note that the value of *i* gets updated in the loop.

```
print(i)
```

```
## [1] 7
```

- *while*

```
j <- 1
while (j < 5) {
  print(j)
  j = j+1
}
```

```
## [1] 1
## [1] 2
## [1] 3
## [1] 4
```

Note that it is not efficient or advisable to use loops in R.

Writing your own functions

```
x=1
```

```
y=100
```

```
result=0
```

```
myfunction <- function(x, y) {
  #function that raises x to the power y
  result <- x^y
  print(paste(x, " raised to the power ", y, " is ", result))
}
```

Note that you need to be careful when you write functions not to re-write any that already exist. For example, if we names our function above `sum` then that would replace the `sum` function in base R. If that ever happens, you can revert back to the original function with `rm(sum)` or `sum <- base::sum`.

Calling the function:

```
myfunction(2,5)
```

```
## [1] "2 raised to the power 5 is 32"
```

TOPIC 3: Working with dataframes

We have the following vectors:

```
a <- c(78, 89.5, 82, 36.5, 74, 39)    #numeric vector
b <- c("one", "two", "three", "six", "five", "four") #character vector
c <- c(TRUE, TRUE, TRUE, FALSE, TRUE, FALSE) #logical vector
d <- c("red", "blue", "yellow", NA, NA, NA) #character vector
```

Create a dataset out of the vectors above

```
mydata <- data.frame(a,b,c,d)
```

Give names to the variables (the vectors)

```
names(mydata) <- c("score", "exam_seat", "passed", "color") #variable names
```

Identify and inspect the elements of a data frame:

```
mydata #prints the data
```

```
##  score exam_seat passed  color
## 1  78.0      one   TRUE   red
## 2  89.5      two   TRUE  blue
## 3  82.0     three   TRUE yellow
## 4  36.5      six  FALSE  <NA>
## 5  74.0     five   TRUE  <NA>
## 6  39.0     four  FALSE  <NA>
```

```
#View(mydata) #opens the dataframe in a new tab. Note the capital V.
```

```
names(mydata) #print the names of the variables
```

```
## [1] "score"      "exam_seat"  "passed"     "color"
```

```
nrow(mydata) #the number of observations in the dataset (number of rows)
```

```
## [1] 6
```

```
ncol(mydata) #the number of variables in the dataset (number of columns)
```

```
## [1] 4
```

```
mydata[2:4] # prints columns 2,3,4 of data frame
```

```
##  exam_seat passed  color
## 1      one   TRUE   red
## 2      two   TRUE  blue
## 3     three   TRUE yellow
## 4      six  FALSE  <NA>
## 5     five   TRUE  <NA>
## 6     four  FALSE  <NA>
```

```
mydata[c("score", "color")] #prints columns id and color from data frame
```

```
##  score  color
```

```
## 1 78.0 red
## 2 89.5 blue
## 3 82.0 yellow
## 4 36.5 <NA>
## 5 74.0 <NA>
## 6 39.0 <NA>

mydata$passed #prints variable 'passed' from the data frame
```

```
## [1] TRUE TRUE TRUE FALSE TRUE FALSE
mydata["passed"]
```

```
## passed
## 1 TRUE
## 2 TRUE
## 3 TRUE
## 4 FALSE
## 5 TRUE
## 6 FALSE
```

Subsetting - selecting only certain observations and/or certain variables.

Select only observations with score less than 80, put them in a new object called 'newdata'

```
newdata1 <- subset(mydata, score<=80)
newdata1
```

```
## score exam_seat passed color
## 1 78.0 one TRUE red
## 4 36.5 six FALSE <NA>
## 5 74.0 five TRUE <NA>
## 6 39.0 four FALSE <NA>
```

Select only observations with score less than 80 and which passed (so passed==TRUE).

```
newdata2 <- subset(mydata, score<80 & passed==TRUE)
newdata2
```

```
## score exam_seat passed color
## 1 78 one TRUE red
## 5 74 five TRUE <NA>
```

Select only observations which passed AND keep only the variables 'score' and 'passed' in the new dataset.

```
newdata3 <- subset(mydata, passed==TRUE, select=c(score,passed))
```

```
newdata3
```

```
## score passed
## 1 78.0 TRUE
## 2 89.5 TRUE
## 3 82.0 TRUE
## 5 74.0 TRUE
```

Aggregate/collapse data

```
#Get the mean score for the group that failed and the one that passed
aggregated_data=aggregate(cbind(score)~passed, data=mydata, mean)
```

More generally, you can do: `aggregate(cbind(var1, var2, var3)~cat_var1+cat_var2, data=mydata, mean)` to get the means of `var1`, `var2` and `var3` BY the categories of categorical variables `cat_var1` and `cat_var2`. Instead of `mean` you can also get `sum`, `median`, `min`, `max`.

You can assign the result of the aggregation to a new dataframe

```
collapsed_data <- aggregate(cbind(score)~passed, data=mydata, mean)
```

Generate and modify variables

To create a new variable in a dataset you just have to declare it.

```
mydata$newvar <- NA #generate a new empty variable (full of NAs).
```

```
ncol(mydata)
```

```
## [1] 5
```

Warning!!! If you declare any variable in the dataset again, it gets overwritten without any warning. So be careful with this!

```
mydata$newvar <- 0 #the already existing variable 'newvar' is changed to all 0s.
```

It's always good practice to copy a variable that you want to make changes into a new variable, and modify the new one.

```
mydata$newvar <- mydata$score #replace the values of 'newvar' with the values of 'score'
```

Assign values to a variable, conditional on something:

```
mydata$newvar[mydata$score!=78] <- 0 #if score different from 78, newvar is 0.
```

```
mydata$newvar[mydata$score==78] <- 1 # if score is equal to 78, newvar is 1.  
 #(You can also use <, <=, >, >=)
```

```
mydata
```

```
##   score exam_seat passed  color newvar  
## 1  78.0      one   TRUE   red     1  
## 2  89.5      two   TRUE  blue     0  
## 3  82.0     three  TRUE yellow     0  
## 4  36.5      six  FALSE  <NA>     0  
## 5  74.0     five   TRUE  <NA>     0  
## 6  39.0     four  FALSE  <NA>     0
```

```
mydata$newvar[mydata$score<=80 & mydata$passed==TRUE] <- 1  
#if score<=80 AND passed is true then newvar becomes 1
```

```
mydata$newvar[is.na(mydata$color)==TRUE | mydata$score==1 | mydata$exam_seat=="three"] <- 2
```

Code above says: if color is missing OR score is 1 OR exam_seat is 'three' then replace the value of newvar with 2.

```
mydata
```

```
##   score exam_seat passed  color newvar  
## 1  78.0      one   TRUE   red     1  
## 2  89.5      two   TRUE  blue     0  
## 3  82.0     three  TRUE yellow     2  
## 4  36.5      six  FALSE  <NA>     2  
## 5  74.0     five   TRUE  <NA>     2
```

```
## 6 39.0      four FALSE  <NA>      2
```

NOTE: You only need quotation marks in the condition for string variables (variables whose values are words). Do NOT use quotation marks for number variables (whose values are either only numbers or a combination of numbers and missing data(NA)), and do NOT use quotation marks for boolean - TRUE/FALSE - variables.

```
A=TRUE
```

```
typeof(A)
```

```
## [1] "logical"
```

Create a dummy variable (a variable with only two categories)

```
mydata$newvar2 <- ifelse(is.na(mydata$color)==TRUE, 1, 0)
mydata
```

```
##   score exam_seat passed  color newvar newvar2
## 1  78.0      one   TRUE   red      1      0
## 2  89.5      two   TRUE  blue      0      0
## 3  82.0     three   TRUE yellow     2      0
## 4  36.5      six  FALSE  <NA>      2      1
## 5  74.0     five   TRUE  <NA>      2      1
## 6  39.0     four  FALSE  <NA>      2      1
```

Create a variable `newvar2` which takes the value 1 if color is not missing, and 0 if color is missing.

What type of variable is `newvar2`?

```
typeof(mydata$newvar2)
```

```
## [1] "double"
```

Change it to integer:

```
mydata$newvar2 <- as.integer(mydata$newvar2)
```

TOPIC 4: Importing and exporting data

Check out the “Import dataset” button in the Environment window. What formats are available for import? Make sure you give the imported data frame a short name. Check ‘Yes’ for heading if the variable names are on the first row.

OR: Set the working directory to be the folder where you have the data:

```
#setwd("C:/Rintro")
```

Now you can read the data directly from the folder. Everything that you save will also go to that folder.

```
data <- read.csv("example_data.csv")
```

Let’s have a look at the data.

```
#View(data)
```

Turning off scientific notation

```
options(scipen=999)
```

Export the data to a new file, called `saved_data.csv`.

```
#write.csv(data, file="saved_data.csv")
```

Exercise:

Replace the “missing” values of the variables in our dataset with NA.

```
typeof(data$reelected)
```

```
## [1] "character"
```

```
data[data=="missing"] <- NA
```

Turn variable to numeric

```
data$reelected=as.numeric(data$reelected)
```

TOPIC 5: Descriptive statistics

Get the mean, median, 25th and 75th quartiles, min, max, number of NA (if no NA given then no NAs exist for that variable).

```
summary(data$spent) # for a single variable
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##      0  172196  256307  351340  393587 3489592
```

```
summary(data) # for ALL the variables in the dataset.
```

```
##      id      region      party      chamber
## Min.   : 1  Length:539  Length:539  Length:539
## 1st Qu.:136  Class :character  Class :character  Class :character
## Median :270  Mode  :character  Mode  :character  Mode  :character
## Mean   :270
## 3rd Qu.:404
## Max.   :539
##
##      spent      raised      reelected      run_again
## Min.   : 0  Min.   : 0  Min.   :0.00  Min.   :0.000
## 1st Qu.:172196 1st Qu.: 276093 1st Qu.:1.00 1st Qu.:0.000
## Median :256307 Median : 455861 Median :2.00 Median :1.000
## Mean   :351340 Mean   : 717580 Mean   :1.96 Mean   :0.614
## 3rd Qu.:393587 3rd Qu.: 755435 3rd Qu.:3.00 3rd Qu.:1.000
## Max.   :3489592 Max.   :9790929 Max.   :5.00 Max.   :1.000
##
##      female      support      incumbent      tweets
## Min.   :0.000  Length:539  Min.   :0.00  Min.   : 20
## 1st Qu.:0.000  Class :character 1st Qu.:0.00 1st Qu.:553
## Median :0.000  Mode  :character Median :1.00 Median :646
## Mean   :0.315  Mean   :0.57 Mean   :600
## 3rd Qu.:1.000 3rd Qu.:1.00 3rd Qu.:706
## Max.   :1.000 Max.   :1.00 Max.   :980
##
##      income      committee
## Min.   : 58372  Length:539
## 1st Qu.: 97574  Class :character
## Median :126581  Mode  :character
## Mean   :165333
```

```
## 3rd Qu.: 165242
## Max.    :1107015
##
```

One measure at a time

```
length(na.omit(data$reelected)) # The number of non-missing observations
```

```
## [1] 537
```

```
mean(data$reelected, na.rm=T) #mean
```

```
## [1] 1.96089
```

```
sd(data$reelected, na.rm=T) #standard deviation
```

```
## [1] 1.56583
```

```
min(data$reelected, na.rm=T) #minimum value
```

```
## [1] 0
```

```
max(data$reelected, na.rm=T) #maximum value
```

```
## [1] 5
```

Set the CRAN mirror. This is where we will install packages from.

```
r <- getOption("repos")
r["CRAN"] <- "http://www.stats.bris.ac.uk/R/"
options(repos = r)
```

Install the package psych if it is not already installed. You only do this once.

```
if(!require(psych)){
  install.packages('psych')
}
```

Loading required package: psych

Load the package into your R session. For every new session.

```
library(psych)
```

Use the describe function on the psych package to get descriptive statistics.

```
describe(data)
```

```
##          vars    n      mean      sd median  trimmed      mad    min
## id          1 539    270.00   155.74   270    270.00    200.15     1
## region*     2 539      3.00    1.42     3      3.00      1.48     1
## party*      3 539      2.00    0.82     2      2.00      1.48     1
## chamber*    4 539      1.19    0.40     1      1.12      0.00     1
## spent       5 539 351340.45 344042.64 256307 288132.41 149300.79     0
## raised      6 539 717579.71 949969.55 455861 529747.58 322852.46     0
## reelected   7 537      1.96    1.57     2      1.85      1.48     0
## run_again    8 539      0.61    0.49     1      0.64      0.00     0
## female      9 539      0.32    0.47     0      0.27      0.00     0
## support*    10 539      1.72    0.89     1      1.66      0.00     1
## incumbent   11 539      0.57    0.50     1      0.59      0.00     0
## tweets     12 539     600.40   159.95   646     625.27    105.26    20
## income     13 539 165332.86 123305.44 126581 138906.98 44172.58 58372
## committee* 14 539      3.32    1.07     4      3.52      0.00     1
```



```
##           max   range skew kurtosis      se
## id         539     538  0.00   -1.21    6.71
## region*      5       4  0.00   -1.31    0.06
## party*       3       2  0.00   -1.50    0.04
## chamber*     2       1  1.54    0.36    0.02
## spent    3489592 3489592  3.61   19.70 14818.97
## raised    9790929 9790929  4.43   26.96 40918.09
## reelected   5       5  0.33   -0.93    0.07
## run_again   1       1 -0.47   -1.78    0.02
## female      1       1  0.79   -1.37    0.02
## support*    3       2  0.57   -1.51    0.04
## incumbent   1       1 -0.28   -1.93    0.02
## tweets      980     960 -1.52    2.34    6.89
## income    1107015 1048643  3.47   15.94  5311.14
## committee*  4       3 -1.23   -0.06    0.05
```

TOPIC 6: Frequency and proportions tables

Frequency table:

```
mytable <- table(data$party,data$reelected)
```

var1 will be on the rows, var2 will be on the columns.

```
mytable  # print table
```

```
##
##           0  1  2  3  4  5
## Centre 43 33 36 37 15 14
## Left   42 32 36 37 16 16
## Right  45 35 30 41 15 14
```

Relative frequency (or proportions) table. Note that this is a table of the frequency table defined above.

```
prop.table(mytable)  #cell percentages
```

```
##
##           0          1          2          3          4          5
## Centre 0.0800745 0.0614525 0.0670391 0.0689013 0.0279330 0.0260708
## Left   0.0782123 0.0595903 0.0670391 0.0689013 0.0297952 0.0297952
## Right  0.0837989 0.0651769 0.0558659 0.0763501 0.0279330 0.0260708
```

```
prop.table(mytable, 1)  #row percentages
```

```
##
##           0          1          2          3          4          5
## Centre 0.2415730 0.1853933 0.2022472 0.2078652 0.0842697 0.0786517
## Left   0.2346369 0.1787709 0.2011173 0.2067039 0.0893855 0.0893855
## Right  0.2500000 0.1944444 0.1666667 0.2277778 0.0833333 0.0777778
```

```
prop.table(mytable, 2)  #column percentages
```

```
##
##           0          1          2          3          4          5
## Centre 0.330769 0.330000 0.352941 0.321739 0.326087 0.318182
## Left   0.323077 0.320000 0.352941 0.321739 0.347826 0.363636
## Right  0.346154 0.350000 0.294118 0.356522 0.326087 0.318182
```

```
table_data=data.frame(mytable)
```

Write the table out to a comma separated file:

```
write.table(table_data, file = 'table_data.csv', sep=',')
```

Topic 7: Simple plots

General syntax: `PlotType(data_to_plot, option1, option2, option3, etc)` where:

PlotType can be: `plot` (does a scatterplot), `barplot`, `pie` (for piechart), `hist` (for histogram), `boxplot`, etc.

`data_to_plot` can be a variable, a table, data coming from some model, etc.

options can be: 'main' and 'sub' for main title and sub-title; 'col' for changing colors; 'ylab' and 'xlab' for labeling the axes; legend; names.arg for labeling plot elements etc.

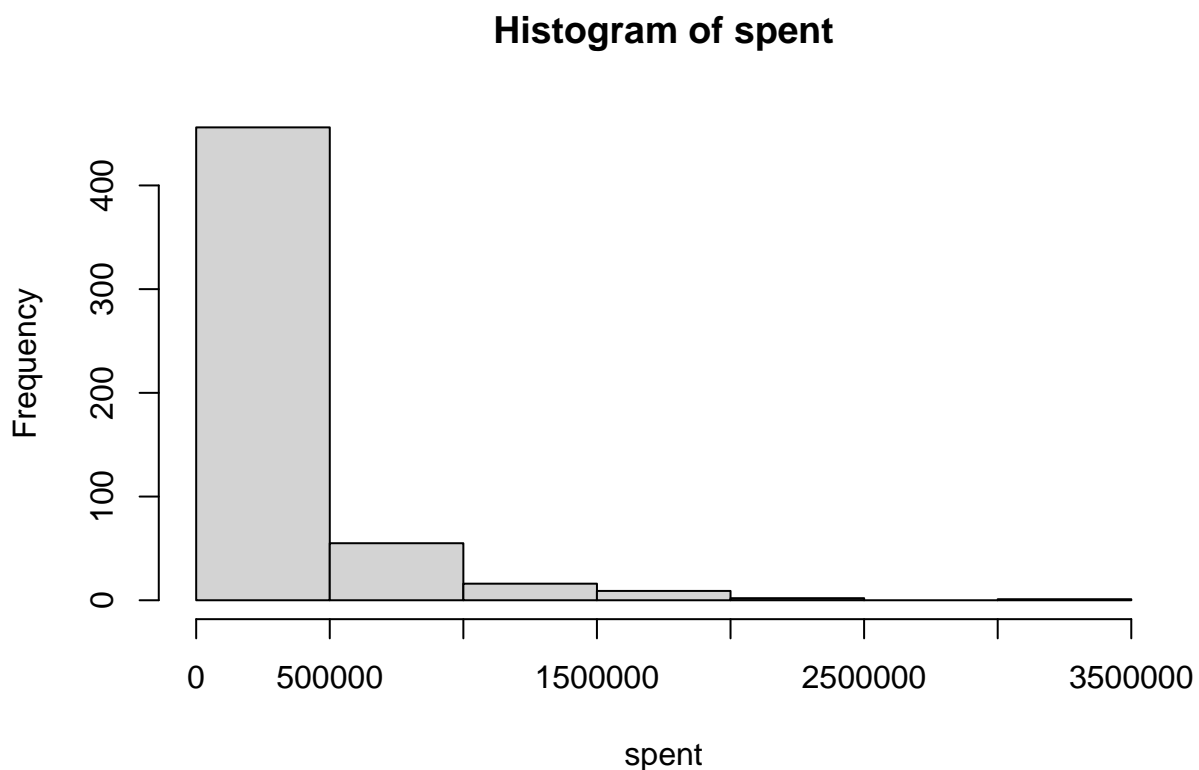
Let's first attach the dataset so we can refer to variable names directly.

```
attach(data)
```

Examples:

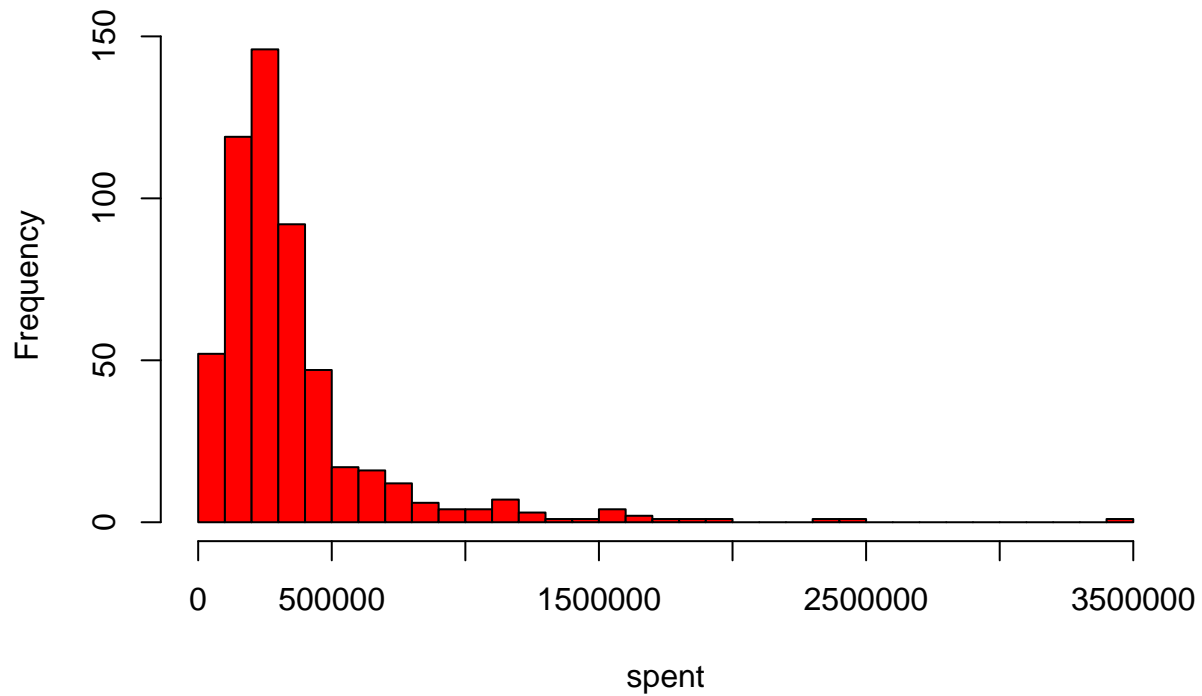
- Histogram:

```
hist(spent) # simple histogram
```



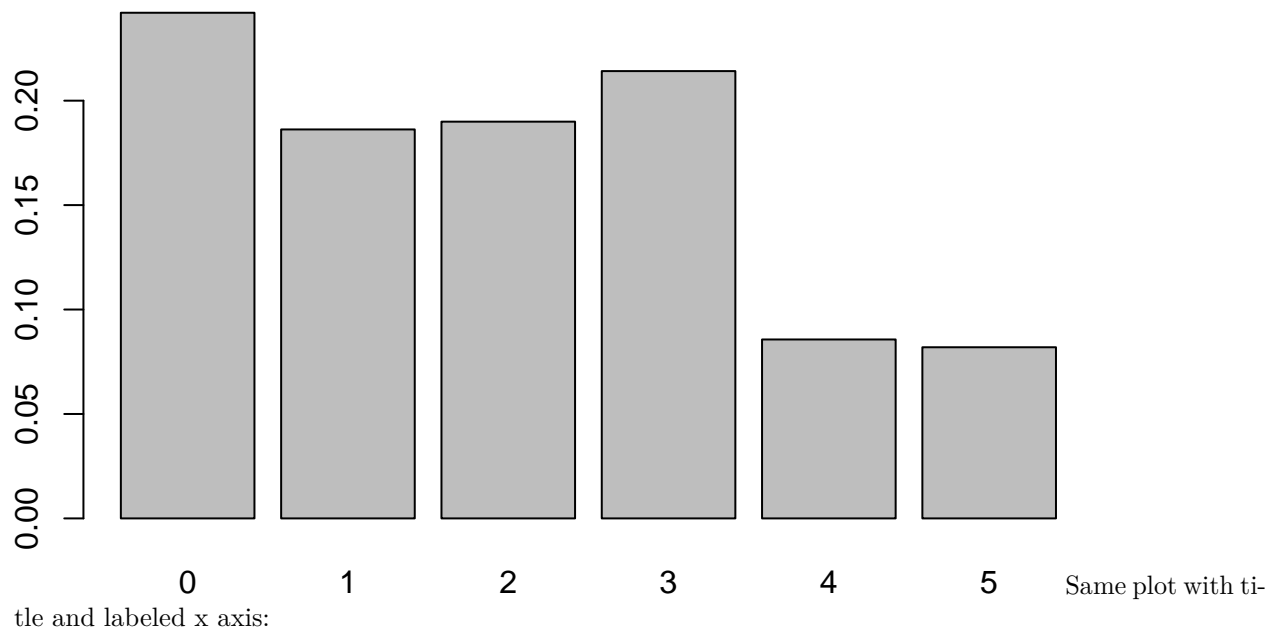
```
hist(spent, breaks=30, col="red") # red histogram with different number of bins
```

Histogram of spent



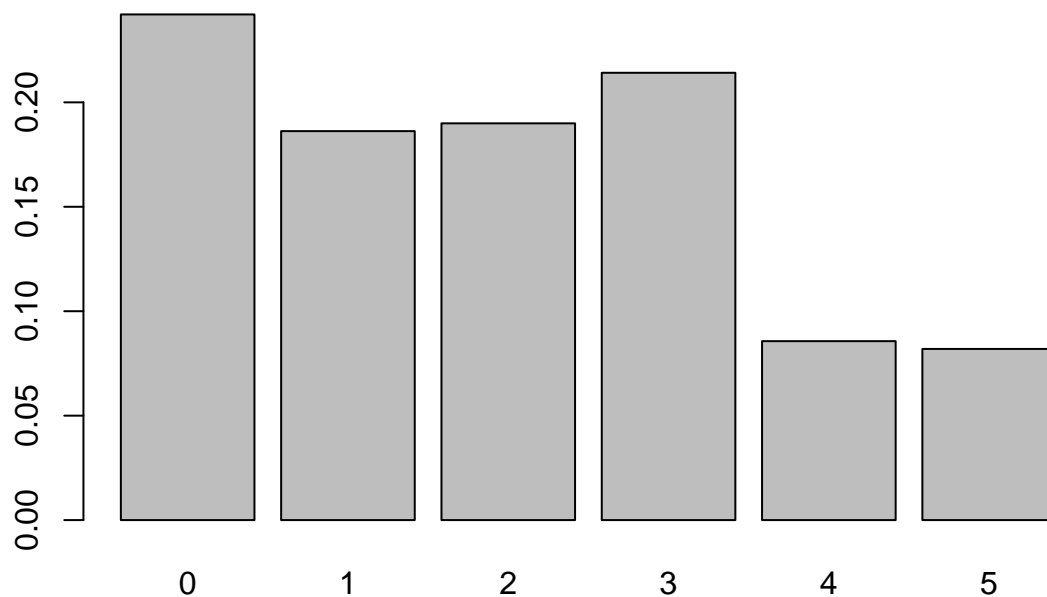
- Barplots: Simple bar plot of a frequency table for a variable (an ordinal or binary variable):

```
barplot(prop.table(table(reelected)))
```



```
barplot(prop.table(table(reelected)), main="Number of observations per category of reelected",
        xlab="Number of times reelected")
```

Number of observations per category of reelected

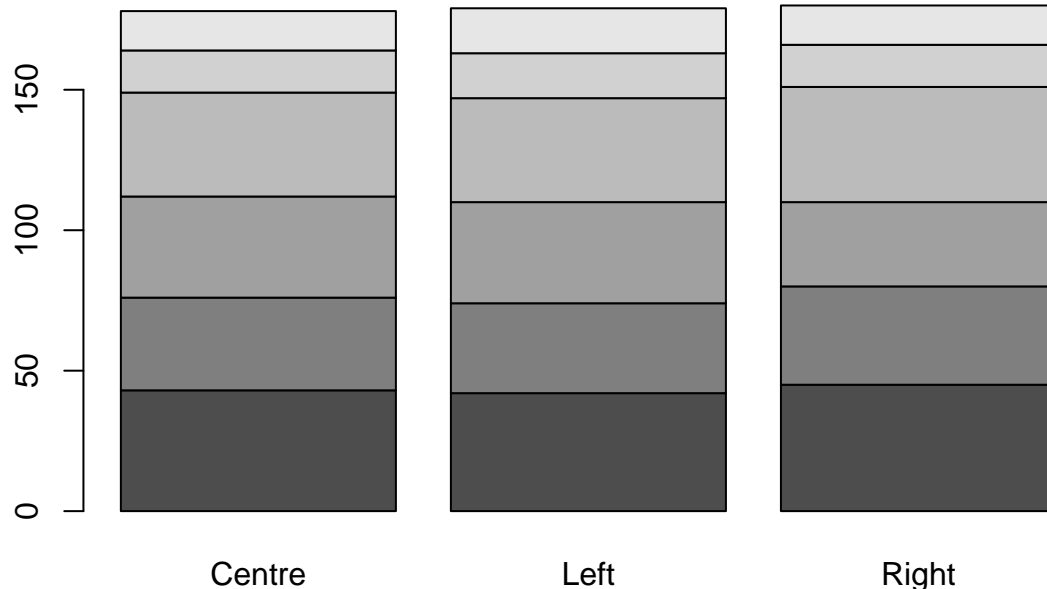


Number of times reelected

Stacked bar plot

of the number of observations by group of var1 and var2

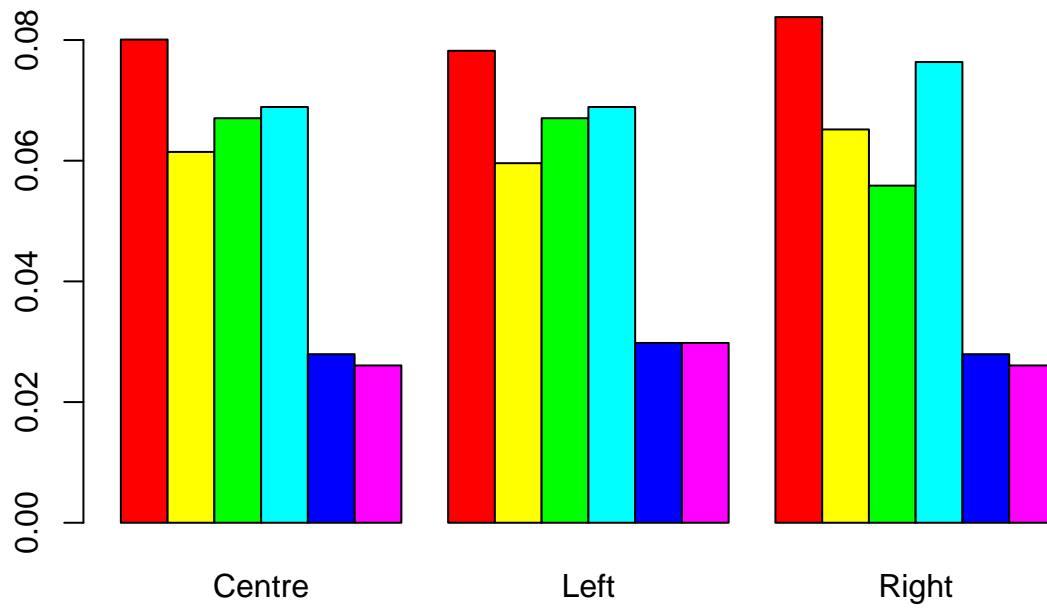
```
barplot(table(reelected,party)) # but so ugly!
```



Same as above,

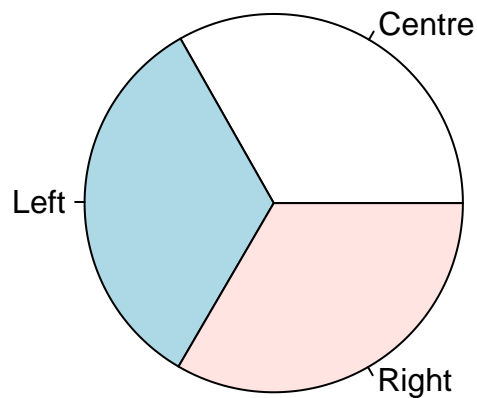
but grouped bar plot (not stacked but aside) and letting R choose the colours:

```
barplot(prop.table(table(reelected,party)), col=rainbow(6), beside=TRUE)
```



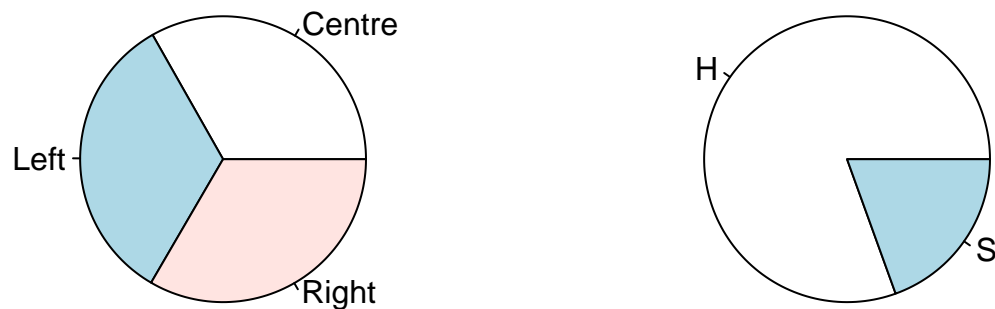
- Piecharts

```
pie(table(party))    #Simple pie chart of var1
```



Put two graphs in the same plot, one next to the other:

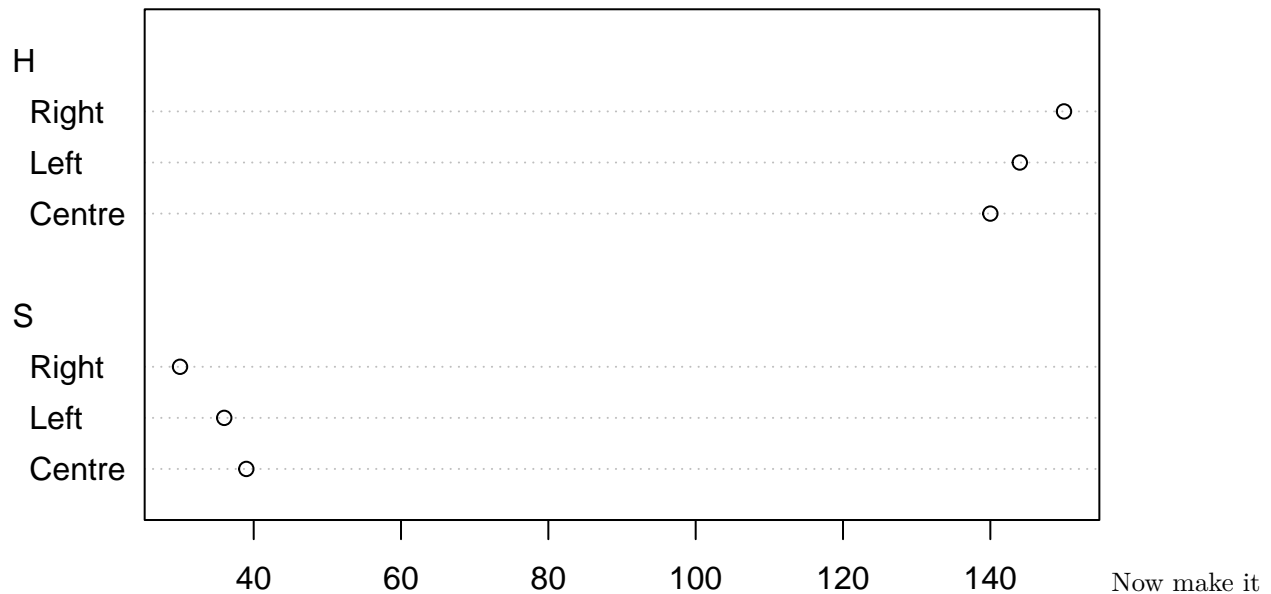
```
par(mfrow=c(1,2))    #This divides the plotting space into two vertical cells.
pie(table(party))     #First plot shows on the left.
pie(table(chamber))   #Second plot shows on the right.
```



```
par(mfrow=c(1,1))    #This turns it back to a single cell
```

- Dotcharts

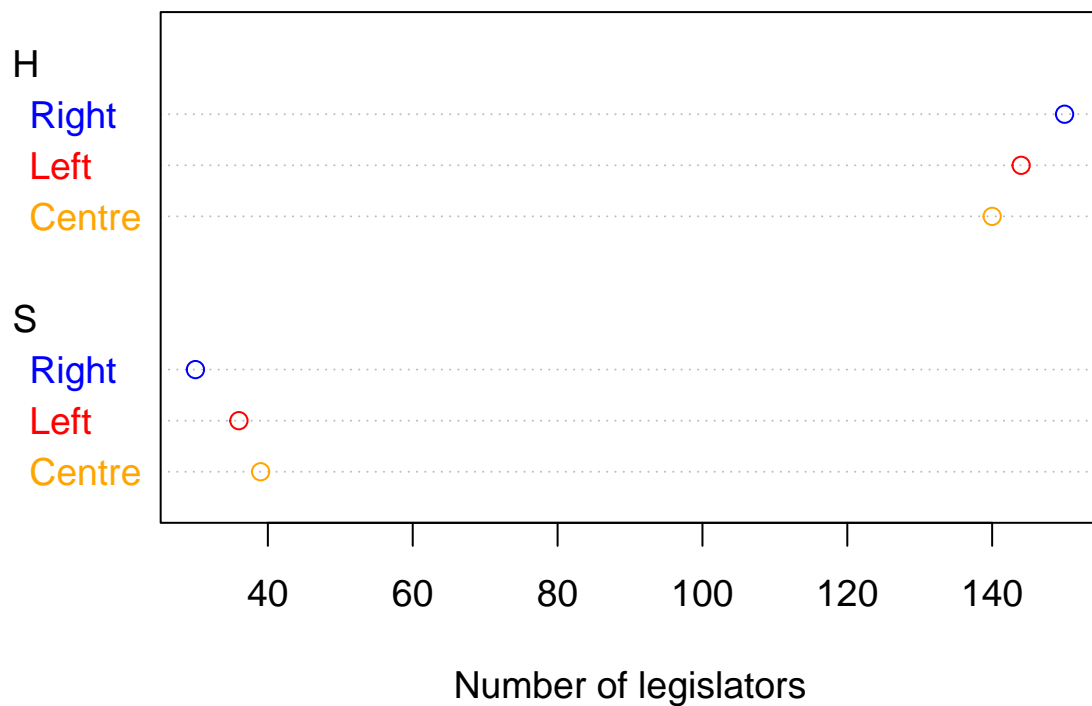
```
dotchart(table(data$party, data$chamber)) #number of legislators by party and chamber
```



pretty:

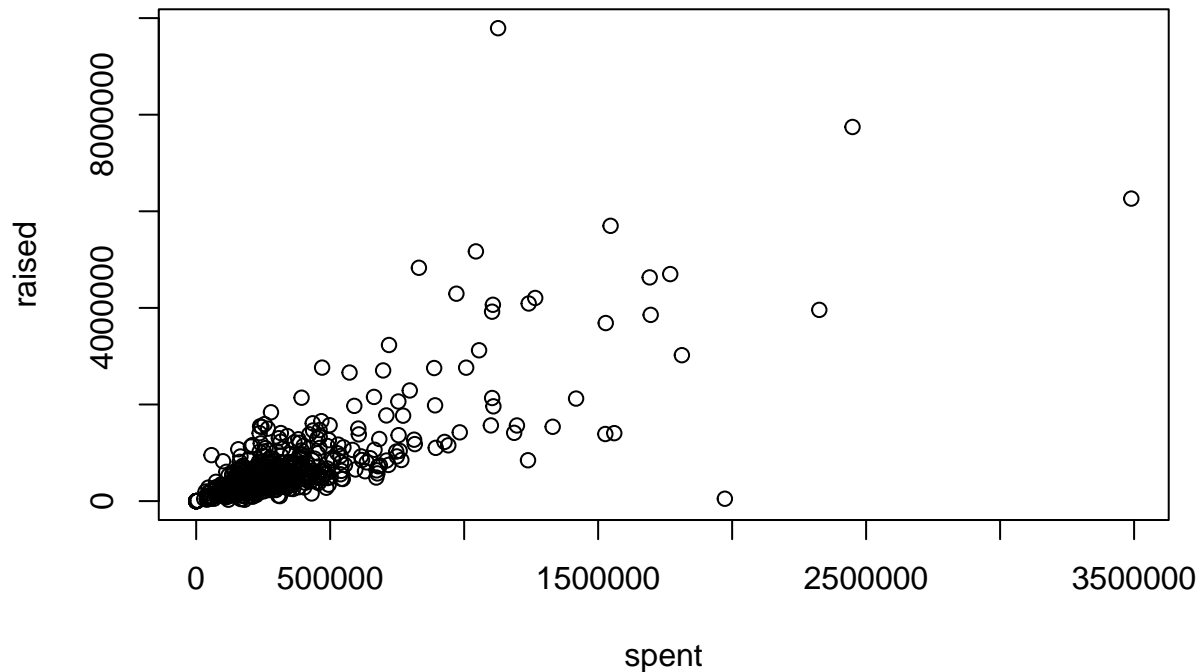
```
dotchart(table(data$party, data$chamber), cex=1.2, color=c("orange", "red", "blue"),
        main="Legislators by house and party",
        xlab="Number of legislators")
```

Legislators by house and party



- Scatterplots

```
plot(spent, raised) # Simple scatterplot of variables x and y.
```

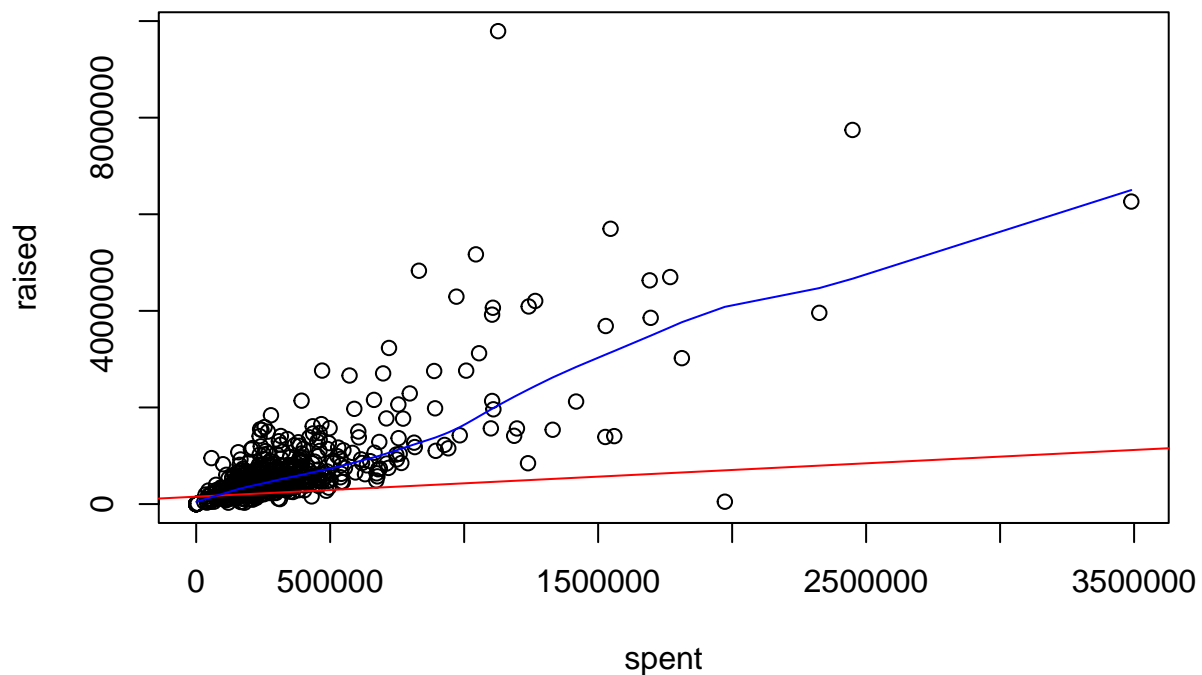


This

order puts x on the horizontal axis and y on the vertical one.

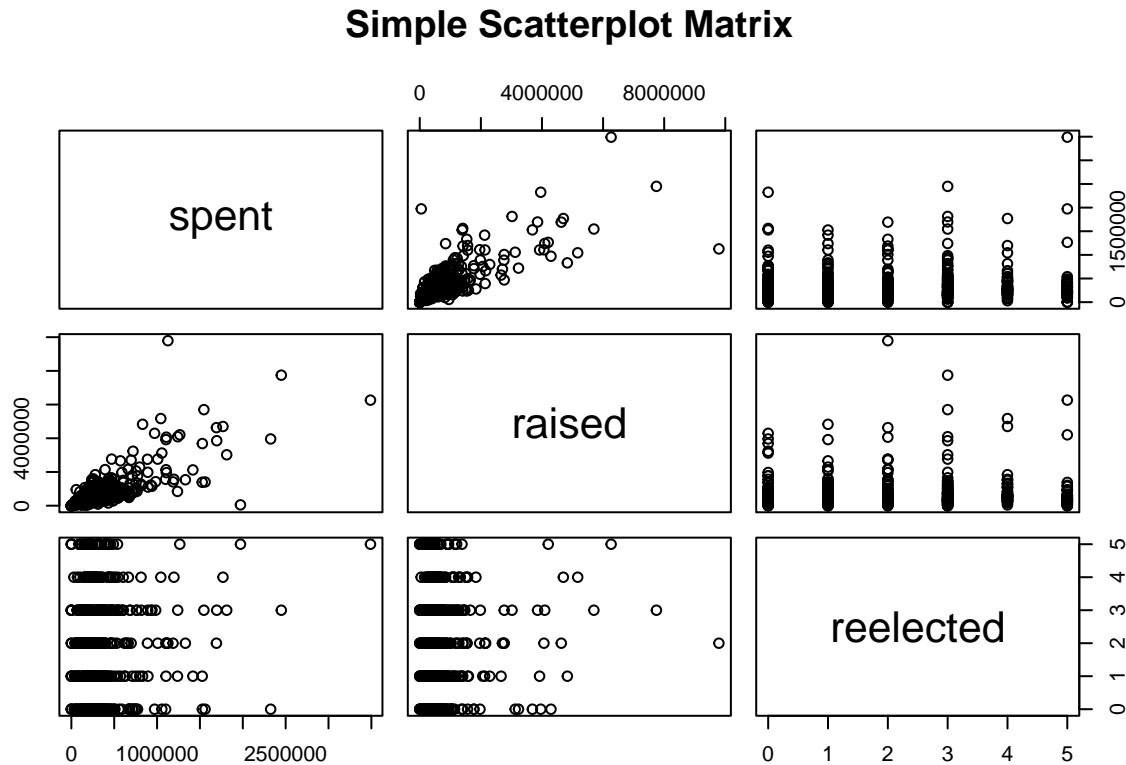
Note that the order is particularly important if you want to add a linear fit line:

```
plot(spent, raised) # Simple scatterplot of variables x and y.
abline(lm(spent~raised), col="red") # add a regression line (y~x)
lines(lowess(spent, raised), col="blue") # add a lowess line
```



Matrices of scatterplots, for multiple variables (two by two)

```
pairs(~spent+raised+reelected, data=data, main="Simple Scatterplot Matrix")
```



Other intro to R resources

Today's workshop relied heavily on "Introduction to Scientific Programming and Simulation Using R" by Jones, Millardet and Robinson

<https://www.crcpress.com/Introduction-to-Scientific-Programming-and-Simulation-Using-R-Second-Edition/Jones-Maillardet-Robinson/9781466569997>.

There are multiple excellent introductions to R online: <http://www.r-tutor.com/r-introduction> <http://tryr.codeschool.com/> <http://www.statmethods.net/>

Very good resource for more advanced programming: Hadley Wickham's Advanced R (which is not THAT advanced) <http://adv-r.had.co.nz/>

Contact info

For any questions about this tutorial, please email me at iulia.cioroianu@gmail.com.