# Data manipulation with R

*Hrvoje Stojic*

*September 15, 2017*

## Data import/export

It is evident that before performing any analysis in R you need to import the data of interest. But data come in many different formats, so you'll have to adapt and learn. Luckily, R has a plethora of input options[1], among them Excel files and databases. Through some useful packages, such as `foreign` or `xlsReadWrite`, many other data formats can be imported as well, such as SAS, SPSS or STATA.

[1] See R Data Import/Export at cran.r-project.org/doc/manuals/r-release/R-data.pdf

   Before getting any data into R, it is advisable to create a data directory in your machine where to store it. We will use it as the working directory - the default path where R will look to when importing/exporting data.

## Importing native .RData files

This is the simplest way, but usually datasets are not disseminated publicly in this format. You can use `load` function to load such a file. Word of caution. Be careful with loading such files, they might contain many R objects, some of them with same names that you have in your current working environment. If this is the case, when you load the file, objects with same names in your environment will be overwritten.

```
load("dataset.RData")
```

## Importing plain text files

One of the standard ways to exchange data (specially in organizations with a low IT profile) are plain text files. There are two paradigms for storing data tables in plain text files: **delimited text** and **fixed width**. In both cases they may or may not have a **header** (column names) and before importing you can explore its contents with a plain text editor or a spreadsheet program (except if the file is huge, some programs may not support too many gigabytes).

   In **delimited text files** the data columns are explicitly delimited, tipically with `;`, `,` or a tabulator (coded as `\t`). You may find other characters as separators. The standard file is **CSV** (comma-separated values, `,`).

   The function `read.table` can import most of the delimited files you will find. Some non-standard files may need other specific

functions, or they might be simply too big. Here we cover only
`read.table`.

   Download the *Demographic data of census sections in Barcelona*
from BCN Open Data.[2] to your working directory. It contains de-
mographic information for each census division. Inspect it with a
plain text editor (if you cannot understand the headers look at their
online description). This is an example of a standard import - it has a
header row (column names) and the field separator is ;.

```
## Import
censusBCN <- read.table("data/MAP_SCENSAL.csv",
    header = TRUE, sep = ";")

## Translating headers into English
names(censusBCN)
names(censusBCN) <- c("Date", "Men", "CensusDivision",
    "Women", "AGE_0_14", "AGE_15_A_24", "AGE_25_A_64",
    "AGE_65_plus", "NATIONALS", "EUCommunity",
    "Overseas")

## Data summary
summary(censusBCN)

## Computing a new column: percent of senior
## citizens
censusBCN$percentSenior <- censusBCN$AGE_65_plus/(censusBCN$Men +
    censusBCN$Women)
summary(censusBCN$percentSenior)

##  [1] "DATA_DADES"    "HOMES"
##  [3] "SECCIO_CENSAL" "DONES"
##  [5] "EDAT_0_A_14"   "EDAT_15_A_24"
##  [7] "EDAT_25_A_64"  "EDAT_65_A_MES"
##  [9] "NACIONALS"     "COMUNITARIS"
## [11] "ESTRANGERS"
##       Date            Men
##  30/03/15:1068   Min.   : 285.0
##                  1st Qu.: 593.0
##                  Median : 683.0
##                  Mean   : 712.2
##                  3rd Qu.: 791.0
##                  Max.   :2085.0
##  CensusDivision     Women
##  Min.   : 1001   Min.   : 261.0
##  1st Qu.: 3041   1st Qu.: 667.0
```

```
##  Median : 6036   Median : 777.0
##  Mean   : 5784   Mean   : 792.4
##  3rd Qu.: 8092   3rd Qu.: 888.0
##  Max.   :10237   Max.   :1607.0
##     AGE_0_14    AGE_15_A_24
##  Min.   : 66   Min.   : 45.0
##  1st Qu.:144   1st Qu.:104.0
##  Median :173   Median :125.0
##  Mean   :188   Mean   :131.3
##  3rd Qu.:215   3rd Qu.:150.0
##  Max.   :652   Max.   :412.0
##   AGE_25_A_64     AGE_65_plus
##  Min.   : 330.0   Min.   : 92.0
##  1st Qu.: 699.0   1st Qu.:275.8
##  Median : 820.0   Median :324.0
##  Mean   : 855.6   Mean   :329.7
##  3rd Qu.: 960.0   3rd Qu.:375.0
##  Max.   :2204.0   Max.   :703.0
##    NATIONALS    EUCommunity
##  Min.   : 462   Min.   :   3.00
##  1st Qu.:1083   1st Qu.: 29.00
##  Median :1240   Median : 53.00
##  Mean   :1263   Mean   : 70.10
##  3rd Qu.:1412   3rd Qu.: 89.25
##  Max.   :2847   Max.   :454.00
##     Overseas
##  Min.   :   10.0
##  1st Qu.:   91.0
##  Median :  134.0
##  Mean   :  171.9
##  3rd Qu.:  199.0
##  Max.   : 1733.0
##    Min. 1st Qu.  Median    Mean 3rd Qu.
## 0.04845 0.19616 0.22296 0.22383 0.25212
##    Max.
## 0.48013
```

Download the Ageing Population file from the UK Open Data site to your working directory[3]. Inspect it with a plain text editor, and you will see the key import characteristics - it has a header row (column names), the field separator is , , and some strings are quoted with double quotation marks. With this in mind we can import it.

[3] opendata.s3.amazonaws.com/aging-population-2008.csv

```
## Import
agingPopulation <- read.table("data/aging-population-2008.csv",
    header = TRUE, sep = ",", quote = "\"")

## Rename last column
names(agingPopulation)
names(agingPopulation)[7] <- "PercentSeniors"

## Data summary
summary(agingPopulation)
sapply(agingPopulation, class)


## [1] "SOA.Code"
## [2] "SOA.Name"
## [3] "Ward.Name"
## [4] "Locality"
## [5] "District.Borough"
## [6] "Total.Population"
## [7] "X..of.Total.Population.aged.60...Females...65...Males."
##       SOA.Code                SOA.Name
##  E01031005:  1   Town Centre      :  2
##  E01031006:  1   Town Centre North:  2
##  E01031007:  1   Whitestone South :  2
##  E01031008:  1   Abbey East       :  1
##  E01031009:  1   Abbey Fields     :  1
##  E01031010:  1   Abbey North      :  1
##  (Other)  :327   (Other)          :324
##         Ward.Name
##  Abbey       : 10
##  Brunswick   :  6
##  Milverton   :  6
##  Warwick North:  6
##  Warwick South:  6
##  Warwick West :  6
##  (Other)      :293
##                               Locality
##  Stratford-on-Avon & Shipston    : 24
##  Alcester, Studley & Henley      : 23
##  Southam, Wellesbourne & Kineton : 22
##  North Leamington                : 18
##  Warwick                         : 18
##  Rugby Rural                     : 17
##  (Other)                         :211
##             District.Borough
```

```
##   North Warwickshire :38
##   Nuneaton & Bedworth:82
##   Rugby              :58
##   Stratford-on-Avon  :71
##   Warwick            :84
##
##
##   Total.Population PercentSeniors
##   1,514   :  4      19.0%  :  6
##   1,594   :  4      21.7%  :  6
##   1,450   :  3      14.2%  :  4
##   1,482   :  3      14.3%  :  4
##   1,516   :  3      15.8%  :  4
##   1,527   :  3      17.0%  :  4
##   (Other):313       (Other):305
##           SOA.Code         SOA.Name
##           "factor"         "factor"
##          Ward.Name         Locality
##           "factor"         "factor"
## District.Borough Total.Population
##           "factor"         "factor"
##     PercentSeniors
##           "factor"
```

The last two columns are numeric, but the format coerces them into character. Later on we will try to convert them to numeric. Crucial detail when importing datasets is that R by default converts character variables into factors. This is governed by argument `stringsAsFactors`. While this was a standard procedure in the times when datasets would be used in linear regression analysis, these days when we use data in many other ways, such transformation is usually unwanted. More so, because transformation into a factor is relatively tricky to undo. As a rule you should set `stringsAsFactors` to FALSE.

```r
agingPopulation <- read.table("data/aging-population-2008.csv",
    header = TRUE, sep = ",", quote = "\"", stringsAsFactors = FALSE)

sapply(agingPopulation, class)
```

```
##                                         SOA.Code
##                                      "character"
##                                         SOA.Name
##                                      "character"
##                                        Ward.Name
```

```
##                                          "character"
##                                            Locality
##                                          "character"
##                                   District.Borough
##                                          "character"
##                                   Total.Population
##                                          "character"
## X..of.Total.Population.aged.60...Females...65...Males.
##                                          "character"
```

**Practice** comma-separated import from a url site. Use `read.table` but instead of naming a file in your working directory use a URL location. Import, for example, the *List of Campus* (Complementary Table 50)[4] from the Upen Data UPF site.[5]

**Fixed-width files** have no delimiter between columns, so you will need a description defining the width of each column. We will work out an example from the INE (Spanish statistical office): Survey on Human Resources in Science and Technology 2009.[6] INE's microdata (individual responses to surveys) are usually stored as fixed/width files accompanied with a metadata spreadsheet. Download both the raw data[7] and its metadata[8], and import it into R.

```
## Metadata for the selected columns
RRHH09Widths <- c(6, 2, 4, 2, 4, 1, 4, 4, 4, 1,
    1, 2, 2, 2, 2, 2, 1, 1)

RRHH09Names <- c("MUIDENT", "CCAARESI", "ANONAC",
    "CCAANAC", "CONTNACIM", "RELA", "CONTNAC1",
    "CONTNAC2", "CONTNAC3", "SEXO", "ESTADOCIVIL",
    "DEPEN5", "DEPEN18", "DEPENMAS", "NIVESTPA",
    "NIVESTMA", "NIVPROFPA", "NIVPROFMA")

## Fixed-width import function
fwfDataFrame <- read.fwf(file = "data/RRHH09.txt",
    n = 100, widths = RRHH09Widths, col.names = RRHH09Names)

## Data summary
summary(fwfDataFrame)

##      MUIDENT            CCAARESI
##   Min.   :10003    Min.    : 8.0
##   1st Qu.:20007    1st Qu.:10.0
##   Median :30008    Median :10.0
##   Mean   :23433    Mean    :10.8
##   3rd Qu.:30049    3rd Qu.:10.0
```

[4] data.upf.edu/en/dataset/listado-de-campus-tabla-complementaria-50

[5] Hint for this dataset: since some strings include the quote '"' as a character, you must turn off the quoting option (quote="").

[6] ine.es/en/prodyser/microdatos_en.htm

[7] `ftp://www.ine.es/temas/recurciencia/micro_recurciencia.zip`

[8] `ftp://www.ine.es/temas/recurciencia/disreg_recurciencia.xls`

```
## Max.   :30096   Max.   :16.0
##
##       ANONAC        CCAANAC   CONTNACIM
## Min.   :1940   10     :39   EU27:98
## 1st Qu.:1961   08     :18   RAME: 2
## Median :1968   16     :15
## Mean   :1966   13     : 9
## 3rd Qu.:1971   09     : 4
## Max.   :1979   14     : 4
##                (Other):11
##       RELA      CONTNAC1   CONTNAC2
## Min.   :1.00   EU27:98   EU27: 2
## 1st Qu.:1.00   RAME: 2   NNNN:98
## Median :1.00
## Mean   :1.04
## 3rd Qu.:1.00
## Max.   :3.00
##
## CONTNAC3         SEXO        ESTADOCIVIL
## NNNN:100   Min.   :1.00   Min.   :1.00
##           1st Qu.:1.00   1st Qu.:1.00
##           Median :1.00   Median :1.00
##           Mean   :1.46   Mean   :1.93
##           3rd Qu.:2.00   3rd Qu.:1.00
##           Max.   :2.00   Max.   :6.00
##
##     DEPEN5        DEPEN18
## Min.   :0.00   Min.   :0.00
## 1st Qu.:0.00   1st Qu.:0.00
## Median :0.00   Median :0.00
## Mean   :0.38   Mean   :0.72
## 3rd Qu.:1.00   3rd Qu.:1.00
## Max.   :2.00   Max.   :4.00
##
##    DEPENMAS       NIVESTPA
## Min.   :0.00   Min.   : 1.0
## 1st Qu.:0.00   1st Qu.: 2.0
## Median :0.00   Median : 3.0
## Mean   :0.27   Mean   : 4.5
## 3rd Qu.:0.00   3rd Qu.: 7.0
## Max.   :3.00   Max.   :10.0
##
##    NIVESTMA       NIVPROFPA
## Min.   :1.00   Min.   :1.00
```

```
##  1st Qu.:2.00   1st Qu.:1.00
##  Median :3.00   Median :1.00
##  Mean   :3.59   Mean   :1.26
##  3rd Qu.:6.00   3rd Qu.:1.00
##  Max.   :9.00   Max.   :3.00
##
##    NIVPROFMA
##  Min.   :1.00
##  1st Qu.:1.00
##  Median :5.00
##  Mean   :3.41
##  3rd Qu.:5.00
##  Max.   :5.00
##
```

We imported only a sample of the data (the first 100 rows and the first 18 columns). It is usually a good idea for an initial exploration of the data.

**Practice** fixed-width import. Import the all the columns of the `RRHH09.txt` file (but only 1000 rows). You will have to look at the metadata for the column widths and names. Summarize the columns, in particular the labor market situation (SITLAB). Is the employment rate among respondents high? (decypher its code values reading the metadata).

### Connecting to databases

A lucky analyst may work for an organization with all of its information stored in databases. This is generally a good thing because importing plain text files requires time to explore the file and fine-tune the importing code. But when standard databases are available, R can connect to them seamlessly and facilitate the reading/writing process.

Working with databases is a topic that you will study in more details in dedicated courses. Connecting to databases from R is more advanced material, however it is important that you know that it is possible to access them from R. Not all databases are supported, but most of the mainstream ones are covered with dedicated packages, both relational (SQLite, Oracle, MySQL, PosgreSQL) and non-relational databases (MongoDB, Cassandra).

We will show you a brief example with SQLite only, mostly because in a basic version we do not have to worry about setting up credentials, while the process for accessing the other databases is similar.

Download the sample database Chinook Sqlite[9], it represents a

[9] chinookdatabase.codeplex.com/downloads/get/557773

digital media store (including tables for artists, albums, media tracks, invoices and customers).

```r
## Require the package for the SQLite DB
## install.packages('RSQLite')
library(RSQLite)


## DB Connection
drv <- dbDriver("SQLite")


# Load the DB driver
con <- dbConnect(drv, dbname = "data/Chinook_Sqlite.sqlite")


# Connect to DB List all tables in the
# connection
dbListTables(con)


## Load a DB table into a data frame
tableSQL <- dbGetQuery(con, "select * from Track")


## ... perform your analyses ...
(results <- head(tableSQL))


## Disconnect & unload
dbGetInfo(con)
dbDisconnect(con)
dbUnloadDriver(drv)

##  [1] "Album"        "Artist"
##  [3] "Customer"     "Employee"
##  [5] "Genre"        "Invoice"
##  [7] "InvoiceLine"  "MediaType"
##  [9] "Playlist"     "PlaylistTrack"
## [11] "Track"
##   TrackId
## 1       1
## 2       2
## 3       3
## 4       4
## 5       5
## 6       6
##                                    Name
## 1 For Those About To Rock (We Salute You)
## 2                        Balls to the Wall
## 3                          Fast As a Shark
```

```
## 4                      Restless and Wild
## 5                   Princess of the Dawn
## 6                   Put The Finger On You
##   AlbumId MediaTypeId GenreId
## 1       1           1       1
## 2       2           2       1
## 3       3           2       1
## 4       3           2       1
## 5       3           2       1
## 6       1           1       1
##                                                               Composer
## 1                          Angus Young, Malcolm Young, Brian Johnson
## 2                                                                 <NA>
## 3                 F. Baltes, S. Kaufman, U. Dirkscneider & W. Hoffman
## 4 F. Baltes, R.A. Smith-Diesel, S. Kaufman, U. Dirkscneider & W. Hoffman
## 5                                          Deaffy & R.A. Smith-Diesel
## 6                          Angus Young, Malcolm Young, Brian Johnson
##   Milliseconds    Bytes UnitPrice
## 1       343719 11170334      0.99
## 2       342562  5510424      0.99
## 3       230619  3990994      0.99
## 4       252051  4331779      0.99
## 5       375418  6290521      0.99
## 6       205662  6713451      0.99
## list()
```

Connectivity to databases not only imports tables into R, you can also run queries into the database, for example, by submitting the SQL code.

### Exporting data

After any serious data analysis we might want to output its results. The most common function to export data in R is write.table. Not by chance, the name reminds of the import function read.table. If no folder is specified, the file will be saved at the working directory.

Before writing data into a plain text file, we must make decisions about its output format: the field separator string, whether to output the column and row names, whether to quote strings or not, or the decimal separator.

Let us see an example by exporting the mtcars data:

```
# Typical csv export
write.table(mtcars, file = "data/mtcars.csv",
    sep = ",", quote = FALSE, row.names = FALSE)
```

```
# Custom export
write.table(mtcars, file = "data/mtcars.dat",
    sep = "\t", row.names = TRUE, dec = ",")
```

Saving in the native .RData file allows you to save multiple objects in any format, while saving in text files you are usually constrained to saving data frames. First argument specifies the object, second the name of the file to be saved.

```
save(mtcars, file="mtcars.RData")
```

Same as with import, the packages like `foreign` and `xlsReadWrite` make it easy to export data in proprietary formats[10], such as MS Excel, SPSS, SAS, or Stata.

[10] statmethods.net/input/exportingdata.html

## Transforming data

Transforming datasets and variables is essential in any data-oriented project. You will need, even for the most basic programming task, to select rows from a data frame or to merge two data sets.

### Subsetting in more details

When we introduced data frames we already saw how to select specific parts of a data set (given certain conditions). We will refresh the basics and dig a little deeper.

Subsetting in data frames uses indices on rows/columns: `[optional rows condition, optional columns condition]`. Note that you can use negative numbers to indicate that the rows/columns with those indices should be *removed*.

Logical conditions are very often used to subset the data. The idea is to produce a logical vector whose length will be the same as the number of rows (if we want to subset according to rows) or the number of columns. Then, those rows indicated with `TRUE` will be produced as an output of subsetting. We have following **logical operators** that we can use in R: <, >, <=, >=, != and ==.

```
# Logical conditions on data frame values
mtcars[mtcars$hp > 200, ][1:5, ]
mtcars[mtcars$cyl == 6, ][1:5, ]
mtcars[mtcars$cyl != 6, ][1:5, ]


##                     mpg cyl disp  hp drat
## Duster 360         14.3   8  360 245 3.21
```

```
## Cadillac Fleetwood  10.4   8  472 205 2.93
## Lincoln Continental 10.4   8  460 215 3.00
## Chrysler Imperial   14.7   8  440 230 3.23
## Camaro Z28          13.3   8  350 245 3.73
##                      wt  qsec vs am gear
## Duster 360          3.570 15.84  0  0    3
## Cadillac Fleetwood  5.250 17.98  0  0    3
## Lincoln Continental 5.424 17.82  0  0    3
## Chrysler Imperial   5.345 17.42  0  0    3
## Camaro Z28          3.840 15.41  0  0    3
##                      carb
## Duster 360              4
## Cadillac Fleetwood      4
## Lincoln Continental     4
## Chrysler Imperial       4
## Camaro Z28              4
##                  mpg cyl  disp  hp drat   wt
## Mazda RX4        21.0   6 160.0 110 3.90 2.620
## Mazda RX4 Wag    21.0   6 160.0 110 3.90 2.875
## Hornet 4 Drive   21.4   6 258.0 110 3.08 3.215
## Valiant          18.1   6 225.0 105 2.76 3.460
## Merc 280         19.2   6 167.6 123 3.92 3.440
##                  qsec vs am gear carb
## Mazda RX4        16.46  0  1    4    4
## Mazda RX4 Wag    17.02  0  1    4    4
## Hornet 4 Drive   19.44  1  0    3    1
## Valiant          20.22  1  0    3    1
## Merc 280         18.30  1  0    4    4
##                   mpg cyl  disp  hp drat
## Datsun 710        22.8   4 108.0  93 3.85
## Hornet Sportabout 18.7   8 360.0 175 3.15
## Duster 360        14.3   8 360.0 245 3.21
## Merc 240D         24.4   4 146.7  62 3.69
## Merc 230          22.8   4 140.8  95 3.92
##                    wt  qsec vs am gear carb
## Datsun 710        2.32 18.61  1  1    4    1
## Hornet Sportabout 3.44 17.02  0  0    3    2
## Duster 360        3.57 15.84  0  0    3    4
## Merc 240D         3.19 20.00  1  0    4    2
## Merc 230          3.15 22.90  1  0    4    2
```

Note that the second brackets are there simply to shorten the output to the first 5 lines. Multiple conditions can be connected with **logical expressions**: !, &, &&, |, || and xor function. Inputs to these

functions need to be logical vectors.

```
# Multiple conditions on data frame values
mtcars[mtcars$hp > 200 & mtcars$mpg > 14, ]
mtcars[mtcars$hp >= 250 | mtcars$hp <= 65, ]
```

```
##                    mpg cyl disp  hp drat
## Duster 360        14.3   8  360 245 3.21
## Chrysler Imperial 14.7   8  440 230 3.23
## Ford Pantera L    15.8   8  351 264 4.22
## Maserati Bora     15.0   8  301 335 3.54
##                     wt  qsec vs am gear
## Duster 360        3.570 15.84  0  0    3
## Chrysler Imperial 5.345 17.42  0  0    3
## Ford Pantera L    3.170 14.50  0  1    5
## Maserati Bora     3.570 14.60  0  1    5
##                   carb
## Duster 360           4
## Chrysler Imperial    4
## Ford Pantera L       4
## Maserati Bora        8
##                mpg cyl  disp  hp drat    wt
## Merc 240D     24.4   4 146.7  62 3.69 3.190
## Honda Civic   30.4   4  75.7  52 4.93 1.615
## Toyota Corolla 33.9  4  71.1  65 4.22 1.835
## Ford Pantera L 15.8  8 351.0 264 4.22 3.170
## Maserati Bora 15.0   8 301.0 335 3.54 3.570
##                qsec vs am gear carb
## Merc 240D     20.00  1  0    4    2
## Honda Civic   18.52  1  1    4    2
## Toyota Corolla 19.90 1  1    4    1
## Ford Pantera L 14.50 0  1    5    4
## Maserati Bora 14.60  0  1    5    8
```

Conditions on both rows and columns.

```
mtcars[row.names(mtcars) %in% c("Fiat 128", "Fiat X1-9"),
    c("mpg", "cyl", "wt")]
```

```
##            mpg cyl    wt
## Fiat 128  32.4   4 2.200
## Fiat X1-9 27.3   4 1.935
```

Conditions using functions.

```
mtcars[substr(row.names(mtcars), 1, 4) == "Fiat",
    c("mpg", "cyl", "wt")]
mtcars[mtcars$hp == max(mtcars$hp), ]
```

```
##              mpg cyl    wt
## Fiat 128    32.4   4 2.200
## Fiat X1-9   27.3   4 1.935
##               mpg cyl disp  hp drat   wt
## Maserati Bora  15   8  301 335 3.54 3.57
##               qsec vs am gear carb
## Maserati Bora 14.6  0  1    5    8
```

Using stored conditions.

```
hpPattern <- mtcars$hp >= 250 | mtcars$hp <= 65
mtcars[hpPattern, ]
```

```
##                 mpg cyl  disp  hp drat    wt
## Merc 240D       24.4   4 146.7  62 3.69 3.190
## Honda Civic     30.4   4  75.7  52 4.93 1.615
## Toyota Corolla  33.9   4  71.1  65 4.22 1.835
## Ford Pantera L  15.8   8 351.0 264 4.22 3.170
## Maserati Bora   15.0   8 301.0 335 3.54 3.570
##                 qsec vs am gear carb
## Merc 240D       20.00  1  0    4    2
## Honda Civic     18.52  1  1    4    2
## Toyota Corolla  19.90  1  1    4    1
## Ford Pantera L  14.50  0  1    5    4
## Maserati Bora   14.60  0  1    5    8
```

*Sorting*

Sorting a vector:

```
sort(mtcars$hp, decreasing = TRUE)
```

```
##  [1] 335 264 245 245 230 215 205 180 180 180
## [11] 175 175 175 150 150 123 123 113 110 110
## [21] 110 109 105  97  95  93  91  66  66  65
## [31]  62  52
```

The subsetting notation can be also used for sorting data frames using the order function:

```
mtcars[order(mtcars$hp), ][1:5, ]
mtcars[order(mtcars$hp, decreasing = TRUE), ][1:5,
    ]
```

```
##                   mpg cyl  disp hp drat    wt
## Honda Civic      30.4    4  75.7 52 4.93 1.615
## Merc 240D        24.4    4 146.7 62 3.69 3.190
## Toyota Corolla 33.9    4  71.1 65 4.22 1.835
## Fiat 128         32.4    4  78.7 66 4.08 2.200
## Fiat X1-9        27.3    4  79.0 66 4.08 1.935
##                  qsec vs am gear carb
## Honda Civic      18.52  1  1    4    2
## Merc 240D        20.00  1  0    4    2
## Toyota Corolla 19.90  1  1    4    1
## Fiat 128         19.47  1  1    4    1
## Fiat X1-9        18.90  1  1    4    1
##                    mpg cyl disp  hp drat
## Maserati Bora      15.0   8  301 335 3.54
## Ford Pantera L     15.8   8  351 264 4.22
## Duster 360         14.3   8  360 245 3.21
## Camaro Z28         13.3   8  350 245 3.73
## Chrysler Imperial 14.7   8  440 230 3.23
##                     wt  qsec vs am gear
## Maserati Bora      3.570 14.60  0  1    5
## Ford Pantera L     3.170 14.50  0  1    5
## Duster 360         3.570 15.84  0  0    3
## Camaro Z28         3.840 15.41  0  0    3
## Chrysler Imperial 5.345 17.42  0  0    3
##                   carb
## Maserati Bora         8
## Ford Pantera L        4
## Duster 360            4
## Camaro Z28            4
## Chrysler Imperial     4
```

Ordering by multiple columns is straightforward (for clarity, we first store the row conditions on a vector):

```
# Index of sorted rows
hpOrder <- order(mtcars$hp, mtcars$mpg, decreasing = TRUE)

# Using the stored order conditions
mtcars[hpOrder, ][1:5, ]
```

```
##                    mpg cyl disp  hp drat
```

```
## Maserati Bora      15.0   8   301 335 3.54
## Ford Pantera L     15.8   8   351 264 4.22
## Duster 360         14.3   8   360 245 3.21
## Camaro Z28         13.3   8   350 245 3.73
## Chrysler Imperial 14.7    8   440 230 3.23
##                       wt  qsec vs am gear
## Maserati Bora      3.570 14.60  0  1    5
## Ford Pantera L     3.170 14.50  0  1    5
## Duster 360         3.570 15.84  0  0    3
## Camaro Z28         3.840 15.41  0  0    3
## Chrysler Imperial 5.345 17.42  0  0    3
##                    carb
## Maserati Bora         8
## Ford Pantera L        4
## Duster 360            4
## Camaro Z28            4
## Chrysler Imperial     4
```

*Appending*

To combine vectors you have already seen that you can use c function. We have used it until now to create atomic vectors, but depending on the objects you are combining, the output might be a list.

```
# combining objects in a vector
c(mtcars[1, 1], mtcars[1, 3])  # atomic vector
c(mtcars[1, 1], mtcars[1:3, 1:3])  # list


## [1]  21 160
## [[1]]
## [1] 21
##
## $mpg
## [1] 21.0 21.0 22.8
##
## $cyl
## [1] 6 6 4
##
## $disp
## [1] 160 160 108
```

Binding together several data frames with a common structure is easy. To combine data frames column-wise and row-wise, you should use functions cbind and rbind.

```
# combining data frames
cbind(mtcars[, 1], mtcars[, 3])[1:5, ]
rbind(mtcars[1:2, ], mtcars[5:6, ])
```

```
##      [,1] [,2]
## [1,] 21.0  160
## [2,] 21.0  160
## [3,] 22.8  108
## [4,] 21.4  258
## [5,] 18.7  360
##                   mpg cyl disp  hp drat
## Mazda RX4         21.0   6  160 110 3.90
## Mazda RX4 Wag     21.0   6  160 110 3.90
## Hornet Sportabout 18.7   8  360 175 3.15
## Valiant           18.1   6  225 105 2.76
##                     wt  qsec vs am gear
## Mazda RX4         2.620 16.46  0  1    4
## Mazda RX4 Wag     2.875 17.02  0  1    4
## Hornet Sportabout 3.440 17.02  0  0    3
## Valiant           3.460 20.22  1  0    3
##                   carb
## Mazda RX4            4
## Mazda RX4 Wag        4
## Hornet Sportabout    2
## Valiant              1
```

Be careful with the **broadcast** feature of R. Usually it is very useful, and you do not even notice it is at work, however, at other times, if you are not careful it can produce an undesired output that you will not notice - R will not show any warning as it assumes you know what you are doing.

```
# broadcasting allows hand abbreviations such
# as
x <- matrix(NA, 4, 4)

# instead of
matrix(rep(NA, 16), 4, 4)

# comes handy in creating data frames
cbind(1, x)

# however, here it will broadcast y to fill
# out the structure if this was not an
```

```
# intention, it will be difficult to detect an
# error
y <- c(1, 2)
cbind(y, x)
```

```
##      [,1] [,2] [,3] [,4]
## [1,]   NA   NA   NA   NA
## [2,]   NA   NA   NA   NA
## [3,]   NA   NA   NA   NA
## [4,]   NA   NA   NA   NA
##      [,1] [,2] [,3] [,4] [,5]
## [1,]    1   NA   NA   NA   NA
## [2,]    1   NA   NA   NA   NA
## [3,]    1   NA   NA   NA   NA
## [4,]    1   NA   NA   NA   NA
##      y
## [1,] 1 NA NA NA NA
## [2,] 2 NA NA NA NA
## [3,] 1 NA NA NA NA
## [4,] 2 NA NA NA NA
```

Note that `cbind` and `rbind` are quite slow operations as they do not change objects in place. If you will be using them in your code really many times it will slow your code significantly. For such occasions you will have to find other solutions.

**Practice** ordering, subsetting, appending. Order the mtcars data frame by ascending number of carburetors and weight, create two data frames with the top 3 and bottom 3 rows according to this order, append them both.

## *Merging*

Adding data from a data frame to another data frame using some joining condition is an essential operation when manipulating data, because most information is stored in tables that relate to each other by some common identifier. A function that should be used for this purpose is `merge`.

```
authors <- data.frame(surname = I(c("Tukey", "Venables",
    "Tierney", "Ripley", "McNeil")), nationality = c("US",
    "Australia", "US", "UK", "Australia"), deceased = c("yes",
    rep("no", 4)))

books <- data.frame(name = I(c("Tukey", "Venables",
    "Tierney", "Ripley", "Ripley", "McNeil", "R Core")),
```

```
    title = c("Exploratory Data Analysis", "Modern Applied Statistics ...",
        "LISP-STAT", "Spatial Statistics", "Stochastic Simulation",
        "Interactive Data Analysis", "An Introduction to R"),
    other.author = c(NA, "Ripley", NA, NA, NA,
        NA, "Venables & Smith"))


merge(authors, books, by.x = "surname", by.y = "name")
```

```
##      surname nationality deceased
## 1    McNeil    Australia       no
## 2    Ripley           UK       no
## 3    Ripley           UK       no
## 4   Tierney           US       no
## 5     Tukey           US      yes
## 6  Venables    Australia       no
##                               title other.author
## 1        Interactive Data Analysis         <NA>
## 2               Spatial Statistics         <NA>
## 3             Stochastic Simulation         <NA>
## 4                        LISP-STAT         <NA>
## 5        Exploratory Data Analysis         <NA>
## 6 Modern Applied Statistics ...       Ripley
```

*Variable transformations*

We have already seen how to create new variables from existing ones. Here we will look at some special (and useful) cases:

Sometimes you need to transform a numerical variable into a categorical one. For example, divide horsepower into 2 categories: low (below average) and high (above average).

A naive approach would be:

```
# Replicate the data frame (for keeping the
# original data unchanged)
mtcarsBis <- mtcars


# Create the above- and below-average bins
mtcarsBis$hpCateg[mtcarsBis$hp < mean(mtcarsBis$hp)] <- "Low"
mtcarsBis$hpCateg[mtcarsBis$hp >= mean(mtcarsBis$hp)] <- "High"
```

A more sophisticated way:

```
mtcarsBis$hpCateg <- ifelse(test = mtcarsBis$hp >
    mean(mtcarsBis$hp), yes = "Low", no = "High")
```

Binning numerical variables into more than 2 categories could be tedious following the previous examples, but the cut function clears the way.

**Practice**: binning into multiple categories with cut function. Bin horsepower (from the mtcarsBis dataset) into 4 categories using the cut function. Add a new column to the dataset with this categorical values.

## *Data display*

The first thing to do when having data at hand data is exploring it. We use the dataset mtcars, one of the several pre-loaded datasets in R, as an introductory example.[11]

First we must make sure the import process was successful (check the number of rows and columns, make sure the numeric fields are not imported as strings, etc.).

[11] stat.ethz.ch/R-manual/R-devel/library/datasets/html/mtcars.html

```
class(mtcars)
str(mtcars)
dim(mtcars)
names(mtcars)
summary(mtcars)
sapply(mtcars, class)


## [1] "data.frame"
## 'data.frame':    32 obs. of  11 variables:
##  $ mpg : num  21 21 22.8 21.4 18.7 18.1 14.3 24.4 22.8 19.2 ...
##  $ cyl : num  6 6 4 6 8 6 8 4 4 6 ...
##  $ disp: num  160 160 108 258 360 ...
##  $ hp  : num  110 110 93 110 175 105 245 62 95 123 ...
##  $ drat: num  3.9 3.9 3.85 3.08 3.15 2.76 3.21 3.69 3.92 3.92 ...
##  $ wt  : num  2.62 2.88 2.32 3.21 3.44 ...
##  $ qsec: num  16.5 17 18.6 19.4 17 ...
##  $ vs  : num  0 0 1 1 0 1 0 1 1 1 ...
##  $ am  : num  1 1 1 0 0 0 0 0 0 0 ...
##  $ gear: num  4 4 4 3 3 3 3 4 4 4 ...
##  $ carb: num  4 4 1 1 2 1 4 2 2 4 ...
## [1] 32 11
##  [1] "mpg"  "cyl"  "disp" "hp"   "drat"
##  [6] "wt"   "qsec" "vs"   "am"   "gear"
## [11] "carb"
##       mpg             cyl
##  Min.   :10.40   Min.   :4.000
##  1st Qu.:15.43   1st Qu.:4.000
```

```
##   Median :19.20   Median :6.000
##   Mean   :20.09   Mean   :6.188
##   3rd Qu.:22.80   3rd Qu.:8.000
##   Max.   :33.90   Max.   :8.000
##       disp            hp
##   Min.   : 71.1   Min.   : 52.0
##   1st Qu.:120.8   1st Qu.: 96.5
##   Median :196.3   Median :123.0
##   Mean   :230.7   Mean   :146.7
##   3rd Qu.:326.0   3rd Qu.:180.0
##   Max.   :472.0   Max.   :335.0
##       drat            wt
##   Min.   :2.760   Min.   :1.513
##   1st Qu.:3.080   1st Qu.:2.581
##   Median :3.695   Median :3.325
##   Mean   :3.597   Mean   :3.217
##   3rd Qu.:3.920   3rd Qu.:3.610
##   Max.   :4.930   Max.   :5.424
##       qsec            vs
##   Min.   :14.50   Min.   :0.0000
##   1st Qu.:16.89   1st Qu.:0.0000
##   Median :17.71   Median :0.0000
##   Mean   :17.85   Mean   :0.4375
##   3rd Qu.:18.90   3rd Qu.:1.0000
##   Max.   :22.90   Max.   :1.0000
##       am             gear
##   Min.   :0.0000   Min.   :3.000
##   1st Qu.:0.0000   1st Qu.:3.000
##   Median :0.0000   Median :4.000
##   Mean   :0.4062   Mean   :3.688
##   3rd Qu.:1.0000   3rd Qu.:4.000
##   Max.   :1.0000   Max.   :5.000
##       carb
##   Min.   :1.000
##   1st Qu.:2.000
##   Median :2.000
##   Mean   :2.812
##   3rd Qu.:4.000
##   Max.   :8.000
##       mpg         cyl        disp          hp
## "numeric"   "numeric"   "numeric"   "numeric"
##      drat          wt        qsec          vs
## "numeric"   "numeric"   "numeric"   "numeric"
##        am        gear        carb
```

```
## "numeric" "numeric" "numeric"
```

Then we can inspect visually our table. Since most of our tables have many more rows than our screens can show we start by looking at the top and bottom rows.

```
head(mtcars)
tail(mtcars)

# you can specify the number of rows
head(mtcars, 10)
```

```
##                    mpg cyl disp  hp drat
## Mazda RX4         21.0   6  160 110 3.90
## Mazda RX4 Wag     21.0   6  160 110 3.90
## Datsun 710        22.8   4  108  93 3.85
## Hornet 4 Drive    21.4   6  258 110 3.08
## Hornet Sportabout 18.7   8  360 175 3.15
## Valiant           18.1   6  225 105 2.76
##                      wt  qsec vs am gear
## Mazda RX4         2.620 16.46  0  1    4
## Mazda RX4 Wag     2.875 17.02  0  1    4
## Datsun 710        2.320 18.61  1  1    4
## Hornet 4 Drive    3.215 19.44  1  0    3
## Hornet Sportabout 3.440 17.02  0  0    3
## Valiant           3.460 20.22  1  0    3
##                   carb
## Mazda RX4            4
## Mazda RX4 Wag        4
## Datsun 710           1
## Hornet 4 Drive       1
## Hornet Sportabout    2
## Valiant              1
##                 mpg cyl  disp  hp drat    wt
## Porsche 914-2  26.0   4 120.3  91 4.43 2.140
## Lotus Europa   30.4   4  95.1 113 3.77 1.513
## Ford Pantera L 15.8   8 351.0 264 4.22 3.170
## Ferrari Dino   19.7   6 145.0 175 3.62 2.770
## Maserati Bora  15.0   8 301.0 335 3.54 3.570
## Volvo 142E     21.4   4 121.0 109 4.11 2.780
##                qsec vs am gear carb
## Porsche 914-2  16.7  0  1    5    2
## Lotus Europa   16.9  1  1    5    2
## Ford Pantera L 14.5  0  1    5    4
## Ferrari Dino   15.5  0  1    5    6
```

```
## Maserati Bora  14.6  0  1   5    8
## Volvo 142E     18.6  1  1   4    2
##                     mpg cyl  disp  hp drat
## Mazda RX4           21.0   6 160.0 110 3.90
## Mazda RX4 Wag       21.0   6 160.0 110 3.90
## Datsun 710          22.8   4 108.0  93 3.85
## Hornet 4 Drive      21.4   6 258.0 110 3.08
## Hornet Sportabout 18.7   8 360.0 175 3.15
## Valiant             18.1   6 225.0 105 2.76
## Duster 360          14.3   8 360.0 245 3.21
## Merc 240D           24.4   4 146.7  62 3.69
## Merc 230            22.8   4 140.8  95 3.92
## Merc 280            19.2   6 167.6 123 3.92
##                      wt  qsec vs am gear
## Mazda RX4           2.620 16.46  0  1   4
## Mazda RX4 Wag       2.875 17.02  0  1   4
## Datsun 710          2.320 18.61  1  1   4
## Hornet 4 Drive      3.215 19.44  1  0   3
## Hornet Sportabout 3.440 17.02  0  0   3
## Valiant             3.460 20.22  1  0   3
## Duster 360          3.570 15.84  0  0   3
## Merc 240D           3.190 20.00  1  0   4
## Merc 230            3.150 22.90  1  0   4
## Merc 280            3.440 18.30  1  0   4
##                     carb
## Mazda RX4              4
## Mazda RX4 Wag          4
## Datsun 710             1
## Hornet 4 Drive         1
## Hornet Sportabout      2
## Valiant                1
## Duster 360             4
## Merc 240D              2
## Merc 230               2
## Merc 280               4
```

After this visual inspection we can describe individual columns, summary tables for categorical data, histograms and descriptive statistics for numeric data, etc.

## Basic graphics

Visualizing results for your analysis is critical for its success - conveying the right message). R is extraordinary powerful at graphing

data, allowing a great degree of personalization and having several state-of-the-art packages.

We will start with the basics[12], and later on we will use `ggplot2` - the package for advanced plotting in R.

**Stripcharts** are one-dimensional scatter plots and provide a (somewhat simplistic) first look at univariate series. Note the optional parameter `xlab` for setting the X-axis title.

```
stripchart(mtcars$mpg, xlab = "Miles per gallon")
```

A **histogram** cuts series in discrete bins, while a continuous distribution varies smoothly along the series. Histogram of mileage in the `mtcars` dataset.



Figure 1: Stripchart: first look at mpg.

```
hist(mtcars$mpg, main = "")
```

The default number of bins may be misleading, we can set more bins (maybe after some trial and error). Using the optional parameter `col` (setting the fill color for the histogram bins) helps focusing on the important message - the distribution.



Figure 2: Histogram of mileage.

```
hist(mtcars$mpg, col = "gray", breaks = 10, main = "")
```

The **kernel density** estimate is a hypothetical continuous distribution generating a univariate series and provides a smooth approximation for the actual distribution. Kernel density estimates are closely related to histograms, but can be endowed with properties such as smoothness or continuity by using a suitable kernel.

Note we must first compute the estimated density with the `density` function.



Figure 3: Smoother histogram of mpg.

```
# Compute the density data
d <- density(mtcars$mpg)

# Graph the results
plot(d, main = "")
```

**Boxplots** summarize univariate series in a single plot (including the range of the variable, its quartiles, and its outliers).[13] In future lectures we will dig deeper on summarizing distributions. Here we will only plot the classical boxplot, grouping by the number of cylinders.



Figure 4: Density estimate of mileage.

```
boxplot(mpg ~ cyl, data = mtcars, ylab = "MPG",
    xlab = "Number of Cylinders")
```

Figure 5: Example of boxplot summary, car milage data

**Scatterplots** display values for two variables for a set of data, and are essential when looking for relationships between them (e.g. linear correlation). In R the simplest way to plot them is the `plot` function, where we can also add a regression line.

```
# Scatterplot
plot(x = mtcars$hp, y = mtcars$wt, ylab = "Horsepower",
    xlab = "Weight (lb/1000)", pch = 16)

# Linear regression line
abline(lm(mtcars$wt ~ mtcars$hp), col = "red")
```



Figure 6: Hp vs weight: scatterplot and linear regression

**Practice** describing BCN census data. Use the Demographic data of census sections in Barcelona we imported earlier and . . .

1. produce stripcharts, histograms and kernel density estimates of a variable of your choice. Be creative: define a new variable combining existing ones, combine colors, explore optional parameters of the functions.
2. scatterplot two variables and add a linear regression line. Try adding a loess regression curve. Choose appropriate point characters and point dimensions (with lots of data points maybe blank-filled smaller dots are most convenient).
3. challenging bit: use the layout function to display all the univariate plots in a single matrix of plots.

*A full example of a customized plot*

Plotting data in R is always easy, but obtaining the format you need almost never is. Some packages will make your life easier, but you need to learn the basics of plot personalization. This is another powerful R feature.

Say we want to plot the miles per gallon of the `mtcars` data set. It is pretty straightforward.

```
plot(mtcars$mpg)
```

But you will agree it is also disappointing - bad axis titles, no main title, no identification of each car etc. But these flaws are also strengths. Plots are objects as any other and all the elements can be personalized and coded. This way, they can be reproduced when data change or if you share your code.
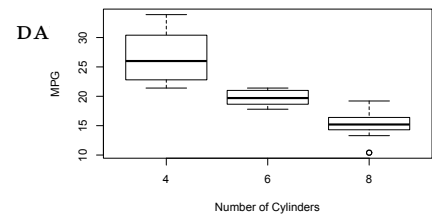
```
plot(mtcars$mpg, main = "Miles per gallon for selected cars",
    ylab = "mpg", pch = 16, cex = 0.8, ylim = c(0,
        max(mtcars$mpg)))
```

Now we have fixed some of these issues: main and ylab improve titles, pch and cex improve formats, ylim lets us set the axis limits. But the X axis remains without the proper car labels. Let us fix that too.

```
plot(mtcars$mpg, main = "Miles per gallon for selected cars",
    ylab = "mpg", pch = 16, cex = 0.8, ylim = c(0,
        max(mtcars$mpg)), xaxt = "n", xlab = "")

axis(side = 1, at = seq_along(mtcars$mpg), labels = rownames(mtcars),
    las = 2, cex.axis = 0.7)
```

With the axis function we can control the position and content of the axes (here the X axis, or side=1).

Maybe adding vertical gridlines will help identifying each car, and coloring according to the cylinders may be informative.

```
plot(mtcars$mpg, main = "Miles per gallon for selected cars",
    ylab = "mpg", pch = 16, cex = 0.8, ylim = c(0,
        max(mtcars$mpg)), xaxt = "n", xlab = "",
    col = mtcars$cyl)

axis(side = 1, at = seq_along(mtcars$mpg), labels = rownames(mtcars),
```
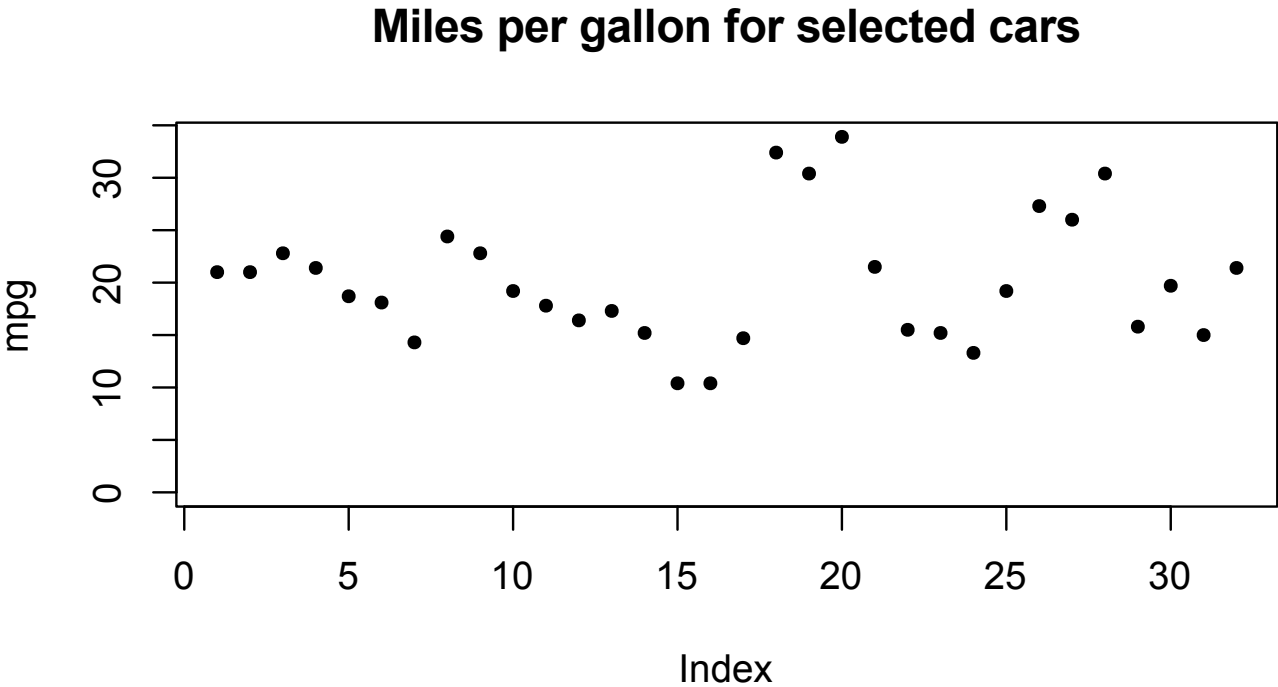
# Miles per gallon for selected cars



Figure 8: Our second attempt.
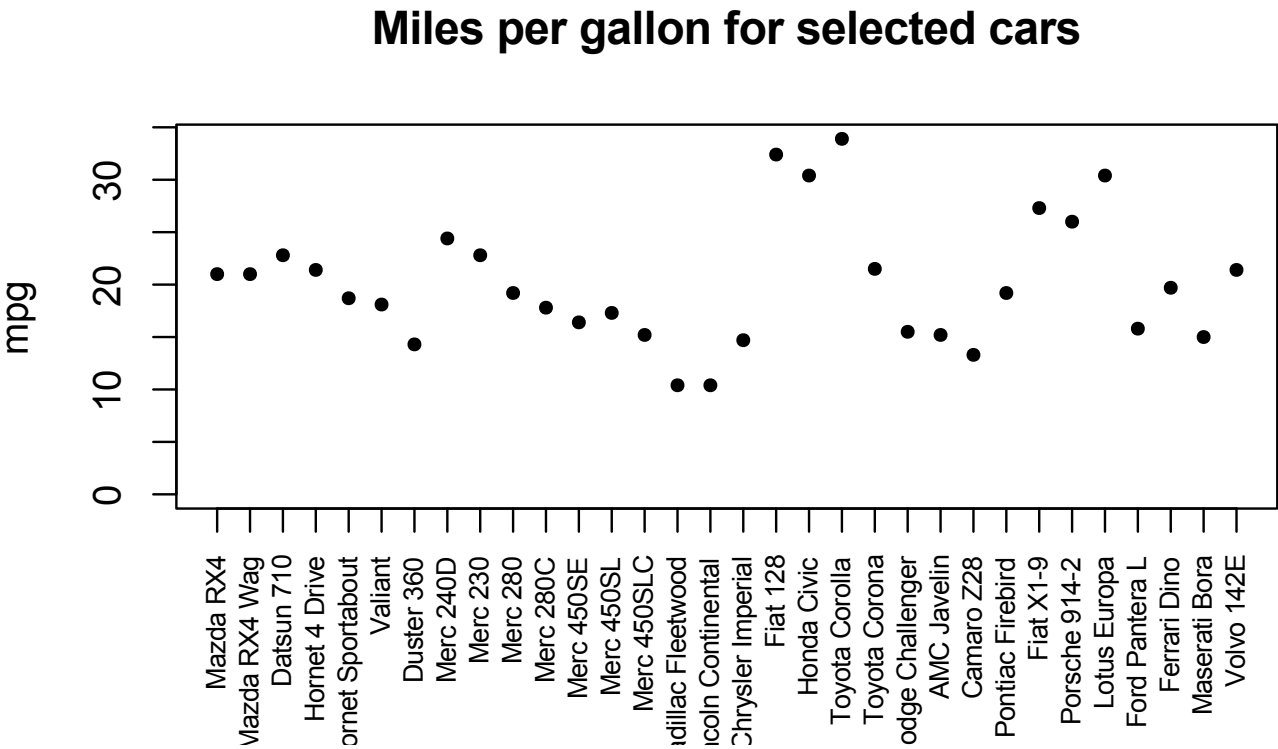
# Miles per gallon for selected cars



Figure 9: Now the X axis is informative.

```
    las = 2, cex.axis = 0.7)

abline(v = seq_along(mtcars$mpg), col = "gray",
    lty = 2)

legend(x = "bottomleft", ncol = 3, cex = 0.6,
    bg = "white", legend = c("4 cyl", "6 cyl",
        "8 cyl"), text.col = "azure4", col = c(4,
        6, 8), pch = 16)
```
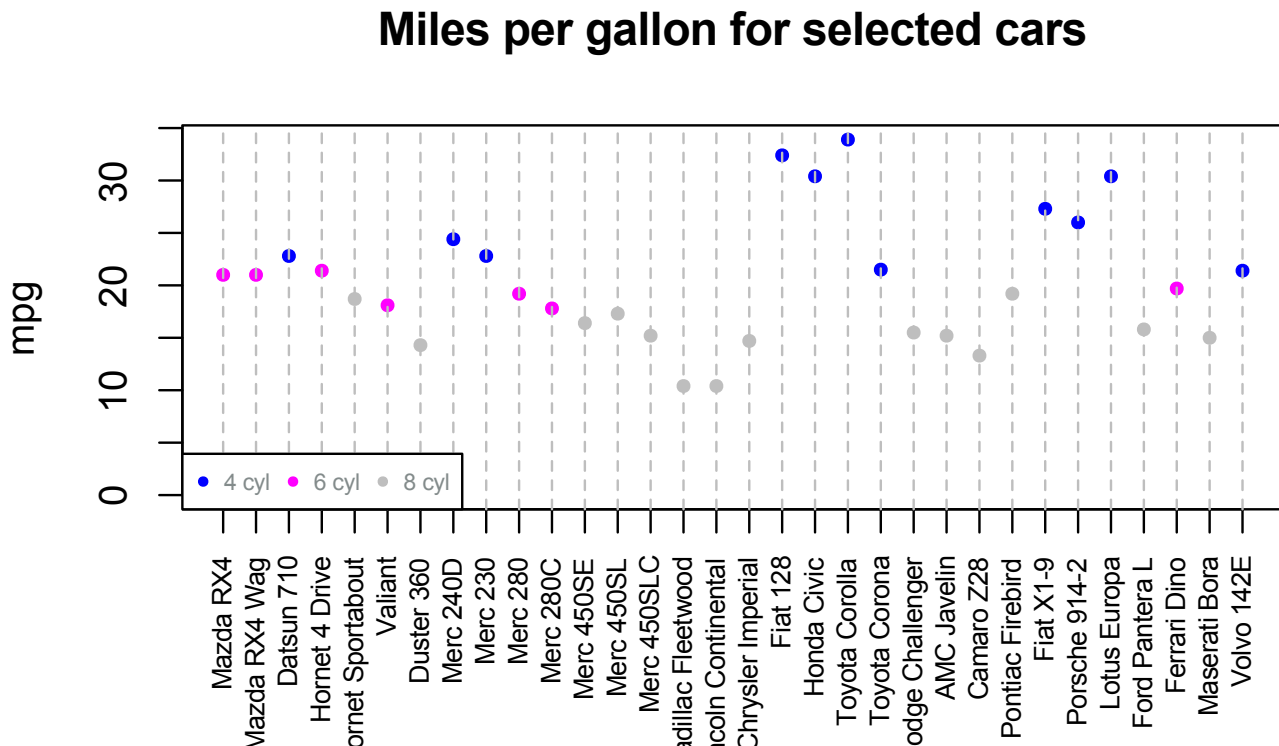


Figure 10: Finally a plot both informative and well formatted.

Changing fonts is a bit more tricky. Moreover, R by default uses Helvetica font in figures and these fonts are not necessarily available everywhere as it is a commercial font. Hence, some pdf readers might not render correctly the figures. See extrafont package for some extra options with fonts.

*Saving plots*

By default, a plot in R opens a device in your desktop (or within RStudio). But usually you will need to save it as a high-quality image. R allows different outpot formats for graphics (pdf, jpg, png). Here we will set an example of pdf output.

```
## Open the device
pdf("mileage.pdf", width = 10, height = 2.5)

## Add the plot
plot(mtcars$mpg)
abline(h = mean(mtcars$mpg),
       col = "lightgray",
       lty = 2)

## Close the device
dev.off()
```

Note that no directory path is specified. By default, the file will be saved at the current working directory, optionally you can set a different output directory.

*Basic text manipulation*

In almost every analysis you need to perform operations on dates and text strings. Here we will take a look at the essential operations on these types of data.

Handling strings in R can sometimes be painful.[14] Several packages exist that ease this pain. Here we look only at R base functions.

The `paste` function is perhaps the most used in R when handling strings. It essentially concatenates strings, but in a generalized way (e.g. you can choose the character separating strings, or it converts non-string objects to characters). By default, it concatenates strings separating them with a blank space.

```
paste("Barcelona", "GSE")
```

```
## [1] "Barcelona GSE"
```

The sep parameter sets a different separator. The `paste0` function is a convenient alternative when you don't want any form of separation, so it is as if you set `sep = ""`.

```
paste("Barcelona", "GSE", sep = "-")
paste0("Barcelona", "GSE")
```

```
## [1] "Barcelona-GSE"
## [1] "BarcelonaGSE"
```

Numeric variables are coerced to strings.

[14] A classical reference for handling text is: Sanchez, G. (2013) Handling and Processing Strings in R. Trowchez Editions. Berkeley. gaston-sanchez.com/Handling_and_Processing_Strings_in_R.pdf

```r
paste("The Life of", pi)
```

```
## [1] "The Life of 3.14159265358979"
```

You can also operate with vectors.

```r
paste("Class of 201", 4:7, sep = "")
```

```
## [1] "Class of 2014" "Class of 2015"
## [3] "Class of 2016" "Class of 2017"
```

Count number of characters: nchar function works both with a single string or with a vector.

```r
nchar(c("How", "many", "characters?"))
nchar("How many characters?")
```

```
## [1]  3  4 11
## [1] 20
```

Convert to lower/upper case with tolower and toupper functions. Again, they also work on vectors.

```r
tolower("Barcelona GSE")
toupper(c("Barcelona", "GSE"))
```

```
## [1] "barcelona gse"
## [1] "BARCELONA" "GSE"
```

Obtain and replace substrings with substr function.

```r
substr("Barcelona GSE", start = 11, stop = 13)
days <- c("Mond", "Tues", "Wedn")
substr(days, 4, 4) <- "."
days
```

```
## [1] "GSE"
## [1] "Mon." "Tue." "Wed."
```

Character translation with chartr.

```r
chartr(old = "4", new = "a", "B4rcelon4 GSE")
chartr(old = "410", new = "aio", "B4rcel0n4 Gr4du4te Sch00l of Ec0n0m1cs")
```

```
## [1] "Barcelona GSE"
## [1] "Barcelona Graduate School of Economics"
```

Uniquely abbreviate strings with abbreviate.

```
abbreviate(c("Statistical Models", "Deterministic Models",
    "Data Warehousing"), minlength = 8)
```

```
##   Statistical Models Deterministic Models
##           "SttstclM"           "DtrmnstM"
##     Data Warehousing
##           "DtWrhsng"
```

**Practice** character to numeric when importing. The Ageing Population data imported earlier has some numerical columns stored as character. Use basic string manipulations to convert them back into numerical.[15]

## *Dates and times in R*

Dates are represented as the number of days since 1970-01-01, with negative values for earlier dates. But R outputs them with the familiar formats (e.g. MMDDYYY). Converting to/from dates and operating with them requires some familiarity with the main R date formats and functions.[16] There are packages like `lubridate` that facilitate handling of dates and time.

**System date and time** are useful for many purposes (e.g. computing execution times, saving files with dynamical names).

```
# System date with default format
Sys.time()

# Time with HH:MM:SS format
format(Sys.time(), "%H:%M:%S")

# Date with YYYYMMDD format
format(Sys.time(), "%Y-%m-%d")

# using system time for measuring difference
x <- Sys.time()
y <- Sys.time()
y - x

# there is a specific function for that
system.time(for (i in 1:100) mad(runif(1000)))
```

```
## [1] "2017-09-16 01:14:51 BST"
## [1] "01:14:51"
## [1] "2017-09-16"
```

```
## Time difference of 0.001306772 secs
##     user  system elapsed
##    0.018   0.000   0.018
```

Output of `system.time` might be confusing. **User CPU time** gives the CPU time spent by the current R session, while **system CPU time** gives the CPU time spent by the operating system on behalf of the R session. The operating system might do additional other operations like opening files, doing input or output, starting other processes, and looking at the system clock, operations that involve resources that many processes must share. **Elapsed time** is the sum of the two.

Converting strings to date/time objects is the name of the game when importing data files. Being familiar with date conversion and formatting is also crucial when reporting results. Some examples.

```
# Input date format
x <- as.Date("20140912", format = "%Y%m%d")
x
class(x)
typeof(x)

# Input time and date
strptime("09/12/11 17.30.00", format = "%m/%d/%y %H.%M.%S")

# convert to string
as.character(Sys.time())
```

```
## [1] "2014-09-12"
## [1] "Date"
## [1] "double"
## [1] "2011-09-12 17:30:00 BST"
## [1] "2017-09-16 01:14:51"
```

**Extracting information from dates.**

```
# Name of weekday
weekdays(Sys.time())

# Name of month
months(Sys.time())

# Number of days since beginning of epoch
julian(Sys.time())
```

```
## [1] "Saturday"
```

```
## [1] "September"
## Time difference of 17425.01 days
```

Julian Day Number (JDN)[17] is the number of days since noon UTC on the first day of 4317 BC.

**Generating sequences of dates**

```r
seq(from = as.Date("2014-09-12"), to = as.Date("2014-09-14"),
    by = "day")

# All days between two dates
seq(from = as.Date("2014-09-12"), to = as.Date("2014-11-12"),
    by = "month")

# All months between two dates
seq(from = as.Date("2014-09-12"), to = as.Date("2014-09-16"),
    length.out = 3)

# Every other day between two dates Next 3
# days
seq.Date(Sys.Date(), length = 3, by = "1 days")

# Next 3 months
seq.Date(Sys.Date(), length = 3, by = "1 months")


## [1] "2014-09-12" "2014-09-13" "2014-09-14"
## [1] "2014-09-12" "2014-10-12" "2014-11-12"
## [1] "2014-09-12" "2014-09-14" "2014-09-16"
## [1] "2017-09-16" "2017-09-17" "2017-09-18"
## [1] "2017-09-16" "2017-10-16" "2017-11-16"
```

Operations with dates.

```r
# Number of days since a given date
julian(Sys.time()) - julian(as.Date("2014-01-01"))

# Adding days
as.Date("2014-09-12") + 30

# Adding months
seq.Date(Sys.Date(), length = 2, by = "3 months")[2]


## Time difference of 1354.01 days
## [1] "2014-10-12"
## [1] "2017-12-16"
```

**Practice** proper formatting of dates in imported data. Create a
new column in the BCN census data containing the day after the first
column date. You must use the `as.Date` function.