# DS116 - Data Visualization
### Intro to R

Habet Madoyan

American University of Armenia

Section 1

**Intro to R**

# Intro to R

- Go to Comprehensive R Archive Network (CRAN) and install version for your operating system
- Go to RStudio and install RStudio IDE
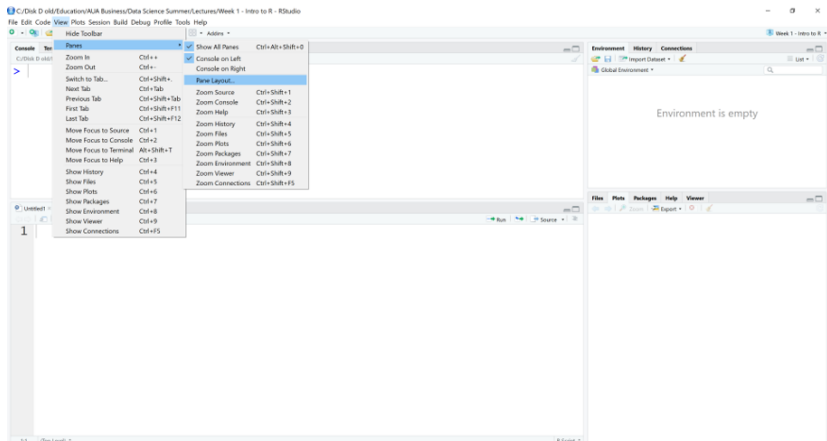- You need to install R before installing RStudio

**Alternatives**
- Install Microsoft R
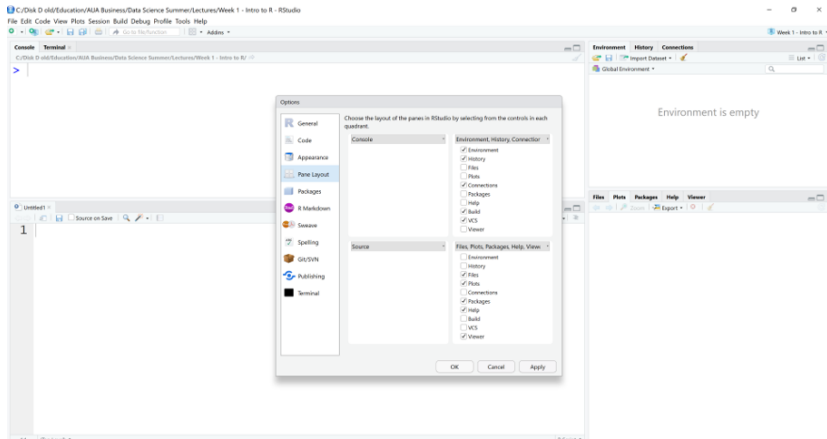- Use Microsoft Visual Studio (free version available)



R-Programming

# Intro to R

Configure the pane layout of RStudio as you wish
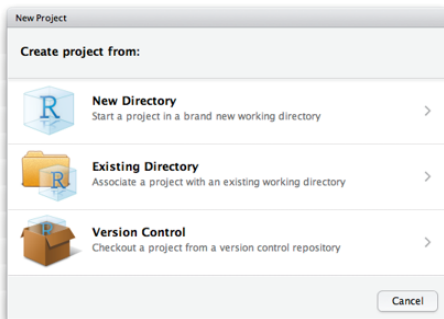
# Intro to R

# Intro to R

RStudio projects make it straightforward to divide your work into multiple contexts, each with their own working directory, workspace, history, and source documents.

**Creating Projects**

- RStudio projects are associated with R working directories. You can create an RStudio project:
  - in a brand new directory;
  - in an existing directory where you already have R code and data;
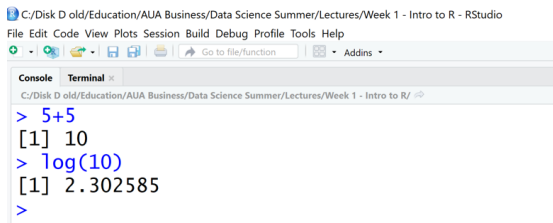  - by cloning a version control (Git or Subversion) repository.

# Intro to R

- Create an empty directory
- Go to File → New Project
- Choose existing directory

# Intro to R

- Now when you have created your project, all files associated with the project need to be saved in the project directory
- To check the project directory do getwd() in Console



C:/Disk D old/Education/AUA Business/Data Science Summer/Lectures/Week 1 - Intro to R - RStudio
File Edit Code View Plots Session Build Debug Profile Tools Help

Console Terminal
C:/Disk D old/Education/AUA Business/Data Science Summer/Lectures/Week 1 - Intro to R/

```
> 5+5
[1] 10
> log(10)
[1] 2.302585
>
```

The name of the project appears at the top of RStudio window

# Why do we need RStudio Project?

The benefits of RStudio project:

- Can have different projects with their own environment.
- While opening a project, R restores previous work and project history includes recent commands which had been used in the project.
- Source pane remembers the files which had been opened.

# Intro to R

When a new project is created, RStudio:

- Creates a project file (with an .Rproj extension) within the project directory. This file contains various project options (discussed below) and can also be used as a shortcut for opening the project directly from the filesystem.
- Creates a hidden directory (named .Rproj.user) where project-specific temporary files (e.g. auto-saved source documents, window-state, etc.) are stored. This directory is also automatically added to .Rbuildignore, .gitignore, etc. if required.
- Loads the project into RStudio and displays its name in the Projects toolbar (which is located on the far right side of the main toolbar).
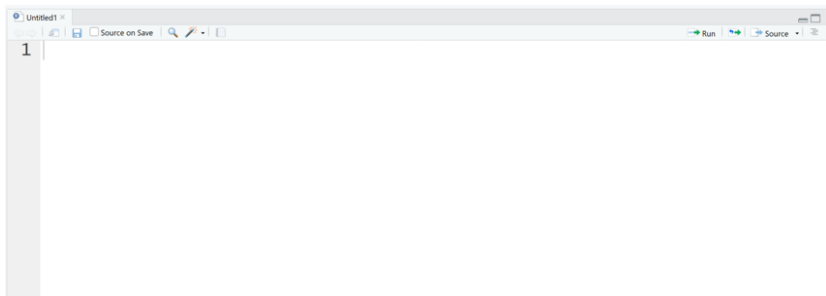
# Intro to R

**Opening Projects**

There are several ways to open a project:

- Using the **Open Project** command (available from both the Projects menu and the Projects toolbar) to browse for and select an existing project file (e.g. MyProject.Rproj).
- Selecting a project from the list of the most recently opened projects (also available from both the Projects menu and toolbar).
- Double-clicking on the project file within the system shell (e.g. Windows Explorer, OSX Finder, etc.).

# Intro to R

**The source editor**
- If you plan to reuse your code, write it in source editor

## Intro to R

Here you can type any valid R command after the $>$ prompt followed by **Enter** and R will execute that command.



```
Console  Terminal
C:/Disk D old/Education/AUA Business/Data Science Summer/Lectures/Week 1 - Intro to R/
> 5+5
[1] 10
> |
```

Use console as a calculator

[1] is the index for the output, just ignore it

# Intro to R

This will create new script file

## Intro to R

- To run the script line from source editor, put the cursor anywhere on the line and hit **Ctrl+Enter**
- You will see the output in the console

# Intro to R

- If you have a piece of text in your editor that is not a code (thus is not executable by R) then you need to comment it, add # before each line
- When you run this line of code, it will be printed in Console as a text
- To comment large chunk of text, use **Ctrl+Shift+C**

```
C:/Disk D old/Education/AUA Business/Data Science Summer/Lectures/Week 1 - Intro to R/
> # This command will calculate the logarithm of 10
>
> log(10)
[1] 2.302585
>
> # This line will calculate factorial of 10
>
> factorial(10)
[1] 3628800
>
```

```
class example.R*
    Source on Save
 1
 2  # This command will calculate the logarithm of 10
 3
 4  log(10)
 5
 6  # This line will calculate factorial of 10
 7
 8  factorial(10)
```

## Intro to R

- The environment window contains objects (data, values, functions) R has currently stored in its memory.
- The history window shows all commands that were executed in the Console.



```
Environment   History
5+5
log(10)
log(10)
log(10)
log(10)
factorial(10)
factorial
```

# Intro to R

Bottom right: files, plots, packages, help, & viewer pane. Here you can open files, view plots, install and load packages, read main pages, and view markdown and other documents in the viewer tab.

# Intro to R

Some useful shortcuts:

- type first few letters of the function/object then hit Tab to open dropdown menu with possible options

## Intro to R

It's often the case that you want to re-execute commands that you previously entered. The RStudio console supports the ability to recall previous commands using the arrow keys:

- **Up** - Recall previous command(s)
- **Down** - Reverse of Up

You can even view a list of your recent commands by pressing **Ctrl+Up** on Windows or **Command+Up** on a Mac.

# Intro to R

- When you download R from the CRAN, you get that "base" R system.
- The base R system comes with basic functionality; implements the R language.
- One reason R is so useful is the large collection of packages that extend the basic functionality of R.
- R packages are developed and published by the larger R community.

# Intro to R

- Packages can be installed with the install.packages() function in R.
- To install a single package, pass the name of the library to the install.packages() function as the first argument.
- The following code installs the devtools package from CRAN.

install.packages("devtools")

# Intro to R

- The package needs to be installed only once
- To load the package into R environment you need to use function library()
- You need to load the library everytime you start a new R

```
library(devtools)
```

## Intro to R

- You can also install packages from github
- The following code will install the following package from github

```
library(devtools)
install_github("christophM/iml")
```

If you do not want to load the entire package but want to use some function from it, use the following command

package::function_name

```
devtools::install_github("christophM/iml")
```

# Intro to R

To access help/documentation on a function from R base package

```
?mean
```

To access help/documentation on function from a library

```
#??geom_path
```

The same

```
#?ggplot2::geom_path
```

Help on the package

```
help(package='ggplot2')
```

# Intro to R

- Each package on CRAN has its own webpage
- This includes documentation and sometimes includes vignettes

- If you have a specific task to do, then look at R Task View
- Here are all the packages and R functionality described for the Time series analysis

Section 2

# Intro to R programming language

# Intro to R programming language

- Anything in R is an object.
- Objects are assigned values using $<-$ . (An equal sign $=$ can also be used.)
  For e.g., the following command assigns value 5 to object x.

```
x <- 5
x
## [1] 5
```

- R is case sensitive, thus Data10 and data10 are two different objects.
- A tidy code requires a space before the assignment operator and a space after.
    - You can see what is inside the object just by simple entering the name of the object in command line
    - When the object is created it should appear in your Environment Window
    - If the object is not in your environment window you cannot work with it
    - The assignment operator works in the opposite direction as well

```
5 -> x
x
## [1] 5
```

# Data Structres in R

R programming supports five basic types of data structure.

- **Homogeneous**

Vector                    Matrix                    Array



- **Heterogeneous**

Data frame                          List

# Data Structures: Vector

**Vector** is a type of data structure that contains similar types of data, i.e., integer, double, logical, complex, etc. In order to create a vector c() function is used.

```
(x <- c(10,5,6))
## [1] 10 5 6
```

```
(x1 <- c(1:10))
## [1] 1 2 3 4 5 6 7 8 9 10
```
The same as (only when you have a sequence)

```
(x1 <- 1:10)
## [1] 1 2 3 4 5 6 7 8 9 10
```

```
class(x)
## [1] "numeric"
```
A character vector

```
y <- c("CS", "DS", "EC")
class(y)
## [1] "character"
```

# Coercion

- All elements of an atomic vector must be the same type, so when you attempt to combine different types they will be coerced to the most flexible type.
- Types from least to most flexible are: logical, numeric, and character.
- For example, combining a character and an integer yields a character:

```
(x3 <- c("A",1))
```
```
## [1] "A" "1"
```

```
class(x3)
```
```
## [1] "character"
```

# Data Structures: Matrix

- **Matrix** is a two-dimensional data structure and can be created using matrix() function.
- Matrix is a collection of vectors with the same length and type

```
m <- matrix(data=1:15, nrow=3)
m
##      [,1] [,2] [,3] [,4] [,5]
## [1,]    1    4    7   10   13
## [2,]    2    5    8   11   14
## [3,]    3    6    9   12   15
```

## Function

- To look at the function arguments, hit Tab
- You need to name the arguments if they are not in the same order as defined within the function
- If the order is the same you can skip the names

# Data Structures

Example: function matrix with arguments flipped

The same result

```
(m <- matrix(1:15, 3))
```
```
##      [,1] [,2] [,3] [,4] [,5]
## [1,]    1    4    7   10   13
## [2,]    2    5    8   11   14
## [3,]    3    6    9   12   15
```

Different result

```
(m <- matrix(3, 1:15))
```
```
##      [,1]
## [1,]    3
```

# Data Structures

Other ways of building a matrix

- Create two vectors of the same length

```
x <- 0:5
y <- -5:0
```

- Combine together by row

```
(m1 <- rbind(x,y))
```
```
##   [,1] [,2] [,3] [,4] [,5] [,6]
## x    0    1    2    3    4    5
## y   -5   -4   -3   -2   -1    0
```

# Data Structures

The matrix has row names but does not have column names

```
colnames(m1)
```
## NULL

```
rownames(m1)
```
## [1] "x" "y"

# Data Structures

- Combine together by column

```
(m2 <- cbind(x,y))
```

```
##      x  y
## [1,] 0 -5
## [2,] 1 -4
## [3,] 2 -3
## [4,] 3 -2
## [5,] 4 -1
## [6,] 5  0
```

Pay attention, the matrix has column names, but does not have row names

```
colnames(m2)
```

```
## [1] "x" "y"
```

```
rownames(m2)
```

```
## NULL
```

# Data Structures

Check if the resulting object is a matrix

```
is.matrix(m1)
```
## [1] TRUE

```
is.matrix(m2)
```
## [1] TRUE

# Data Structures

Set the rownames

```r
rownames(m2) <- c('A', 'B', 'C', 'D', 'E', 'F')
m2
```

```
##   x  y
## A 0 -5
## B 1 -4
## C 2 -3
## D 3 -2
## E 4 -1
## F 5  0
```

# Data Structures

Subset by either integer position, or by rowname

```
m2[1:3,]
```

```
##   x  y
## A 0 -5
## B 1 -4
## C 2 -3
```

```
m2[c('A', 'B', 'C'),]
```

```
##   x  y
## A 0 -5
## B 1 -4
## C 2 -3
```

# Data Structures: Data frame

- A **data frame** is the most common way of storing data in R, and if used systematically makes data analysis easier.
- Under the hood, a data frame is a list of equal-length vectors that can also have different types.

```
df <- data.frame(x=1:10, y=11:20)
head(df, n=5)
```
```
##   x  y
## 1 1 11
## 2 2 12
## 3 3 13
## 4 4 14
## 5 5 15
```

# Data Structures

- Dataframe always has both column and row names
- If you don't specify row names it is just the row number

```
colnames(df)
```
## [1] "x" "y"

```
rownames(df)
```
## [1] "1" "2" "3" "4" "5" "6" "7" "8" "9" "10"

# Data Structures

- You can combine vectors with different type into a dataframe
- First create a matrix with two vectors of different types

```
Club <- c("Juventus", "Napoli", "Roma", "Inter")
Points <- c(95,91,77,72)
```

```
s_matrix <- cbind(Club, Points)
str(s_matrix)
## chr [1:4, 1:2] "Juventus" "Napoli" "Roma" "Inter" "95" "91" "77" "72"
## - attr(*, "dimnames")=List of 2
## ..$ : NULL
## ..$ : chr [1:2] "Club" "Points"
```

```
class(s_matrix)
## [1] "matrix" "array"
```

```
typeof(s_matrix)
## [1] "character"
```

# Data Structures

Create a dataframe

```
df_seriea <- data.frame(Club, Points)
str(df_seriea)
## 'data.frame':    4 obs. of  2 variables:
## $ Club  : chr  "Juventus" "Napoli" "Roma" "Inter"
## $ Points: num  95 91 77 72
```

Two columns (vectors) in the data frame have different types

# Data Structures: List

**List**

- In R lists act as containers.
- Unlike vectors, the contents of a list are not restricted to a single type and can encompass any mixture of data types.
- Lists are sometimes called generic vectors, because the elements of a list can be of any type of R object, even lists containing further lists.
- This property makes them fundamentally different from atomic vectors.
- A list is a special type of vector. Each element can be a different type.

# Data Structures

List containing dataframe and a matrix

```
my_list <- list(df_seriea, m1)
str(my_list)
```

```
## List of 2
## $ :'data.frame':    4 obs. of  2 variables:
##  ..$ Club  : chr [1:4] "Juventus" "Napoli" "Roma" "Inter"
##  ..$ Points: num [1:4] 95 91 77 72
## $ : int [1:2, 1:6] 0 -5 1 -4 2 -3 3 -2 4 -1 ...
##  ..- attr(*, "dimnames")=List of 2
##  .. ..$ : chr [1:2] "x" "y"
##  .. ..$ : NULL
```

# Data Structures

Look what is inside

```
my_list
```

```
## [[1]]
##       Club Points
## 1 Juventus    95
## 2   Napoli    91
## 3    Roma     77
## 4   Inter     72
##
## [[2]]
##   [,1] [,2] [,3] [,4] [,5] [,6]
## x    0    1    2    3    4    5
## y   -5   -4   -3   -2   -1    0
```

# Data Structures

The unlist() function will flatten the list into a vector

```
z <- unlist(my_list)
z
## Club1 Club2 Club3 Club4 Points1 Points2 Points3
## "Juventus" "Napoli" "Roma" "Inter" "95" "91" "77"
## Points4
## "72" "0" "-5" "1" "-4" "2" "-3"
##
## "3" "-2" "4" "-1" "5" "0"
```

# Data Structures

**Attributes**

- All objects can have arbitrary additional attributes, used to store metadata about the object.
- Attributes can be thought of as a named list (with unique names).
- Attributes can be accessed individually with attr() or all at once (as a list) with attributes().

```
attr(df_seriea, "topic") <- "sports"
df_seriea
```

```
##       Club Points
## 1 Juventus     95
## 2   Napoli     91
## 3     Roma     77
## 4    Inter     72
```

```
attr(df_seriea, "topic")
```

```
## [1] "sports"
```

Section 3

# Data Types

# Data Types

Data types used in R

- Logical
- Numeric
- Factor
- Character

# Data Types

Logical data type is one of the frequently used data types usually used for comparing two values. Values that a logical data type takes is TRUE or FALSE.

```
a <- 10
a > 5
## [1] TRUE
```

```
log1 <- c(5, 6, TRUE)
typeof(log1)
## [1] "double"
```

```
log1
## [1] 5 6 1
```

Coercion of numeric and logical values will result as numeric.

# Data Types

String literals or string values are stored as Character objects in R.

```r
b <- c("Armenia", "Georgia", "Azerbaijan")
typeof(b)
```
```
## [1] "character"
```

```r
(b1 <- c(b, TRUE))
```
```
## [1] "Armenia"    "Georgia"    "Azerbaijan" "TRUE"
```

# Data Types

- A **factor** is a vector that can contain only predefined values, and is used to store categorical data.
- Factors are built on top of integer vectors using two attributes: the class, factor, which makes them behave differently from regular integer vectors, and the levels, which defines the set of allowed values.

```
b <- as.factor(c("Armenia", "Georgia", "Azerbaijan"))
```

Factors in R are stored as a vector of integer values with a corresponding set of character values used when the factor is displayed.

```
typeof(b)
## [1] "integer"
```

```
as.numeric(b)
## [1] 1 3 2
```

```
b
## [1] Armenia    Georgia    Azerbaijan
## Levels: Armenia Azerbaijan Georgia
```

# Data Types

You can relevel the factor variable by changing the reference value

```
b<- relevel(b, ref='Azerbaijan')
levels(b)
## [1] "Azerbaijan" "Armenia"    "Georgia"
```

```
as.numeric(b)
## [1] 2 3 1
```

Section 4

**Special Values in R**

# Special Values in R

**Missing Values**

The examples and reasons of having missing values:

- The information was not collected (e.g. some people decline to give their age or weight).
- Some attributes are not applicable to all objects (e.g. forms have conditional parts that are filled out only when a person answers a previous question in a certain way).
- The information is not imported in the database (system missing).
- In R missing or undefined values are represented by NA.

# Special Values in R

**NA**

- In R, the NA values are used to represent missing values. (NA stands for "not available").
- You may encounter NA values in text loaded into R (to represent missing values) or in data loaded from databases (to replace NULL values).

```r
v1 <- c(1,2,4,NA,5)
```

```r
is.na(v1)
## [1] FALSE FALSE FALSE  TRUE FALSE
```

# Special Values in R

```r
v2 <- c(10, "A", 20)
```

```r
as.numeric(v2)
```
```
## Warning: NAs introduced by coercion
## [1] 10 NA 20
```

# Special Values in R

**Inf and -Inf**

If a computation results in a number that is too big, R will return Inf for a positive number and -Inf for a negative number (meaning positive and negative infinity, respectively).

```
120/0
## [1] Inf
```

```
-120/0
## [1] -Inf
```
Too big to show

```
45^12500
## [1] Inf
```

# Special Values in R

**NaN**

Sometimes, a computation will produce a result that makes little sense. In these cases, R will often return NaN (meaning "Not a Number").

```
0/0
## [1] NaN
```

# Special Values in R

**NULL**

- There is a null object in R, represented by the symbol NULL.
- NULL is often used as an argument in functions to mean that no value was assigned to the argument.
- Some functions may return NULL.
- NULL is not the same as NA, Inf, -Inf, or NaN.
- NULL represents the null object in R.
- NULL is used mainly to represent the lists with zero length, and is often returned by expressions and functions whose value is undefined.

Section 5

**Importing Data to R**

# Importing Data to R

**Flat files**

- A flat file database is a database that stores data in a plain text file.
- Each line of the text file holds one record with fields separated by delimiters, such as commas or tabs.
- While it uses a simple structure, a flat file database cannot contain multiple tables like a relational database can.

# Importing Data to R

**Delimiter-Separated Values**

- Formats that use delimiter-separated values (also DSV) store two-dimensional arrays of data by separating the values in each row with specific delimiter characters.
- Most database and spreadsheet programs are able to read or save data in a delimited format.
- Most widely used delimiters are:
  - comma (CSV – comma separated values)
  - tab (TSV – tab separated values)

# Importing Data to R

- File winter.csv contains data on winter Olympic games from 1924 to 2014.
- The file is comma-separated.
- Use readLines() with n=5 to look at the first five lines of the text file.
- You can see that the values are separated by comma.

```
readLines('Data/winter.csv', n=5)
```
```
## [1] "Year,City,Sport,Discipline,Athlete,Country,Gender,Event,Medal"
## [2] "1924,Chamonix,Biathlon,Biathlon,\"BERTHET, G.\",FRA,Men,Military Patrol,Bronze"
## [3] "1924,Chamonix,Biathlon,Biathlon,\"MANDRILLON, C.\",FRA,Men,Military Patrol,Bronze"
## [4] "1924,Chamonix,Biathlon,Biathlon,\"MANDRILLON, Maurice\",FRA,Men,Military Patrol,Bronze"
## [5] "1924,Chamonix,Biathlon,Biathlon,\"VANDELLE, André\",FRA,Men,Military Patrol,Bronze"
```

# Importing Data to R

Use the function read.csv() to load the file into R environment

```
winter <- read.csv('Data/winter.csv')
str(winter)
```

```
## 'data.frame':    5770 obs. of  9 variables:
##  $ Year      : int  1924 1924 1924 1924 1924 1924 1924 1924 1924 1924 ...
##  $ City      : chr  "Chamonix" "Chamonix" "Chamonix" "Chamonix" ...
##  $ Sport     : chr  "Biathlon" "Biathlon" "Biathlon" "Biathlon" ...
##  $ Discipline: chr  "Biathlon" "Biathlon" "Biathlon" "Biathlon" ...
##  $ Athlete   : chr  "BERTHET, G." "MANDRILLON, C." "MANDRILLON, Maurice" "VANDELLE, André" ...
##  $ Country   : chr  "FRA" "FRA" "FRA" "FRA" ...
##  $ Gender    : chr  "Men" "Men" "Men" "Men" ...
##  $ Event     : chr  "Military Patrol" "Military Patrol" "Military Patrol" "Military Patrol" ...
##  $ Medal     : chr  "Bronze" "Bronze" "Bronze" "Bronze" ...
```

Be sure that you can see the new object in the environment. Otherwise you cannot work with it

# Importing Data to R

- Look the help for read.csv (?read.csv), several options on how the file is imported.
- stringsAsFactors = F – the strings are loaded as a text rather than as a factor.

```r
winter <- read.csv('Data/winter.csv', stringsAsFactors = F)
str(winter)
```
```
## 'data.frame':    5770 obs. of  9 variables:
##  $ Year      : int  1924 1924 1924 1924 1924 1924 1924 1924 1924 1924 ...
##  $ City      : chr  "Chamonix" "Chamonix" "Chamonix" "Chamonix" ...
##  $ Sport     : chr  "Biathlon" "Biathlon" "Biathlon" "Biathlon" ...
##  $ Discipline: chr  "Biathlon" "Biathlon" "Biathlon" "Biathlon" ...
##  $ Athlete   : chr  "BERTHET, G." "MANDRILLON, C." "MANDRILLON, Maurice" "VANDELLE, André" ...
##  $ Country   : chr  "FRA" "FRA" "FRA" "FRA" ...
##  $ Gender    : chr  "Men" "Men" "Men" "Men" ...
##  $ Event     : chr  "Military Patrol" "Military Patrol" "Military Patrol" "Military Patrol" ...
##  $ Medal     : chr  "Bronze" "Bronze" "Bronze" "Bronze" ...
```

# Importing Data to R

- The same data is saved in tab delimited file.
- A tab-separated values (TSV) file is a simple text format for storing data in a tabular structure, e.g., database table or spreadsheet data, and a way of exchanging information between databases. Each record in the table is one line of the text file.

```r
readLines('Data/winter.txt', n=5)
```

```
## [1] "\"Year\"\t\"City\"\t\"Sport\"\t\"Discipline\"\t\"Athlete\"\t\"Country\"\t\"Gender\"\t\"Event\"\t\"Medal\""
## [2] "\"1\"\t1924\t\"Chamonix\"\t\"Biathlon\"\t\"Biathlon\"\t\"BERTHET, G.\"\t\"FRA\"\t\"Men\"\t\"Military Patrol\"\t\"Bronze\""
## [3] "\"2\"\t1924\t\"Chamonix\"\t\"Biathlon\"\t\"Biathlon\"\t\"MANDRILLON, C.\"\t\"FRA\"\t\"Men\"\t\"Military Patrol\"\t\"Bronze\""
## [4] "\"3\"\t1924\t\"Chamonix\"\t\"Biathlon\"\t\"Biathlon\"\t\"MANDRILLON, Maurice\"\t\"FRA\"\t\"Men\"\t\"Military Patrol\"\t\"Bronz
## [5] "\"4\"\t1924\t\"Chamonix\"\t\"Biathlon\"\t\"Biathlon\"\t\"VANDELLE, André\"\t\"FRA\"\t\"Men\"\t\"Military Patrol\"\t\"Bronze\""
```

# Importing Data to R

You can use the same read.csv() file but need to specify the separator (delimiter)

```
winter <- read.csv('Data/winter.txt', sep="\t", stringsAsFactors = FALSE)
str(winter)
```
```
## 'data.frame':    5770 obs. of  9 variables:
##  $ Year      : int  1924 1924 1924 1924 1924 1924 1924 1924 1924 1924 ...
##  $ City      : chr  "Chamonix" "Chamonix" "Chamonix" "Chamonix" ...
##  $ Sport     : chr  "Biathlon" "Biathlon" "Biathlon" "Biathlon" ...
##  $ Discipline: chr  "Biathlon" "Biathlon" "Biathlon" "Biathlon" ...
##  $ Athlete   : chr  "BERTHET, G." "MANDRILLON, C." "MANDRILLON, Maurice" "VANDELLE, André" ...
##  $ Country   : chr  "FRA" "FRA" "FRA" "FRA" ...
##  $ Gender    : chr  "Men" "Men" "Men" "Men" ...
##  $ Event     : chr  "Military Patrol" "Military Patrol" "Military Patrol" "Military Patrol" ...
##  $ Medal     : chr  "Bronze" "Bronze" "Bronze" "Bronze" ...
```

# Importing Data to R

R can work with other data files as well

- STATA files
- SPSS
- SAS
- MS Excel
- Databases
- XML
- HTML
- Etc.

Section 6

**Subsetting Data**

# Subsetting Data: Vectors

Create named vector

```
Club <- c("Juventus", "Napoli", "Roma", "Inter")
Points <- c(95,91,77,72)
names(Points) <- Club
Points
## Juventus   Napoli     Roma    Inter
##       95       91       77       72
```

# Subsetting Data: Vectors

- In R (unlike Python) the indexing starts with 1.
- Subsetting by index location.
- The vector has one dimension, so to subset it you need to specify the location of 1 index only.

```
Points[2]
```
```
## Napoli
##    91
```

Subset by name

```
Points["Juventus"]
```
```
## Juventus
##       95
```

# Subsetting Data: Vectors

- Subset several elements
- First two elements

```
Points[1:2]
```

```
## Juventus   Napoli
##      95       91
```

- 1 and 3 elements

```
Points[c(1,3)]
```

```
## Juventus     Roma
##      95       77
```

- Subsetting by name

```
Points[c("Napoli","Inter")]
```

```
## Napoli  Inter
##     91     72
```

# Subsetting Data: Vectors

Why R you getting an eRRoR ?

```
Club[1,2]
```
```
## Error in Club[1, 2]: incorrect number of dimensions
```

```
Club[2,4]
```
```
## Error in Club[2, 4]: incorrect number of dimensions
```

# Subsetting Data: Data frames

Load nba dataset

```
load('Data/nba2009_2018.rda')
summary(nba2009_2018)
```

```
##   SEASON_ID          GAME_DATE         home.TEAM_ABBREVIATION
## Length:12060       Min.   :2009-10-27   Length:12060
## Class :character   1st Qu.:2012-03-05   Class :character
## Mode  :character   Median :2014-11-13   Mode  :character
##                    Mean   :2014-08-10
##                    3rd Qu.:2017-01-24
##                    Max.   :2019-04-10
## home.TEAM_NAME       home.PTS      away.TEAM_ABBREVIATION away.TEAM_NAME
## Length:12060       Min.   :  0.0   Length:12060           Length:12060
## Class :character   1st Qu.: 95.0   Class :character       Class :character
## Mode  :character   Median :103.0   Mode  :character       Mode  :character
##                    Mean   :103.6
##                    3rd Qu.:112.0
##                    Max.   :161.0
##    away.PTS        home.WL
## Min.   :  0.0   Length:12060
## 1st Qu.: 92.0   Class :character
## Median :100.0   Mode  :character
## Mean   :100.9
## 3rd Qu.:109.0
## Max.   :168.0
```

# Subsetting Data: Data frames

Data frame has two dimensions: Rows (first dimension) and Columns (second dimension)

# dataframe[ , ]

Placeholder for rows

Placeholder for columns

# Subsetting Data: Data frames

Will return the value on first row and forth column

# dataframe[1, 4]

Will return the value of the first 10 rows and columns 4,5,6

# dataframe[1:10, 4:6]

# Subsetting Data: Data frames

Will return the first 10 rows and all columns

## dataframe[1:10, ]

If you want to select all elements for the given index, just leave the placeholder empty.

What will this command return?

## dataframe[, c(2,4:6)]

# Subsetting Data: Data frames

Indexing by column names

```
colnames(nba2009_2018)
```
```
## [1] "SEASON_ID"        "GAME_DATE"        "home.TEAM_ABBREVIATION"
## [4] "home.TEAM_NAME"   "home.PTS"         "away.TEAM_ABBREVIATION"
## [7] "away.TEAM_NAME"   "away.PTS"         "home.WL"
```

```
nba1 <- nba2009_2018[,c("home.PTS", "away.PTS")]
colnames(nba1)
```
```
## [1] "home.PTS" "away.PTS"
```

```
dim(nba1)
```
```
## [1] 12060     2
```

# Subsetting Data: Data frames

- Negative indexing is used to exclude certain records from the dataframe.
- This does not work with column names indexing.

```
nba2 <- nba2009_2018[,-c(1,2,4:6)]
head(nba2)
```

```
##   home.TEAM_ABBREVIATION      away.TEAM_NAME away.PTS home.WL
## 1                    CLE      Boston Celtics       95       L
## 2                    DAL  Washington Wizards      102       L
## 3                    POR     Houston Rockets       87       W
## 4                    LAL          LA Clippers      92       W
## 5                    ATL       Indiana Pacers      109       W
## 6                    ORL  Philadelphia 76ers      106       W
```

# Subsetting Data: Data frames

Exercises

- Create new data frame from nba dataset
- Include first 100 rows and columns 2, 3, 5
- Exclude rows 250, 300 to 350 and exclude column 5

# Subsetting Data: Data frames

You can access specific column in data frame by using dollar sign

```
mean(nba2009_2018$home.PTS)
```
```
## [1] 103.6102
```

```
table(nba2009_2018$home.WL)
```
```
##
##    L    W
## 4949 7110
```

# Subsetting Data: Data frames

Subset one variable

- First approach gives you a dataframe with one columns
- Second gives you a vector

```
home_pts <- nba2009_2018['home.PTS']
str(home_pts)
```
```
## 'data.frame':    12060 obs. of  1 variable:
## $ home.PTS: num  89 91 96 99 120 120 101 92 115 74 ...
```

```
home_pts <- nba2009_2018[,'home.PTS']
str(home_pts)
```
```
## num [1:12060] 89 91 96 99 120 120 101 92 115 74 ...
```

# Subsetting Data: Data frames

Conditional indexing

- Create new dataframe with games only from season 2009
- We need all the rows where the value for SEASON_ID is 2009
- Note that the type for column SEASON_ID is character

```
nba4 <- nba2009_2018[nba2009_2018$SEASON_ID=='2009',]
```

Check if everything is done right

```
table(nba4$SEASON_ID)
```

```
##
## 2009
## 1230
```

# Subsetting Data: Data frames

- Take only seasons 2009 and 2010
- As the SEASON_ID is character we will do the following

```
nba5 <- nba2009_2018[nba2009_2018$SEASON_ID %in% c("2009", "2010"),]
table(nba5$SEASON_ID)
```
```
##
## 2009 2010
## 1230 1230
```

# Subsetting Data: Data frames

- Workaround: make SEASON_ID numeric vector

```
nba2009_2018$SEASON_ID <- as.numeric(nba2009_2018$SEASON_ID)
nba5 <- nba2009_2018[nba2009_2018$SEASON_ID < 2011,]
table(nba5$SEASON_ID)
```
```
##
## 2009 2010
## 1230 1230
```

# Subsetting Data: Data frames

Logical operators in R

| Operator | Description |
|----------|-------------|
| < | less than |
| <= | less than or equal to |
| > | greater than |
| >= | greater than or equal to |
| == | exactly equal to |
| != | not equal to |
| !x | Not x |
| x \| y | x OR y |
| x & y | x AND y |

# Subsetting Data: Data frames

- Indexing on multiple conditions
- Take all the home game records for Detroit Pistons for seasons 2010 and 2011

```
nba6 <- nba2009_2018[nba2009_2018$SEASON_ID %in% c(2010,2011) &
                     nba2009_2018$home.TEAM_NAME =="Detroit Pistons",]

table(nba6$home.TEAM_NAME, nba6$SEASON_ID)
##
##                 2010 2011
##  Detroit Pistons   41   33
```

# Subsetting Data: Data frames

Using OR (|) operator in R

```
nba7 <- nba2009_2018[nba2009_2018$away.TEAM_NAME=="Detroit Pistons" |
                      nba2009_2018$home.TEAM_NAME =="Detroit Pistons",]
```

```
head(nba7)
```

```
##    SEASON_ID GAME_DATE home.TEAM_ABBREVIATION   home.TEAM_NAME home.PTS
## 10      2009 2009-10-28                  MEM Memphis Grizzlies       74
## 23      2009 2009-10-30                  DET   Detroit Pistons       83
## 36      2009 2009-10-31                  MIL   Milwaukee Bucks       96
## 54      2009 2009-11-03                  DET   Detroit Pistons       85
## 61      2009 2009-11-04                  TOR   Toronto Raptors      110
## 74      2009 2009-11-06                  ORL     Orlando Magic      110
##    away.TEAM_ABBREVIATION      away.TEAM_NAME away.PTS home.WL
## 10                    DET     Detroit Pistons       96       L
## 23                    OKC Oklahoma City Thunder      91       L
## 36                    DET     Detroit Pistons       85       W
## 54                    ORL       Orlando Magic       80       W
## 61                    DET     Detroit Pistons       99       W
## 74                    DET     Detroit Pistons      103       W
```

# Subsetting Data: Data frames

Adding new variable in dataframe
- Point differential

```
nba2009_2018$Diff <- nba2009_2018$home.PTS-nba2009_2018$away.PTS
hist(nba2009_2018$Diff)
```

**Histogram of nba2009_2018$Diff**

# Indexing List

```
(seriea <- data.frame(Club, Points))
```

```
##              Club Points
## Juventus Juventus     95
## Napoli     Napoli     91
## Roma         Roma     77
## Inter       Inter     72
```

- We have created a dataframe. Note that the dataframe has rownames
- Create a list

```
(list1 <- list(Club, "Italy", 2017, seriea))
```

```
## [[1]]
## [1] "Juventus" "Napoli"   "Roma"     "Inter"
##
## [[2]]
## [1] "Italy"
##
## [[3]]
## [1] 2017
##
## [[4]]
##              Club Points
## Juventus Juventus     95
## Napoli     Napoli     91
## Roma         Roma     77
## Inter       Inter     72
```

# Indexing List

- If you index the list with [], the result is a list.
- If you index the list with [[]] the result has the same structure as the list element

```
list1[4]
```
```
## [[1]]
##                Club Points
## Juventus Juventus    95
## Napoli     Napoli     91
## Roma         Roma     77
## Inter       Inter     72
```

```
list1[[4]]
```
```
##                Club Points
## Juventus Juventus    95
## Napoli     Napoli     91
## Roma         Roma     77
## Inter       Inter     72
```

```
is.data.frame(list1[[4]])
```
```
## [1] TRUE
```

```
is.data.frame(list1[4])
```
```
## [1] FALSE
```

# Indexing List

Two elements

```
list1[1:2]
```
```
## [[1]]
## [1] "Juventus" "Napoli"   "Roma"     "Inter"
##
## [[2]]
## [1] "Italy"
```

Slicing element from the element

This will bring the second element from the first element in the list

```
list1[[1]][2]
```
```
## [1] "Napoli"
```

# Indexing List

Named elements in the list

```
list1 <- list(Teams = Club, country = "Italy",
              year = 2017, standings = seriea)
list1
```

```
## $Teams
## [1] "Juventus" "Napoli"   "Roma"     "Inter"
##
## $country
## [1] "Italy"
##
## $year
## [1] 2017
##
## $standings
##               Club Points
## Juventus Juventus     95
## Napoli     Napoli     91
## Roma         Roma     77
## Inter       Inter     72
```

# Indexing List

Indexing by name

## `list1$standings`

```
##                Club Points
## Juventus Juventus    95
## Napoli     Napoli    91
## Roma         Roma    77
## Inter       Inter    72
```

## `list1$year`

```
## [1] 2017
```

Section 7

**Functions**

# Functions

- In programming, you use functions to incorporate sets of instructions that you want to use repeatedly or that, because of their complexity, are better self-contained in a sub program and called when needed.
- A function is a piece of code written to carry out a specified task; it can or can not accept arguments or parameters and it can or can not return one or more values.
- In fact, there are several possible formal definitions of 'function' spanning from mathematics to computer science.
- **Generically, its arguments constitute the input and return values as output.**

# Functions

- Functions have named arguments which potentially have default values.
- The formal arguments are the arguments included in the function definition.
- The formals function returns a list of all the formal arguments of a function.
- Not every function call in R makes use of all the formal arguments.
- Function arguments can be missing or might have default values.

# Functions

The function runif() generates random numbers from continuous uniform distribution

Arguments are:

- n, number of observations, with no default - so needs to be specified
- min, the minimum
- max, the maximum

```
formals(runif)
```
```
## $n
##
##
## $min
## [1] 0
##
## $max
## [1] 1
```

```
args(runif)
```
```
## function (n, min = 0, max = 1)
## NULL
```

# Functions

R functions' arguments can be matched **positionally** or **by name**. So the following calls to runif() are all equivalent (you can mix them as well).

- Matching by name

```r
runif(min=1, n=10, max=2)
```

```
## [1] 1.120291 1.039831 1.040716 1.710022 1.647759 1.442891 1.491432 1.471562
## [9] 1.741866 1.996648
```

- Matching by position

```r
runif(10,1,2)
```

```
## [1] 1.267439 1.770249 1.640040 1.426882 1.050750 1.019916 1.773276 1.147958
## [9] 1.733208 1.351730
```

Note, as the process is random you might get different results.

# Functions

Matched by position, this will give an error

```
runif(1,10,2)
```
```
## Warning in runif(1, 10, 2): NAs produced
## [1] NaN
```

## Functions

- If the argument has a default value and is not defined then the default value is used.
- Generated 10 numbers with min=0, max=1

```
runif(10)
```
```
## [1] 0.882818504 0.169275823 0.657059403 0.137945563 0.616535203 0.008407116
## [7] 0.115444378 0.709159585 0.393922719 0.049844107
```

# Functions

You can look at the arguments of the function and their description by hitting **Tab** inside the brackets

```
> runif()
```

| | |
|---|---|
| ⬧ n = | **n** |
| ⬧ min = | number of observations. If length(n) > 1, the length is taken to be the number required. |
| ⬧ max = | Press F1 for additional help |

# Functions

**General rule: if you have to copy a code more than twice, write a function.**

f <- function(a, b = 1, c = 2, d = NULL) {
what needs to be done
}

In addition to not specifying a default value, you can also set an argument value to NULL.

## Functions

- User-defined functions are stored in the global environment and can be accessed easily.
- The following function will calculate the x power of y (default value of y is 2).

```
foo <- function (x, y=2){
  x^y
}

foo(4)
## [1] 16
```

## Functions

Arguments to functions are evaluated lazily, so they are evaluated only as needed.

```
foo <- function (x,y){
  x^2
}

foo(3)
## [1] 9
```

This function never actually uses the argument y, so calling f(2) will not produce an error because the 2 gets positionally matched to x.

# Functions

- The return value of a function is the last expression in the function body to be evaluated.
- It can be also specified with the function return()

# Functions

The following function will calculate the z-score of a vector

```
norm <- function(x) {
  return((x-mean(x))/sd(x))
}

norm(mtcars$mpg)
## [1]  0.15088482  0.15088482  0.44954345  0.21725341 -0.23073453 -0.33028740
## [7] -0.96078893  0.71501778  0.44954345 -0.14777380 -0.38006384 -0.61235388
## [13] -0.46302456 -0.81145962 -1.60788262 -1.60788262 -0.89442035  2.04238943
## [19]  1.71054652  2.29127162  0.23384555 -0.76168319 -0.81145962 -1.12671039
## [25] -0.14777380  1.19619000  0.98049211  1.71054652 -0.71190675 -0.06481307
## [31] -0.84464392  0.21725341
```

Section 8

**Dates and Times in R**

# Dates in R

- R provides several options for dealing with date and date/time data.
- The built in function as.Date() handles dates (without times)
- Library chron handles dates and times, but does not control for time zones
- The POSIXct and POSIXlt classes allow for dates and times with control for time zones.
- The general rule for date/time data in R is to use the simplest technique possible.
- Thus, for date only data, as.Date() will usually be the best choice. If you need to handle dates and times without time-zone information, the chron library is a good choice; the POSIX classes are especially useful when time-zone manipulation is important.

# Dates in R

- Except for the **POSIXlt** class, dates are stored internally as the number of days or seconds from some reference date.
- The function as.Date() stores the number of days passed from 1970-01-01

```r
x <- as.Date("2018-06-18")
class(x)
```
```
## [1] "Date"
```

```r
as.numeric(x)
```
```
## [1] 17700
```

```r
as.numeric(as.Date("1970-01-01"))
```
```
## [1] 0
```

# Dates in R

- The date can come in different formats.
- Usually, you need to tell R in which format the date is.

| Symbol | Meaning | Example |
|--------|---------|---------|
| **%d** | day as a number (0-31) | 01-31 |
| **%a** | abbreviated weekday | Mon |
| **%A** | unabbreviated weekday | Monday |
| **%m** | month (00-12) | 00-12 |
| **%b** | abbreviated month | Jan |
| **%B** | unabbreviated month | January |
| **%y** | 2-digit year | 07 |
| **%Y** | 4-digit year | 2007 |

```
as.Date("2018/01/15", format = "%Y/%m/%d")
```
## [1] "2018-01-15"

```
as.Date("01-15-2018", format = "%m-%d-%Y")
```
## [1] "2018-01-15"

# Dates in R

```
oil <- read.csv("Data/oil.csv", stringsAsFactors = FALSE)
head(oil)
```

```
##       DATE   OPEN   HIGH    LOW  CLOSE   VOL
## 1 1-Oct-12 112.14 113.27 110.76 111.40 80055
## 2 2-Oct-12 111.40 111.70 110.55 110.55 29332
## 3 3-Oct-12 110.55 110.59 106.95 107.34 56307
## 4 4-Oct-12 107.44 111.79 107.24 111.36 61664
## 5 5-Oct-12 111.27 112.09 109.64 111.13 51704
## 6 7-Oct-12 111.11 111.13 111.02 111.02    57
```

```
str(oil)
```

```
## 'data.frame':    1643 obs. of  6 variables:
##  $ DATE : chr  "1-Oct-12" "2-Oct-12" "3-Oct-12" "4-Oct-12" ...
##  $ OPEN : num  112 111 111 107 111 ...
##  $ HIGH : num  113 112 111 112 112 ...
##  $ LOW  : num  111 111 107 107 110 ...
##  $ CLOSE: num  111 111 107 111 111 ...
##  $ VOL  : int  80055 29332 56307 61664 51704 57 40638 61077 68836 68708 ...
```

# Dates in R

Format the DATE column with as.Date

```
oil$DATE<-as.Date(oil$DATE, format="%d-%b-%y")
str(oil)
## 'data.frame':    1643 obs. of  6 variables:
## $ DATE : Date, format: "2012-10-01" "2012-10-02" ...
## $ OPEN : num  112 111 111 107 111 ...
## $ HIGH : num  113 112 111 112 112 ...
## $ LOW  : num  111 111 107 107 110 ...
## $ CLOSE: num  111 111 107 111 111 ...
## $ VOL  : int  80055 29332 56307 61664 51704 57 40638 61077 68836 68708 ...
```

# Dates in R

- Now we can subset the data frame with DATE column.
- The following code is going to create a data frame with cases only from 2014

```
oil_2014 <- oil[oil$DATE >= "2014-01-01" & oil$DATE < "2015-01-01",]
head(oil_2014)
```

```
##          DATE   OPEN   HIGH    LOW  CLOSE   VOL
## 389 2014-01-02 110.74 111.07 107.49 107.65 32977
## 390 2014-01-03 107.66 108.42 106.44 106.65 35605
## 391 2014-01-06 106.68 107.60 106.28 106.67 35795
## 392 2014-01-07 106.73 107.33 106.56 107.01 30686
## 393 2014-01-08 107.08 107.54 106.60 106.83 36728
## 394 2014-01-09 106.83 107.71 105.70 106.00 55777
```

# Dates in R

You can do arithmetic operations with class Date

```
as.Date("2018-06-18")+1
```
## [1] "2018-06-19"

```
as.Date("2018-06-18")-as.Date("2018-05-18")
```
## Time difference of 31 days

# Dates in R

Extracting weekday

```
oil$WEEKDAY <- weekdays(oil$DATE)
head(oil)
##         DATE   OPEN   HIGH    LOW  CLOSE   VOL   WEEKDAY
## 1 2012-10-01 112.14 113.27 110.76 111.40 80055    Monday
## 2 2012-10-02 111.40 111.70 110.55 110.55 29332   Tuesday
## 3 2012-10-03 110.55 110.59 106.95 107.34 56307 Wednesday
## 4 2012-10-04 107.44 111.79 107.24 111.36 61664  Thursday
## 5 2012-10-05 111.27 112.09 109.64 111.13 51704    Friday
## 6 2012-10-07 111.11 111.13 111.02 111.02    57    Sunday
```

# Dates in R

Extract month

```
oil$MONTH <- months(oil$DATE)
head(oil)
```

```
##          DATE   OPEN   HIGH    LOW  CLOSE   VOL   WEEKDAY   MONTH
## 1 2012-10-01 112.14 113.27 110.76 111.40 80055    Monday October
## 2 2012-10-02 111.40 111.70 110.55 110.55 29332   Tuesday October
## 3 2012-10-03 110.55 110.59 106.95 107.34 56307 Wednesday October
## 4 2012-10-04 107.44 111.79 107.24 111.36 61664  Thursday October
## 5 2012-10-05 111.27 112.09 109.64 111.13 51704    Friday October
## 6 2012-10-07 111.11 111.13 111.02 111.02    57    Sunday October
```

# Dates in R

- To extract the day of the month, use format()
- The returned column is character, but you can make it numeric

```r
oil$DAY <- format(oil$DATE, "%d")
str(oil)
## 'data.frame':    1643 obs. of  9 variables:
## $ DATE   : Date, format: "2012-10-01" "2012-10-02" ...
## $ OPEN   : num  112 111 111 107 111 ...
## $ HIGH   : num  113 112 111 112 112 ...
## $ LOW    : num  111 111 107 107 110 ...
## $ CLOSE  : num  111 111 107 111 111 ...
## $ VOL    : int  80055 29332 56307 61664 51704 57 40638 61077 68836 68708 ...
## $ WEEKDAY: chr  "Monday" "Tuesday" "Wednesday" "Thursday" ...
## $ MONTH  : chr  "October" "October" "October" "October" ...
## $ DAY    : chr  "01" "02" "03" "04" ...
```

# Dates in R

Make the variable numeric

```
oil$DAY <- as.numeric(oil$DAY)
head(oil)
```

```
##          DATE   OPEN   HIGH    LOW  CLOSE   VOL   WEEKDAY   MONTH DAY
## 1 2012-10-01 112.14 113.27 110.76 111.40 80055    Monday October   1
## 2 2012-10-02 111.40 111.70 110.55 110.55 29332   Tuesday October   2
## 3 2012-10-03 110.55 110.59 106.95 107.34 56307 Wednesday October   3
## 4 2012-10-04 107.44 111.79 107.24 111.36 61664  Thursday October   4
## 5 2012-10-05 111.27 112.09 109.64 111.13 51704    Friday October   5
## 6 2012-10-07 111.11 111.13 111.02 111.02    57    Sunday October   7
```

Section 9

**R Markdown**

# R Markdown

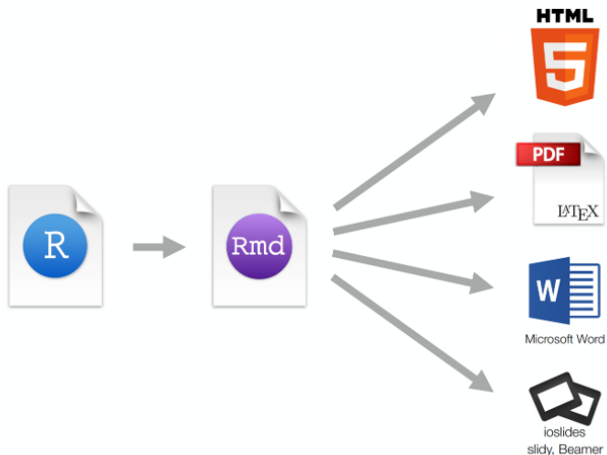- R Markdown documents are fully reproducible.
- It allows to combine text and code together in one report.
- You can use multiple languages within markdown file (R, python, SQL).
- R Markdown supports dozens of static and dynamic output formats including HTML, PDF, MS Word, Beamer, HTML5 slides, etc.
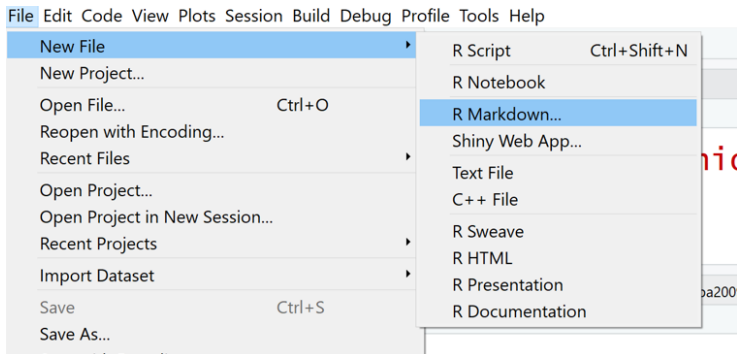- R markdown is run in its own environment.

# R Markdown

The Workflow

# R Markdown

To open new markdown file

# R Markdown

The document below is a template R Markdown document. It includes the most familiar parts of an R Markdown document:

- A YAML header that contains some metadata (1)
- Narrative text written in Markdown (2)
- R code chunks surrounded by triple backticks {r} and followed by triple backticks; a syntax that comes from the knitr package (3)

```
1  ---
2  title: "Example"
3  author: "Habet Madoyan"
4  date: "June 16, 2018"
5  output: pdf_document
6  ---
7
8  An R markdown document with some text and code for Data Science class
9
0
1  ```{r}
2  head(mtcars, n=5)
3  ```
```

# R Markdown

Labeling and reusing code chunks

- Apart from the popular code chunk options you have learned by now, you can define even more things in the curly braces that follow the triple backticks.
- An interesting feature available in knitr is the labeling of code snippets. The code chunk below would be assigned the label simple_sum:

```
```{r simple_sum, results = 'hide'}
2 + 2
```
```

- However, because the results option is equal to hide, no output is shown.

# R Markdown

- You can embed R code into the text of your document with the r syntax. Be sure to include the lower case r in order for this to work properly.
- R Markdown will run the code and replace it with its result, which should be a piece of text, such as a character string or a number.
- For example, the line below uses embedded R code to create a complete sentence:

```
The factorial of four is `r factorial(4)`.
```

- When you render the document the result will appear as:

The factorial of four is 24.

- Inline code provides a useful way to make your reports completely automated.

# R Markdown

**LaTeX equations**

- You can also use the Markdown syntax to embed latex math equations into your reports. To embed an equation in its own centered equation block, surround the equation with two pairs of dollar signs like this,

`$$1 + 1 = 2$$`

- To embed an equation inline, surround it with a single pair of dollar signs, like this:

`$1 + 1 = 2$`

- You can use all of the standard latex math symbols to create attractive equations.

# R Markdown

LaTex formula example

$$\frac{n!}{k!(n-k)!} = \binom{n}{k}$$

In Tex

```
\frac{n!}{k!(n-k)!} = \binom{n}{k}
```

Tex script needs to be written in the text portion of the document

# R Markdown

**Lists in R Markdown**

To make a bullet list in Markdown, place each item on a new line after an asterisk or hyphen and a space, like this:

```
* item 1
* item 2
* item 3
```

You can make an ordered list by placing each item on a new line after a number followed by a period followed by a space, like this

```
1. item 1
2. item 2
3. item 3
```

In each case, you need to place a blank line between the list and any paragraphs that come before it.

# R Markdown

**R code chunks**

You can embed R code into your R Markdown report with the knitr syntax. To do this, surround your code with two lines: one that contains triple backticks {r} and the below one contains triple backticks. The result is a code chunk that looks like this:

```r
```{r}
# some code
```
```

When you render the report, R will execute the code. If the code returns any results, R will add them to your report.

# R Markdown

Customize R code chunks

- You can customize each R code chunk in your report by providing optional arguments after the r in "'{r}, which appears at the start of the code chunk. Let's look at one set of options.
- R functions sometimes return messages, warnings, and even error messages. By default, R Markdown will include these messages in your report. You can use the **message, warning and error** options to prevent R Markdown from displaying these. If any of the options are set to FALSE, R Markdown will not include the corresponding type of message in the output. Packages often generate messages when you first load them with library().
- For example, R Markdown would ignore any messages or warnings generated by the chunk below.

```r
```{r, warning = FALSE, message = FALSE}
library(dplyr)
```
```

# R Markdown

Three of the most popular chunk options are echo, eval and results.

- If **echo = FALSE**, R Markdown will not display the code in the final document (but it will still run the code and display its results unless told otherwise).
- If **eval = FALSE**, R Markdown will not run the code or include its results, (but it will still display the code unless told otherwise).
- If **results = 'hide'**, R Markdown will not display the results of the code (but it will still run the code and display the code itself unless told otherwise).

# R Markdown

Other important stuff

Emphasis

```
*italic*   **bold**

_italic_   __bold__
```

Headers

```
# Header 1

## Header 2

### Header 3
```

If you want to start from a new page in pdf file, put in text part of the markdown

```
\newpage

\pagebreak
```

Section 10

**Data Transformation: dplyr**

# Grammar of data manipulation

dplyr is a grammar of data manipulation, providing a consistent set of verbs that help you solve the most common data manipulation challenges:

- **mutate()** adds new variables that are functions of existing variables
- **select()** picks variables based on their names.
- **filter()** picks cases based on their values.
- **summarise()** reduces multiple values down to a single summary.
- **arrange()** changes the ordering of the rows.

These all combine naturally with group_by() which allows you to perform any operation "by group".

# Grammar of data manipulation

- Not necessary, but dplyr works the best with pipe like operator from **magrittr** package
- %>% operator takes the object from its left-hand side and use it as an argument in the function on the right-hand side

%>%
magrittr

*Ceci n'est pas un pipe.*

# Grammar of data manipulation

- The "pipe" operation is a handy tool to make your code more legible: %>%.

Key points:

- It takes the output of your previous operation and uses it as an input to your next operation.
- You can determine where the previous argument goes with the period symbol, . , which acts as a placeholder.
- Understand how to use it by replacing the pipe operation with "then" (in your mind, not in the code). For example, **filter(data, …) %>% select(…)** filters first then selects columns from the output of filter.

# Grammar of data manipulation

Filtering USA games only

```
summer <- read.csv("Data/summer.csv", stringsAsFactors = F)

summer_usa <- summer %>% filter(Country=="USA")

table(summer_usa$Country, summer_usa$Medal)

##
##      Bronze Gold Silver
##  USA   1098 2235   1252
```

# Grammar of data manipulation

Filter and group_by

```
summer %>%
  filter(Country %in% c("USA", "FRA", "GBR")) %>%
  group_by(Country) %>%
  summarise(Count=n())
## # A tibble: 3 x 2
##   Country Count
##   <chr>   <int>
## 1 FRA      1396
## 2 GBR      1720
## 3 USA      4585
```

- The result is a tibble
- Count=n() creates a new variable named Count with frequencies, n() calculates frequencies

# Grammar of data manipulation

Number of medals by country (is this long or wide format?)

```
summer %>%
  filter(Country %in% c("USA", "FRA", "GBR")) %>%
  group_by(Country, Medal) %>%
  summarise(Count=n())
```

```
## # A tibble: 9 x 3
## # Groups:   Country [3]
##   Country Medal  Count
##   <chr>   <chr>  <int>
## 1 FRA     Bronze   497
## 2 FRA     Gold     408
## 3 FRA     Silver   491
## 4 GBR     Bronze   553
## 5 GBR     Gold     546
## 6 GBR     Silver   621
## 7 USA     Bronze  1098
## 8 USA     Gold    2235
## 9 USA     Silver  1252
```

# Grammar of data manipulation

Using ggplot with dplyr

```
summer %>%
  filter(Country %in% c("USA", "FRA", "GBR")) %>%
  group_by(Country, Medal) %>%
  summarise(Count=n())
```

```
## # A tibble: 9 x 3
## # Groups:   Country [3]
##   Country Medal  Count
##   <chr>   <chr>  <int>
## 1 FRA     Bronze   497
## 2 FRA     Gold     408
## 3 FRA     Silver   491
## 4 GBR     Bronze   553
## 5 GBR     Gold     546
## 6 GBR     Silver   621
## 7 USA     Bronze  1098
## 8 USA     Gold    2235
## 9 USA     Silver  1252
```