



IPS Instituto
Politécnico de Setúbal
Escola Superior de
Tecnologia de Setúbal

Session 2

EDA and Data Cleaning in R



Course Name:

Supervised Machine Learning

Copyright © 2019-2020 Vala A. Rohani

By Prof. Vala A. Rohani
vala.ali.rohani@estsetubal.ips.pt

Session outline

1. Reviewing the essential topics in R Programming
2. Exploratory Data Analysis in R
3. Preparation and Data Cleaning for ML (Machine Learning) in R



Reviewing the essential topics in R Programming

R is a programming language for **statistical analysis** and **data visualization**. Using R, you can take raw data and understand the trends and patterns in it.

You can also use R to build and validate the machine learning models.



RStudio is an integrated development environment (IDE) for R.



1. Installing R

Latest version: 4.2.3 for windows

<https://cran.r-project.org/bin/windows/base/>

2. Installing R Studio

RStudio Desktop 2023.03.0+386

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

2: Install RStudio

DOWNLOAD RSTUDIO DESKTOP FOR WINDOWS

Size: 208.08 MB | [SHA-256: 885432DB](#) | Version: 2023.03.0+386 |

Released: 2023-03-16

RStudio environment [you can change the theme through **Tools > Global Option > Appearance**]

The screenshot displays the RStudio interface with the following components:

- Source Editor:** Contains R code for calculating distances and applying functions to data frames. The code includes comments and function definitions for `transportcostDP` and `transportcostCU`.
- Environment Pane:** Shows the global environment with variables `dfcus` (305 obs. of 4 variables), `dfDPS` (2513 obs. of 4 variables), and `dfWHS` (34 obs. of 4 variables). It also lists values for `ncu`, `ndp`, and `nwh`, and functions `transportcostCU` and `transportcostDP`.
- Files Pane:** Shows the file explorer with a list of folders including `.Rhistory`, `Custom Office Templates`, `Python Scripts`, and `R`.
- Console:** Displays the output of the R script, showing the execution of the `ncu` variable and the application of the `transportcostDP` and `transportcostCU` functions to the data frames.

First things first!

Create the course folder in your laptop:

/IPS-ESCE SP ML

/Session 2

/Session 3

/Session 4

/Session 5

/Session 6

Your first line of code in R!

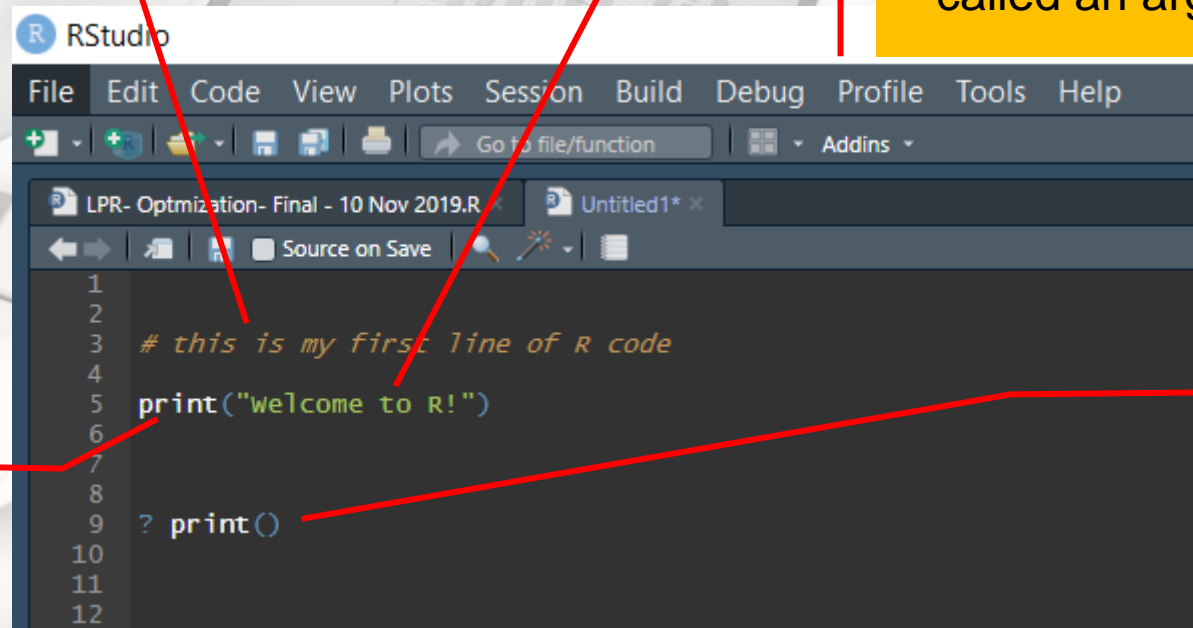
Ctrl+ENTER to run each line

Comments helps to make your codes more readable

"Welcome to R!" is called an argument

? Print() shows help about print() function

Print() is a function

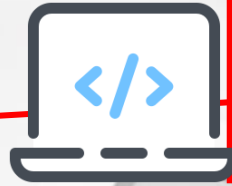


The screenshot shows the RStudio interface with a script editor containing the following code:

```
1  
2  
3 # this is my first line of R code  
4 print("Welcome to R!")  
5  
6  
7  
8  
9 ? print()  
10  
11  
12
```

Red lines from the surrounding text boxes point to specific parts of the code: from the 'Comments' box to line 3; from the '"Welcome to R!"' box to the string argument in line 4; from the '? Print()' box to line 9; and from the 'Print() is a function' box to the function name in line 9.

Variables



Your turn! Create a variable named "aSentence" and store a sentence in it

```
12 # in R, you can store data in a variable by using "<-" . You can name a  
13 # variable anything you want as long as it's not already the name of something else.  
14 # I find that a short phrase (without spaces) is generally best.  
15  
16 textToPrint <- "this is some text to print"  
17  
18 # if you give R the name of a variable, it will print whatever is in that variable  
19  
20 textToPrint  
21  
22 # note that capitalization does matter! this line will generate an error because  
23 # there is nothing called "texttoprint"  
24  
25 texttoprint  
26  
27 ##### Your turn! Create a variable named "aSentence" and store a sentence in it  
28  
29 # our old friend print()  
30 print(textToPrint)  
31  
32 # the nchar() function tells you the number of characters in a variable  
33 nchar(textToPrint)  
34  
35 # the c() function concatenates (strings together) all its arguments  
36 c(textToPrint, textToPrint, textToPrint)  
37  
38
```



Lab Activity:

Your turn! Try the following three exercises on your own.

- 1) print the variable "**aSentence**" you made previously
- 2) find out the number of characters in "**aSentence**"
- 3) concatenate the "**textToPrint**" and "**aSentence**" variables

Data Types

R has 6 basic data types:

- **character:** "a", "swc"
- **numeric:** 2, 15.5
- **integer:** 2L (the L tells R to store this as an integer)
- **logical:** TRUE, FALSE
- **complex:** 1+4i (complex numbers with real and imaginary parts)

Data Types

```
37  
38 # what we've seen so far are characters. This is the type of data you'll use for text  
39 varStr <- "someText"  
40  
41 # we can check the data type of a variable using the function str() (like "structure")  
42 str(varStr)  
43 # we can tell this is a character because it's structure is "chr"  
44
```



Your turn!
Check the data type of the "aSentence" variable you made above.

Data Types

```
51  
52 # let's create some numeric variables  
53 hoursPerDay <- 24  
54 daysPerWeek <- 7  
55  
56 # we can check to make sure that these actually are numeric  
57 class(hoursPerDay)  
58 class(daysPerWeek)  
59  
60 # since this is numeric data, we can do math with it!  
61 # "*" is the symbol for multiplication  
62 hoursPerWeek <- hoursPerDay * daysPerWeek  
63 hoursPerWeek  
64
```



Your turn!

Create a numeric variable "minutesPerHour" and use it to calculate a new variable called "minutesPerWeek" that has the number of minutes per week in it

Data Types



Your turn!

Now try this:

```
a <- 5
```

```
b <- "6"
```

```
a * b
```



```
75  
76 # You can change character data to numeric data using the as.numeric() function.  
77 # This will let you do math with it again. :)  
78 a * as.numeric(b)  
79  
80 a * b  
81  
82  
83 # check out the structure: note that b changes from "chr" to "num"  
84 str(b)  
85 str(as.numeric(b))  
86  
87 # to fix b to be a number permanently  
88 b <- as.numeric(b)  
89  
90 str (b)  
91
```


Data Types

So far we've learned about two data types: **character** and **numeric**. But there's a third common data type you'll encounter a lot in R: **logical or boolean** data. Booleans can only take two values, TRUE and FALSE.

In [15]:

```
91  
92 #--- Boolean types  
93  
94 # You'll get a boolean back if you ask R "are these two things the same?" using "=="  
95  
96 var1 <- "a" == "b"  
97 var2 <- 1 == 1  
98  
99 var1  
100 var2  
101  
102 str (var1)  
103
```



Your turn!

First, take a guess: will `6 == "6"` return TRUE or FALSE?
Then test your prediction. Are you surprised by the outcome?
What does this tell you about datatypes in R?

Vectors

In programming, a vector is **list of data that is all of the same data type.**



```
104
105 # vectors
106
107 # let's make a vector!
108 listOfNumbers <- c(1,5,91,42.8,100008.41)
109 listOfNumbers
110
111 str (listOfNumbers)
112
113 # because this is a numeric vector, we can do math on it! When you do math to a vector,
114 # it happens to every number in the vector. (If you're familiar with matrix
115 # multiplication, it's the same thing as multiplying a 1x1 matrix by a 1xN matrix.)
116
117 # multiply every number in the vector by 5
118 5 * listOfNumbers
119
120
121 # add one to every number in the vector
122 listOfNumbers + 1
123
124 listOfNumbers
125 |
126
```

Your turn!

divide every number in the vector
listOfNumbers by 2 and store it in v2.

[In R, you do division using the / symbol]

Vectors

How to access the elements of vector?

You can look inside a vector using square brackets (`[]`). The number inside the square brackets tells R what **index** to look at. *In R, indexes start at one.*



Your turn!

1. Print the sum of first and second item of "listOfNumbers" and,
2. Check if the first item of "listOfNumbers" is bigger than the second one or not?

```
129  
130 # get the third item from "listOfNumbers"  
131 listOfNumbers[3]  
132  
133 # and store it in another variable  
134  
135 varA <- listOfNumbers[3]  
136  
137 varA  
138  
139
```

Managing Data Frames

The **data frame** is a key data structure in statistics and in R. The basic structure of a data frame is that there is one observation **per row** and each **column** represents a variable of that observation.

The **dplyr** package is designed to provide a highly optimized set of routines specifically for dealing with data frames.

So first,
Let's see how to install
packages in R

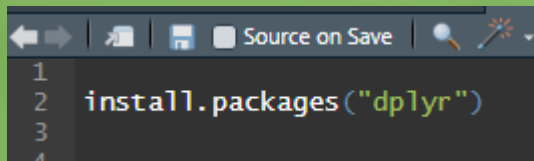


Managing Data Frames

Packages are collections of R functions, data, and compiled code in a well-defined format.

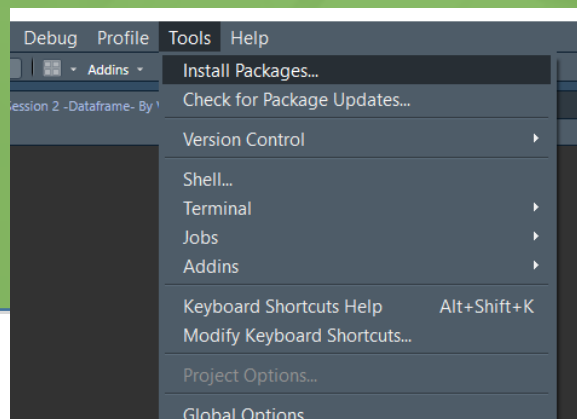
The directory where packages are stored is called the library. R comes with a standard set of packages. Others are available for download and installation. Once installed, they have to be loaded into the session to be used.

How to install a new package?

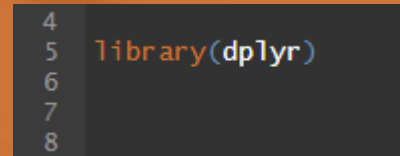


```
1  
2 install.packages("dplyr")  
3  
4
```

Or,



How to use an installed package?



```
4  
5 library(dplyr)  
6  
7  
8
```

Managing Data Frames

dplyr Grammar:

Some of the key “verbs” provided by the dplyr package are:

select: return a subset of the columns of a data frame, using a flexible notation

filter: extract a subset of rows from a data frame based on logical conditions

arrange: reorder rows of a data frame

rename: rename variables in a data frame

mutate: add new variables/columns or transform existing variables

summarise / summarize: generate summary statistics of different variables in the data frame, possibly within strata

%>%: the “pipe” operator is used to connect multiple verb actions together into a pipeline

Managing Data Frames

Now, we need to read the data :

For the examples in this chapter we will be using a dataset containing air pollution and temperature data for the city of Chicago in the U.S. The dataset is available in the shared folder.

1. Setting working directory:

```
7  
8 # ----- Setting the working directory  
9  
10 # First let's see the current working directory  
11  
12 getwd()  
13  
14 # now, we can set our desired directory as the working directory  
15  
16 setwd("D:/Academia/IPS/2019 Analytics for Master Students/Training Material/Seasion 2/data/chicago_data")  
17  
18 # let's check it  
19  
20 getwd()  
21
```

2. Reading data into the dataframe

```
22  
23  
24 dfChicago <- readRDS("chicago.rds")  
25  
26
```

Managing Data Frames

Let's see some basic characteristics of the dataset with the **dim()** and **str()** functions.

```
26  
27  
28 dim(dfChicago)  
29  
30 str (dfChicago)  
31  
32
```

31:1 (Top Level) ▾

Console Terminal × Jobs ×

D:/Academia/IPS/2019 Analytics for Master Students/Training Material/Seasion 2/data/chicago_data/ ↗

```
> dim(dfChicago)  
[1] 6940    8  
> str (dfChicago)  
'data.frame':   6940 obs. of  8 variables:  
 $ city      : chr  "chic" "chic" "chic" "chic" ...  
 $ tmpd      : num  31.5 33 33 29 32 40 34.5 29 26.5 32.5 ...  
 $ dptp      : num  31.5 29.9 27.4 28.6 28.9 ...  
 $ date      : Date, format: "1987-01-01" "1987-01-02" "1987-01-03" "1987-01-04" ...  
 $ pm25tmean2: num  NA NA NA NA NA NA NA NA NA NA ...  
 $ pm10tmean2: num  34 NA 34.2 47 NA ...  
 $ o3tmean2  : num  4.25 3.3 3.33 4.38 4.75 ...  
 $ no2tmean2 : num  20 23.2 23.8 30.4 30.3 ...  
> |
```

Managing Data Frames

The **select()** function can be used to select columns of a data frame that you want to focus on. Often you'll have a large data frame containing "all" of the data, but any given analysis might only use a subset of variables or observations. The **select()** function allows you to get the few columns you might need.

```
32 #--- select
33
34 names(dfChicago)
35
36 # -- Results: "city"      "tmpd"      "dptp"      "date"      "pm25tmean2" "pm10tmean2" "o3tmean2" "no2tmean2"
37
38 dfSubset <- select(dfChicago, city:dptp)
39
40 names(dfSubset)
41
42 # -- Results: "city" "tmpd" "dptp"
43
44 head(dfSubset)
45
46
```

You can also **omit variables** using the **select()** function by using the negative sign.

```
47 # -- Omit columns
48
49 select(dfChicago, -(city:dptp))
50
```

Managing Data Frames



Your turn!

1. Write a line of code to select columns ("city", "date", and "pm10tmean2") and store them in df3
2. Can you implement it with pipe (%>%)

Managing Data Frames

The **filter()** function is used to extract subsets of rows from a data frame.

```
51  
52 #--- Filter  
53  
54 dfFiltered <- filter(dfChicago, pm25tmean2 > 30)  
55  
56 head(dfFiltered)  
57  
58 dfFiltered2 <- filter(dfChicago, pm25tmean2 > 30 & tmpd > 80)  
59  
60 head(dfFiltered2)  
61  
62
```

Managing Data Frames

The **arrange()** function is used to sort rows of a data frame according to one of the variables/columns.

```
62  
63 #---- Sort  
64  
65 dfsorted <- arrange(dfChicago,date)  
66  
67 dfsorted2 <- arrange(dfChicago, desc (date))  
68  
69
```

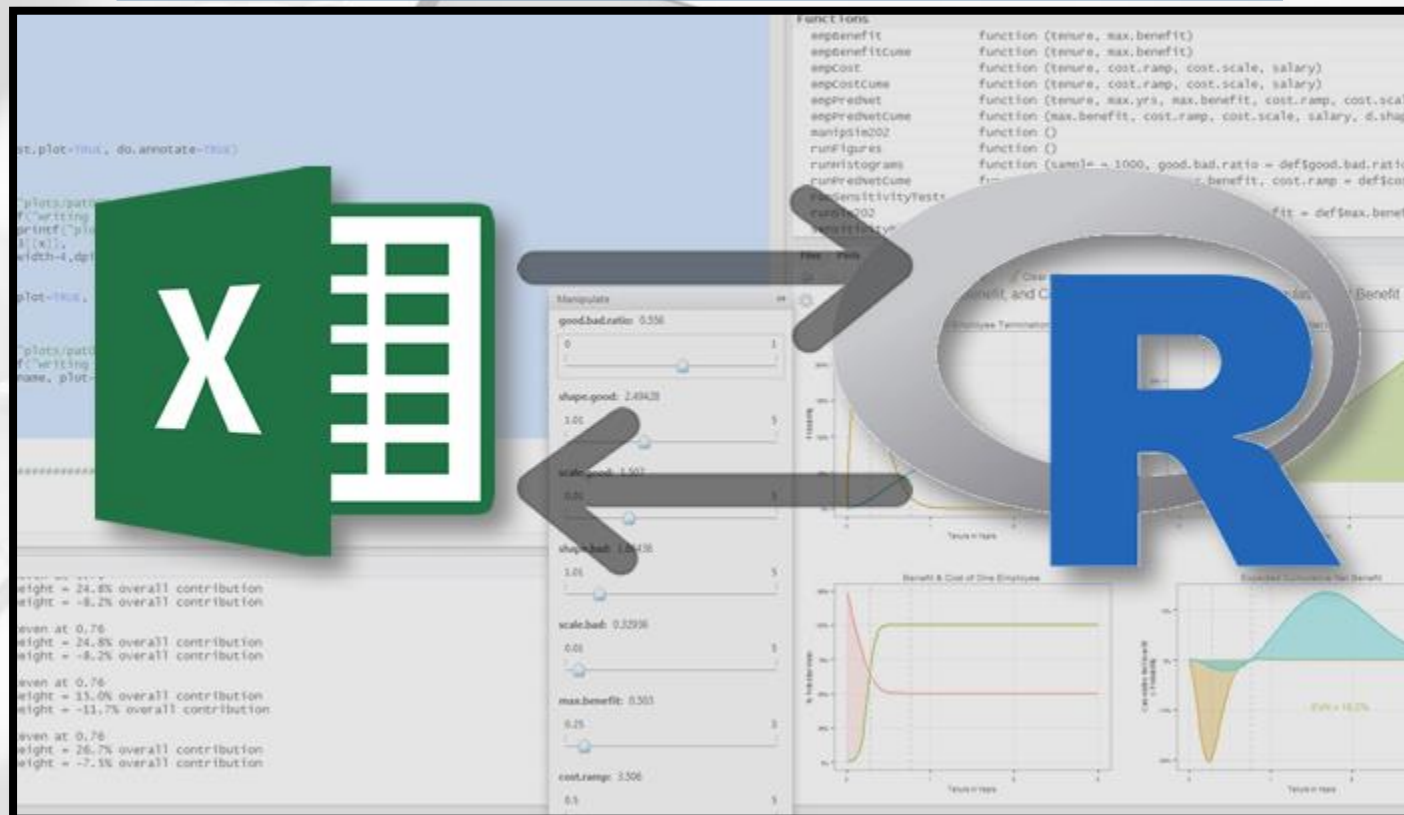


Your turn!

Check the functions **rename()** and **mutate()** from dplyr package.
What is their usage? Apply them on dfChicago

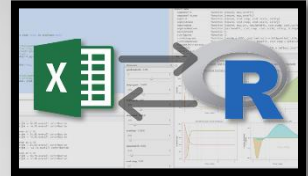
Managing Data Frames

How to communicate with excel?



Managing Data Frames

To communicate with excel you need **readxl** package, so let's first install it:



```
73
74 # ----- Read and write to Excel/CSV
75
76 install.packages("readxl")
77
78 library("readxl")
79
80
81 # now, we can set our desired directory as the working directory
82
83 setwd("D:/Academia/IPS/2019 Analytics for Master Students/Training Material/Seasion 2/data/chocolate-bar-ratings")
84
85 getwd()
86
87 #-- reading excel file
88
89 dfCacao <- read_excel("flavors_of_cacao.xlsx")
90
91
92 #--- Analysing the data
93
94
95 #--- Storing the dada in csv format
96
97 write.csv(dfCacao, file = "flavors_of_cacao_editted.csv")
98
99
```



Lab assignment:

1. Read the flavors_of_cacao. Xlsx into the dfCacao data frame.
2. How many columns and rows does it have? What are the column names? What is the datatype of each column?
3. Sort the dataset descending based on Ratings. Store them in separate data frame .
4. How many records exists in Portugal, Spain, Germany, and USA? Store them in separate data frame .
5. How many records exists in Switzerland in 2010 with rating of 3.5? Store them in separate data frame .
6. Store the records of above countries (in 4th step) in separate data frames and compare the average ratings for each of them.
7. Store all generated datafarmes in excel format.

Data Cleaning > Graphs Principles > Graphs in R

Data cleaning is one of the most important aspects of data science.

As a data scientist, you can expect to spend up to **80% of your time cleaning data**.

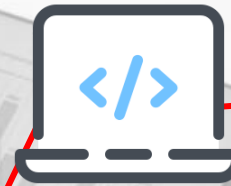
In this post you'll learn how to detect missing values using the [tidyr](#) and [dplyr](#) packages from the [Tidyverse](#).

Data Cleaning > Graphs Principles > Graphs in R

NA Stands for not available. NA is a placeholder for a missing value.

what is NA?

```
NA + 1
sum(c(NA, 1, 2))
median(c(NA, 1, 2, 3), na.rm = TRUE)
length(c(NA, 2, 3, 4))
3 == NA
NA == NA
TRUE | NA
```

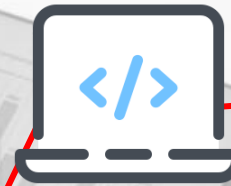


Your turn!

Try the following codes and see the results

Data Cleaning > Graphs Principles > Graphs in R

NULL means no class (its class is NULL) and has length 0 so it does not take up any space in a vector.



Your turn!

Try the following codes and see the results

```
# what is NULL?
```

```
length(c(1, 2, NULL, 4))
```

```
sum(c(1, 2, NULL, 4))
```

```
x <- NULL
```

```
c(x, 2)
```


Data Cleaning > Graphs Principles > Graphs in R

In statistics, missing data, or missing values, occur when no data value is stored for the variable in an observation. Missing data are a common occurrence and can have a significant effect on the conclusions that can be drawn from the data.

Missing Values

Statistical analysis error

```
39  
40 # ---- Statistical analysis error  
41  
42  
43 age <- c(23, 16, NA)  
44 mean(age)  
45  
46 ## [1] NA  
47  
48 mean(age, na.rm = TRUE)  
49 ## [1] 19.5  
50
```

Data Cleaning > Graphs Principles > Graphs in R

Identify and remove the NA value

Missing Values

```
51  
52 # --- Identify the NA value  
53  
54  
55 complete.cases(age)  
56  
57 #or  
58  
59 is.na(age)  
60  
61  
62 # --- Remove NA values  
63  
64 na.omit(age)  
65  
66 age  
67
```

Data Cleaning > Graphs Principles > Graphs in R

recode missing values with the mean

```
68  
69 #--- Recode missing values with mean  
70  
71 x <- c(1:5, NA, 9:11, NA)  
72 is.na(x)  
73  
74 df <- data.frame(col1 = c(1:3, NA),  
75                   col2 = c("this", NA, "is", "text"),  
76                   col3 = c(TRUE, FALSE, TRUE, TRUE),  
77                   col4 = c(2.5, 4.2, 3.2, NA),  
78                   stringsAsFactors = FALSE)  
79 is.na(df)  
80  
81 # -- To identify the location of the NA.  
82  
83  
84 which(is.na(x))  
85 sum(is.na(df))  
86  
87 #or for data frame  
88  
89 colSums(is.na(df))  
90  
91  
92 # recode missing values with the mean  
93 x  
94  
95 x[is.na(x)] <- mean(x, na.rm = TRUE)  
96  
97 # do it for dataframe  
98  
99 df  
100  
101 df$col4[is.na(df$col4)] <- mean(df$col4, na.rm = TRUE)  
102  
103
```

Missing Values

Data Cleaning > Graphs Principles > Graphs in R

showing complete and incomplete observations

Missing Values

```
103  
104 #-- For complete observations.  
105  
106  
107 df  
108  
109 complete.cases(df)  
110  
111  
112 # subset with complete.cases to get complete cases  
113 df[complete.cases(df), ]  
114  
115 # or use na.omit() to get same as above  
116 na.omit(df)  
117  
118 # or subset with `!` operator to get incomplete cases  
119 df[!complete.cases(df), ]  
120  
121
```

What can we do with missing values in datasets?

Missing Values

1. Deleting the observations

- Have sufficient data points, so the model doesn't lose power.
- Not to introduce bias (meaning, disproportionate or non-representation of classes).

2. Deleting the variable

- When we have lots of missing values for a specific variable

3. Imputation with mean / median / mode

- Replacing the missing values with the mean / median / mode is a crude way of treating missing values. Depending on the context, like if the variation is low or if the variable has low leverage over the response, such

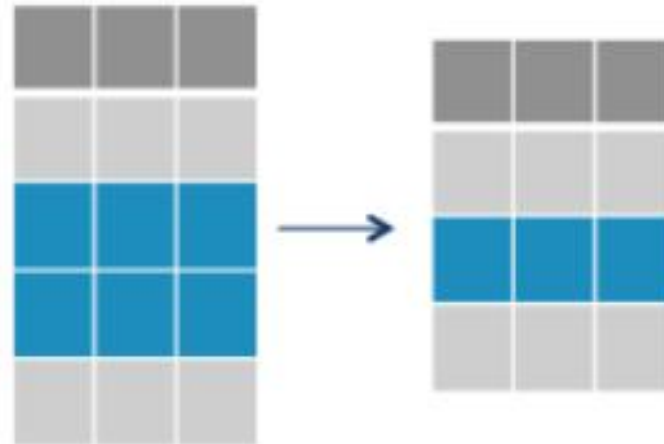
4. Prediction

- Prediction is most advanced method to impute your missing values and includes different approaches such as: kNN Imputation, rpart, and mice.

Data Cleaning > Graphs Principles > Graphs in R

Duplicate values

Remove Duplicate Data in R



`duplicated()`: Identify duplicate elements (R base)

`unique()`: Keep only unique elements (R base)

`distinct()`: Efficient solution to remove duplicate in a data table (dplyr)

Data Cleaning > Graphs Principles > Graphs in R

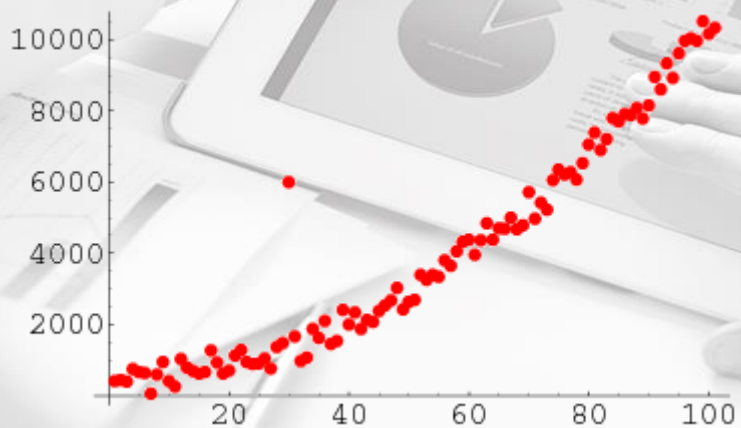
Duplicate values

```
124
125 ##### Removing duplicates #####
126
127 #--- Given the following vector:
128
129 x <- c(1, 1, 4, 5, 4, 6)
130
131 #--- To find the position of duplicate elements in x, use this:
132
133 duplicated(x)
134
135 # --- Extract duplicate elements:
136
137 x[duplicated(x)]
138
139 #--- If you want to remove duplicated elements, use !duplicated()
140
141 x[!duplicated(x)]
142
143 x <- x[!duplicated(x)]
144
145 x
146
147
148 #-- Following this way, you can remove duplicate rows from a data frame based on a column values, as follow:
149
150 library(tidyverse)
151
152 my_data <- as_tibble(iris)
153 my_data
154
155 # Remove duplicates based on Sepal.Width columns
156
157 my_data[!duplicated(my_data$Sepal.Width), ]
158
159
160 # --- Extract unique elements
161
162 x <- c(1, 1, 4, 5, 4, 6)
163
164 unique(x)
165
166 unique(my_data)
167
168
169 #--- Remove duplicate rows in a data frame
170
171
172 my_data %>% distinct()
173
174 my_data %>% distinct(Sepal.Length, .keep_all = TRUE)
175
```

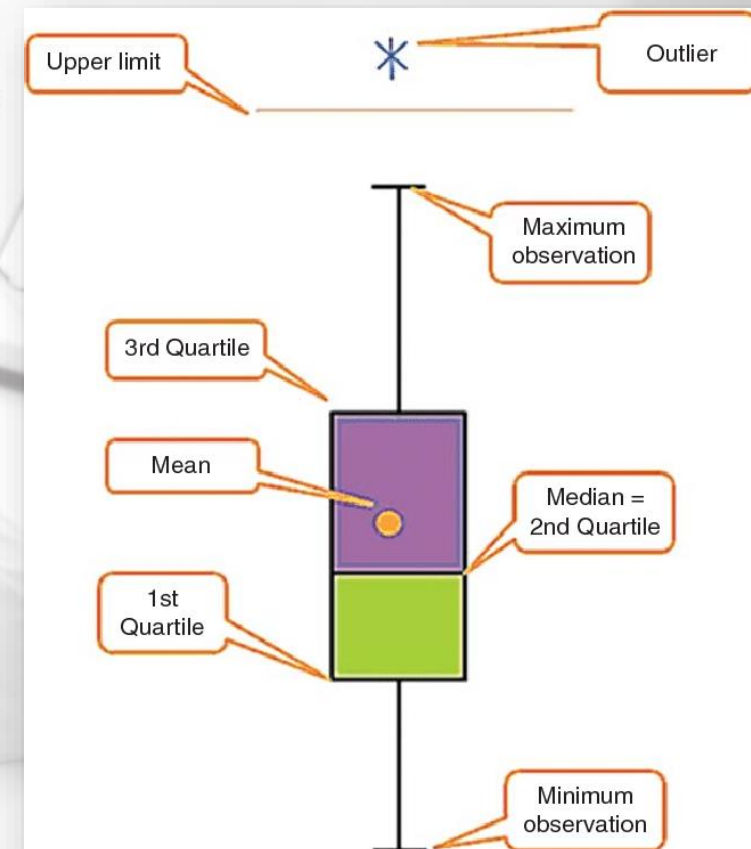

Data Cleaning > Graphs Principles > Graphs in R

In statistics, an **outlier** is a data point that differs significantly from other observations.

An outlier may be due to variability in the measurement or it may indicate experimental error. An outlier can cause serious problems in statistical analyses.



Outliers



Data Cleaning > Graphs Principles > Graphs in R

In this part, we use mtcars dataset. It comes with the base package, so no need to import anything)

Outliers

mtcars {datasets}

R Documentation

Motor Trend Car Road Tests

Description

The data was extracted from the 1974 *Motor Trend* US magazine, and comprises fuel consumption and 10 aspects of automobile design and performance for 32 automobiles (1973–74 models).

Usage

`mtcars`

Format

A data frame with 32 observations on 11 (numeric) variables.

[, 1] mpg Miles/(US) gallon
[, 2] cyl Number of cylinders
[, 3] disp Displacement (cu.in.)
[, 4] hp Gross horsepower
[, 5] drat Rear axle ratio
[, 6] wt Weight (1000 lbs)
[, 7] qsec 1/4 mile time
[, 8] vs Engine (0 = V-shaped, 1 = straight)
[, 9] am Transmission (0 = automatic, 1 = manual)
[,10] gear Number of forward gears
[,11] carb Number of carburetors

Data Cleaning > Graphs Principles > Graphs in R

Outliers

```
177
178 ▾ ##### Removing outliers #####
179
180 # First of all, we insert a couple of outliers to the $disp column of the mtcars dataset
181 # (mtcars comes with the base package, so no need to import anything)
182 # In order to have a couple of outliers in this dataset, we simply multiply the values in mtcars$disp that are higher than 420 by *2
183
184 mtcars$disp[which(mtcars$disp >420)] <- c(mtcars$disp[which(mtcars$disp >420)]*2)
185
186 # (This is just a random way of inserting a couple of outlier values, you could also assign a couple of high values in a million different
187
188 # Now we have a look at $disp column of the mtcars dataset with boxplot
189
190 boxplot(mtcars$disp)
191
192 # You can get the actual values of the outliers with this
193
194 boxplot(mtcars$disp)$out
195
196 # Now you can assign the outlier values into a vector
197
198 outliers <- boxplot(mtcars$disp, plot=FALSE)$out
199
200 # Check the results
201
202 print(outliers)
203
204
205 ##### Removing the outliers
206
207 mtcars[which(mtcars$disp %in% outliers),]
208
209
210 # Now you can remove the rows containing the outliers, one possible option is:
211
212 mtcars <- mtcars[-which(mtcars$disp %in% outliers),]
213
214 # If you check now with boxplot, you will notice that those pesky outliers are gone
215
216 boxplot(mtcars$disp)
217
```

Data Cleaning > Graphs Principles > Graphs in R



Lab Activity:

Data cleaning the Singapore Airbnb dataset

You can download the raw data in the shared folder :

/Session 3/data

Source:

<https://www.kaggle.com/jojoker/singapore-airbnb>



Data Cleaning > Graphs Principles > Graphs in R



Lab Activity:

Data cleaning the Singapore Airbnb dataset

idroom id
nameroom names
host_idhost id
host_namehost names
neighbourhood_groupSingapore regions
neighbourhoodspecific place
latitudelatitude
longitudelongitude
room_typeroom type

pricesingapore dollar per night
minimum_nightsminimum nights
number_of_reviewsnumber of review
last_reviewlast review
reviews_per_monthI don't know exactly
calculated_host_listings_counttotal room or house in host
catalog on Airbnb
availability_365availability



Lab Activity #1:

Data cleaning the Singapore Airbnb dataset

Apply all required data cleaning techniques that you learned in this session to clean the Singapore Airbnb dataset.

Then submit the following files:

- your source code (name: Session3-[your name]) in R
- The cleaned excel file

Some useful sources for further reading:

1. Exploratory Data Analysis with R

<https://bookdown.org/rdpeng/exdata/>

2. Programming with R

<https://swcarpentry.github.io/r-novice-inflammation/>

Exploratory Data Analysis with R



Roger D. Peng

Any
questions ?

Vala@data-corner.com

