



R: BASIC CONCEPTS

Eduardo Ogasawara
eogasawara@ieee.org
<https://eic.cefet-rj.br/~eogasawara>



Introduction to R

- R is a programming language and free software environment for statistical computing
 - Supported by the R Foundation for Statistical Computing
- Created by Ross Ihaka and Robert Gentleman at Auckland University, New Zealand
- R was derived by S (Bell Laboratories - AT&T)
- R is a language broadly used by statisticians, data miners, and data scientists

R Console

```
R version 4.0.0 (2020-04-24) -- "Arbor Day"
Copyright (C) 2020 The R Foundation for Statistical Computing
Platform: x86_64-apple-darwin17.0 (64-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY.
You are welcome to redistribute it under certain conditions.
Type 'license()' or 'licence()' for distribution details.

Natural language support but running in an English locale

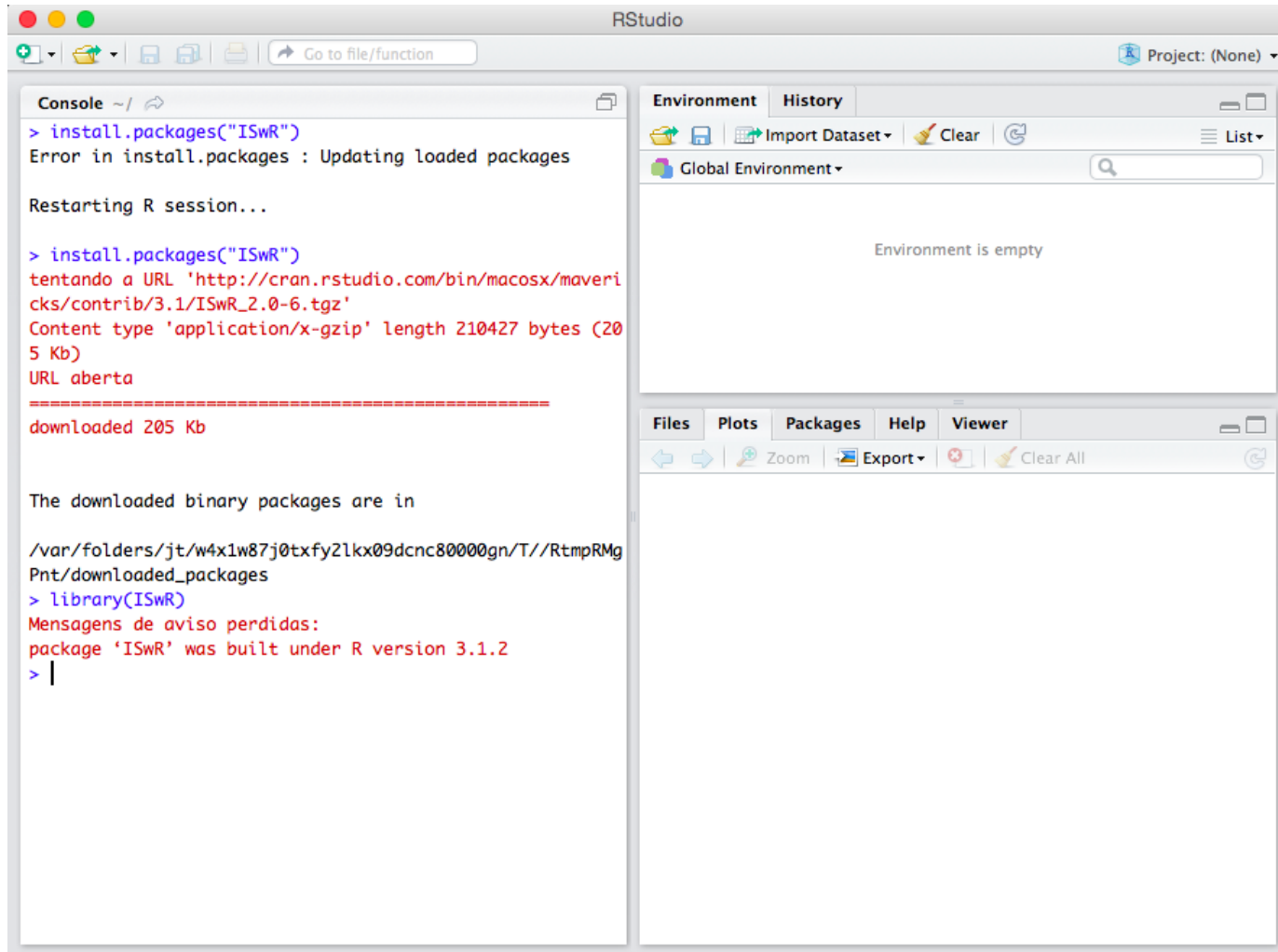
R is a collaborative project with many contributors.
Type 'contributors()' for more information and
'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.

During startup - Warning messages:
1: Setting LC_CTYPE failed, using "C"
2: Setting LC_COLLATE failed, using "C"
3: Setting LC_TIME failed, using "C"
4: Setting LC_MESSAGES failed, using "C"
5: Setting LC_MONETARY failed, using "C"
[R.app GUI 1.71 (7827) x86_64-apple-darwin17.0]

WARNING: You're using a non-UTF8 locale, therefore only ASCII characters will work.
Please read R for Mac OS X FAQ (see Help) section 9 and adjust your system preferences accordingly.
> |
```

R Studio



Exercise

- Take your time to install R and R Studio
 - <https://www.r-project.org>
 - <http://www.rstudio.com>

Basic concepts

- Assignment
- Value display
- Logical test
- Vector definition
 - Computing Body Mass Index (BMI)
- Printing values

```
x <- 2 # variable assignment
x # variable evaluation
is.numeric(x) # variable

weight = c(60, 72, 57, 90, 95, 72) # vector with six observations
height = c(1.75, 1.80, 1.65, 1.90, 1.74, 1.91)

bmi = weight/height^2
print(bmi)

print(sprintf("%.2f +/- %.2f", mean(bmi), sd(bmi)))
```

2

TRUE

```
[1] 19.59184 22.22222 20.93664 24.93075 31.37799 19.73630
[1] "23.13 +/- 4.49"
```

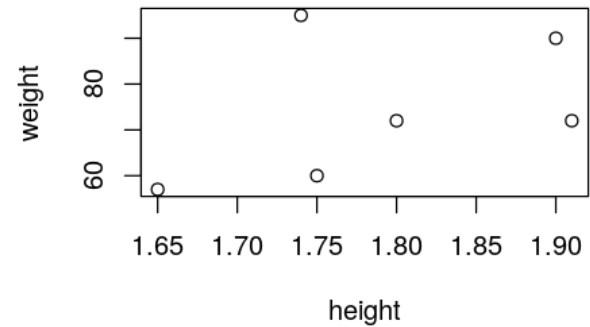
Exercise

- Knowing that standard deviation of a sample is
 - $\sigma(X) = \sqrt{\frac{\sum_{i=1}^n (x_i - \mu(X))^2}{n-1}}$
- using basic R commands:
 - compute the standard deviation of BMI using loops
 - compute the standard deviation of BMI using vector operations

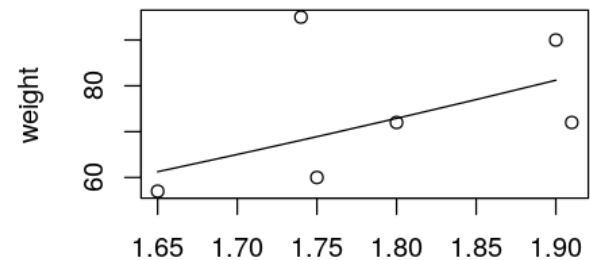
Plotting graphics

- Plotting a scatter graphics
 - Canvas is active until the next plot

```
plot(height, weight)
```



```
plot(height, weight)  
hh <- c(1.65, 1.70, 1.75, 1.80, 1.85, 1.90)  
lines(hh, 22.5 * hh^2)
```



Default arguments and help for functions

- Functions have default values
- View parameters of the function
- Use online help

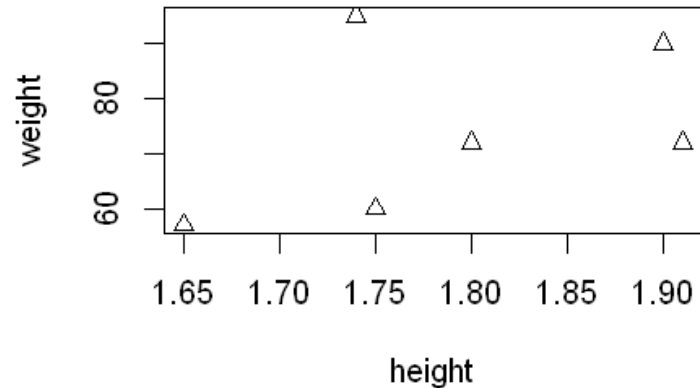
```
plot(height, weight, pch=2)
```

```
args(plot.default)
```

```
?graphics::plot
```

```
function (x, y = NULL, type = "p", xlim = NULL, ylim = NULL,  
         log = "", main = NULL, sub = NULL, xlab = NULL, ylab = NULL,  
         ann = par("ann"), axes = TRUE, frame.plot = axes, panel.first = NULL,  
         panel.last = NULL, asp = NA, ...)
```

```
NULL
```



Statistical analysis

- Test theoretical value of BMI equals to 22.5
 - Null hypothesis: no difference observed (p-value > 5%)
 - Alternative hypothesis: they are different

```
t.test(bmi, mu=22.5)
```

One Sample t-test

```
data: bmi
t = 0.34488, df = 5, p-value = 0.7442
alternative hypothesis: true mean is not equal to 22.5
95 percent confidence interval:
 18.41734 27.84791
sample estimates:
mean of x
 23.13262
```

More about vectors

- Operations with NA
- Name of observations
- Scalar multiplication

```
x <- c(A=1, B=NA, C=3)
```

```
mean(x)
```

```
mean(x, na.rm=TRUE)
```

```
names(x)
```

```
x["B"] <- 2
```

```
x["B"]*x
```

```
<NA>
```

```
2
```

```
'A' 'B' 'C'
```

A	2
B	4
C	6

Strings in R

- A string is a sequence of characters
- It is enclosed inside single quotes('This is a string') or inside the double quotes("This is also a string")
- It is internally represented in double-quotes

Strings

```
1
2 weight = c(60, 72, 57, 90, 95, 72) # vector with six observations
3
4 height = c(1.75, 1.80, 1.65, 1.90, 1.74, 1.91)
5
6 subject = c("A", "B", "C", "D", "E", "F")
7
8
9 mensagem <- rep("", length(subject))
10 for (i in 1:length(subject)) {
11   mensagem[i] <- sprintf(
12     "The height of %s is %.2f. The weight of %s is %d.",
13     subject[i], height[i], subject[i], weight[i])
14 }
15
16 print(mensagem)
17
18 rm(mensagem)
19
20 mensagem <- sprintf(
21   "The height of %s is %.2f. The weight of %s is %d.",
22   subject, height, subject, weight)
23 print(mensagem)
24
```

Factors

- Factors are variables in R that refer to categorical data
- Factors in R are stored as a vector of integer values with a corresponding set of character values to use when the factor is displayed
- Both numeric and character variables can be made into factors, but a factor's levels are always character values

Factors

```
pain = c(0,3,2,2,1)
fpain = factor(pain,levels=0:3)
levels(fpain) = c("none","mild","medium","severe")
```

```
fpain
```

```
as.numeric(fpain)
```

```
levels(fpain)
```

```
none severe medium medium mild
```

► **Levels:**

```
1 4 3 3 2
```

```
'none' 'mild' 'medium' 'severe'
```

Matrix

- Creation
- Creation by rows
- Names for rows and columns
- Transpose
- Determinant

```
m <- 1:9
dim(m) <- c(3,3)
m

mb <- matrix(1:9, nrow=3,byrow=TRUE)
rownames(mb) = LETTERS[1:3]
mb

t(m)

m*x

det(m)
```

```
1 4 7
```

```
2 5 8
```

```
3 6 9
```

```
A 1 2 3
```

```
B 4 5 6
```

```
C 7 8 9
```

```
1 2 3
```

```
4 5 6
```

```
7 8 9
```

```
1 4 7
```

```
4 10 16
```

```
9 18 27
```

```
0
```


Matrix

- Matrix is a two-dimensional data structure in R programming

Exercise

- Knowing that the determinant of a matrix is computed by function `det`, compute the determinant for the following matrix

$$\begin{bmatrix} 1 & 9 & 5 \\ 3 & 7 & 8 \\ 10 & 4 & 2 \end{bmatrix}$$

- Multiplying a Vector by a Matrix
 - Knowing that `%*%` computes matrix multiplication

- Find Ay , where $A = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix}$, and $y = \begin{bmatrix} 2 \\ 1 \\ 3 \end{bmatrix}$

Lists

- Lists are the R objects which contain elements of different types, such as numbers, strings, vectors, matrix, data frame, and another list inside it
- A list can also contain a matrix or a function as its elements
- A list is created using the `list()` function
- List manipulation
 - Slicing a list `[]`
 - Accessing an element inside a list `[[[]]`

Lists

```
x = c(5260, 5470, 5640, 6180, 6390,  
      6515, 6805, 7515, 7515, 8230, 8770)  
y = c(3910, 4220, 3885, 5160, 5645,  
      4680, 5265, 5975, 6790, 6900, 7335)  
  
lst <- list(A=x, B=y)  
  
lst  
  
lst$A
```

\$A

5260 5470 5640 6180 6390 6515 6805 7515 7515 8230 8770

\$B

3910 4220 3885 5160 5645 4680 5265 5975 6790 6900 7335

5260 5470 5640 6180 6390 6515 6805 7515 7515 8230 8770

Exercise

- Create a list (main list) with two lists inside
 - The first list contains all elements present in vector `intake.pre` of previous slide
 - The second list contains all elements present in vector `intake.pos` of previous slide
 - Filter the first list inside the main list such elements inside them are between 6000 and 7000

Data frames

- A data frame is a table where each column corresponds to attributes, and each row corresponds to a tuple (object)

Data frames

```
d <- data.frame(A=1st$A,B=1st$B)
d
df <- d[d$A > 7000 | d$A < 6000,]
df
```

	A	B
	5260	3910
	5470	4220
	5640	3885
	6180	5160
	6390	5645
	6515	4680
	6805	5265
	7515	5975
	7515	6790
	8230	6900
	8770	7335

	A	B
1	5260	3910
2	5470	4220
3	5640	3885
8	7515	5975
9	7515	6790
10	8230	6900
11	8770	7335

Exercise

- Create a data frame from the weight and height
- Filter the data frame such that all tuples are between 1.70 and 1.90
- Computes the bmi for filtered data frame

Implicitly Loops – sapply, lapply

- lapply, sapply executes a function for each column
 - The first character defines the return type
 - l – list, s – simple (vector or matrix)
 - The second parameter is the function to invoke
 - Following parameters are passed to the invoked function
- apply is the generic function
 - The second parameter defines if it calls the function for each row (1) or each column (2)

Implicitly Loops – *sapply*, *lapply*

```
lapply(d, min, na.rm=TRUE)
sapply(d, min, na.rm=TRUE)
apply(d, 1, min)
apply(d, 2, min)
```

SA

5260

SB

3885

A 5260

B 3885

3910 4220 3885 5160 5645 4680 5265 5975 6790 6900 7335

A 5260

B 3885

Exercise

- Create a data frame for weight and height and using apply, compute the bmi
 - Hint: Use a custom function for this goal

Sort and order

```
sort(d$B)
o <- order(d$B)
o
ds <- d[o,]
ds
```

3885 3910 4220 4680 5160 5265 5645 5975 6790 6900 7335

3 1 2 6 4 7 5 8 9 10 11

	A	B
3	5640	3885
1	5260	3910
2	5470	4220
6	6515	4680
4	6180	5160
7	6805	5265
5	6390	5645
8	7515	5975
9	7515	6790
10	8230	6900
11	8770	7335

Exercise

- Create a data frame from the weight and height
- Sort it according to the weight

Loading and saving files

```
wine = read.table("http://archive.ics.uci.edu/ml/machine-learning-databases/wine/wine.data",
  header = TRUE, sep = ",")
head(wine)
save(wine, file="wine.RData")

rm(wine)

load("wine.RData")
write.table(wine, file="wine.csv", row.names=FALSE, quote = FALSE)
```

X1	X14.23	X1.71	X2.43	X15.6	X127	X2.8	X3.06	X.28	X2.29	X5.64	X1.04	X3.92	X1065
1	13.20	1.78	2.14	11.2	100	2.65	2.76	0.26	1.28	4.38	1.05	3.40	1050
1	13.16	2.36	2.67	18.6	101	2.80	3.24	0.30	2.81	5.68	1.03	3.17	1185
1	14.37	1.95	2.50	16.8	113	3.85	3.49	0.24	2.18	7.80	0.86	3.45	1480
1	13.24	2.59	2.87	21.0	118	2.80	2.69	0.39	1.82	4.32	1.04	2.93	735
1	14.20	1.76	2.45	15.2	112	3.27	3.39	0.34	1.97	6.75	1.05	2.85	1450
1	14.39	1.87	2.45	14.6	96	2.50	2.52	0.30	1.98	5.25	1.02	3.58	1290

Exercise

- Create a data frame from the weight and height
- Save it in Rdata
- Save it in CSV
- Save it in RDS
- Read saved RDS
- Read saved Rdata
- Read saved CSV

Creating functions

```
create_dataset <- function() {  
  data <- read.table(text = "Year Months Flights Delays  
    2016 Jan-Mar 11 6  
    2016 Apr-Jun 12 5  
    2016 Jul-Sep 13 3  
    2016 Oct-Dec 12 5  
    2017 Jan-Mar 10 4  
    2017 Apr-Jun 9 3  
    2017 Jul-Sep 11 4  
    2017 Oct-Dec 25 15  
    2018 Jan-Mar 14 3  
    2018 Apr-Jun 12 5  
    2018 Jul-Sep 13 3  
    2018 Oct-Dec 15 4",  
    header = TRUE, sep = "")  
  data$OnTime <- data$Flights - data$Delays  
  data$Perc <- round(100 * data$Delays / data$Flights)  
  return(data)  
}  
  
data <- create_dataset()  
head(data)
```

Year	Months	Flights	Delays	OnTime	Perc
2016	Jan-Mar	11	6	5	55
2016	Apr-Jun	12	5	7	42
2016	Jul-Sep	13	3	10	23
2016	Oct-Dec	12	5	7	42
2017	Jan-Mar	10	4	6	40
2017	Apr-Jun	9	3	6	33

Pipelines

```
loadlibrary("dplyr")

data_sd <- create_dataset() %>%
  select(variable=Months, value=Delays) %>%
  group_by(variable) %>%
  summarize(sd = sd(value), value = mean(value))

data_sd$variable <- factor(data_sd$variable,
  levels = c('Jan-Mar', 'Apr-Jun', 'Jul-Sep', 'Oct-Dec'))

head(data_sd)
```

variable	sd	value
Apr-Jun	1.1547005	4.333333
Jan-Mar	1.5275252	4.333333
Jul-Sep	0.5773503	3.333333
Oct-Dec	6.0827625	8.000000

Exercise

- Create a data frame from the weight and height
- Filter the data frame such that all tuples are between 1.70 and 1.90
 - Use function `dplyr::filter`

Melt function

```
loadlibrary("reshape")
data <- create_dataset()
head(data)
data <- melt(data[,c('Year', 'Months', 'Flights', 'Delays', 'OnTime', 'Perc')],
             id.vars = c(1,2))
head(data)
```

A data.frame: 6 × 6

Year	Months	Flights	Delays	OnTime	Perc
<int>	<fct>	<int>	<int>	<int>	<dbl>
2016	Jan-Mar	11	6	5	55
2016	Apr-Jun	12	5	7	42
2016	Jul-Sep	13	3	10	23
2016	Oct-Dec	12	5	7	42
2017	Jan-Mar	10	4	6	40
2017	Apr-Jun	9	3	6	33

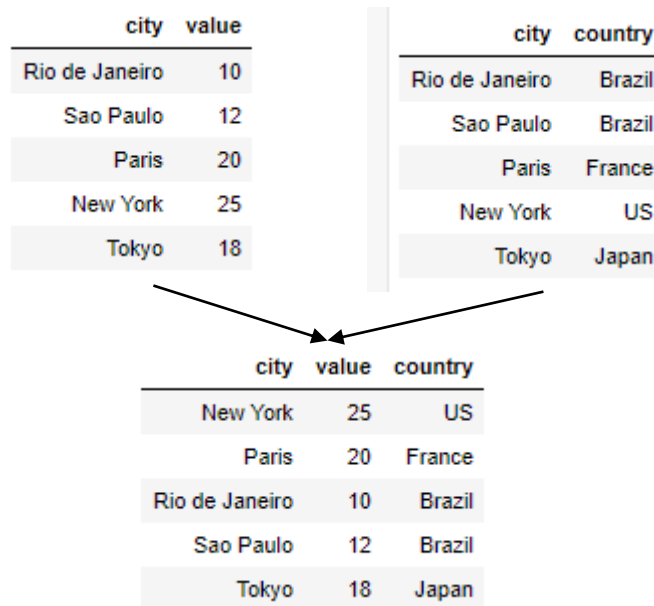
A data.frame: 6 × 4

Year	Months	variable	value
<int>	<fct>	<fct>	<dbl>
2016	Jan-Mar	Flights	11
2016	Apr-Jun	Flights	12
2016	Jul-Sep	Flights	13
2016	Oct-Dec	Flights	12
2017	Jan-Mar	Flights	10
2017	Apr-Jun	Flights	9

The **melt** function transforms columns values into rows grouped by **id.vars**.

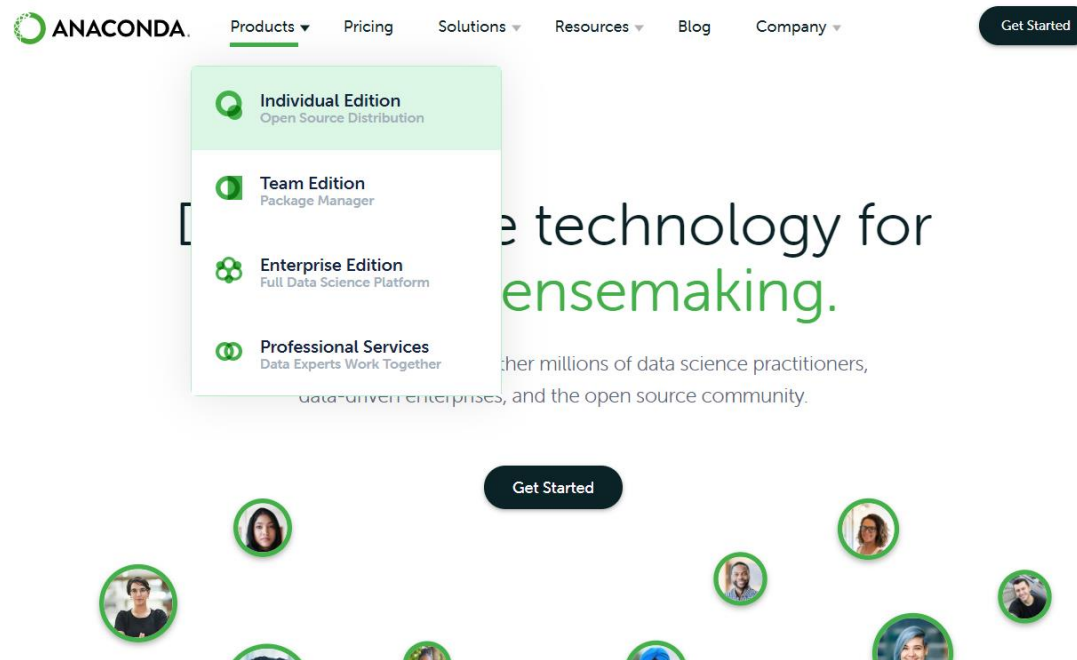
The name of columns is used to fill the **variable** attribute created during the **melt**.

Joining data frames



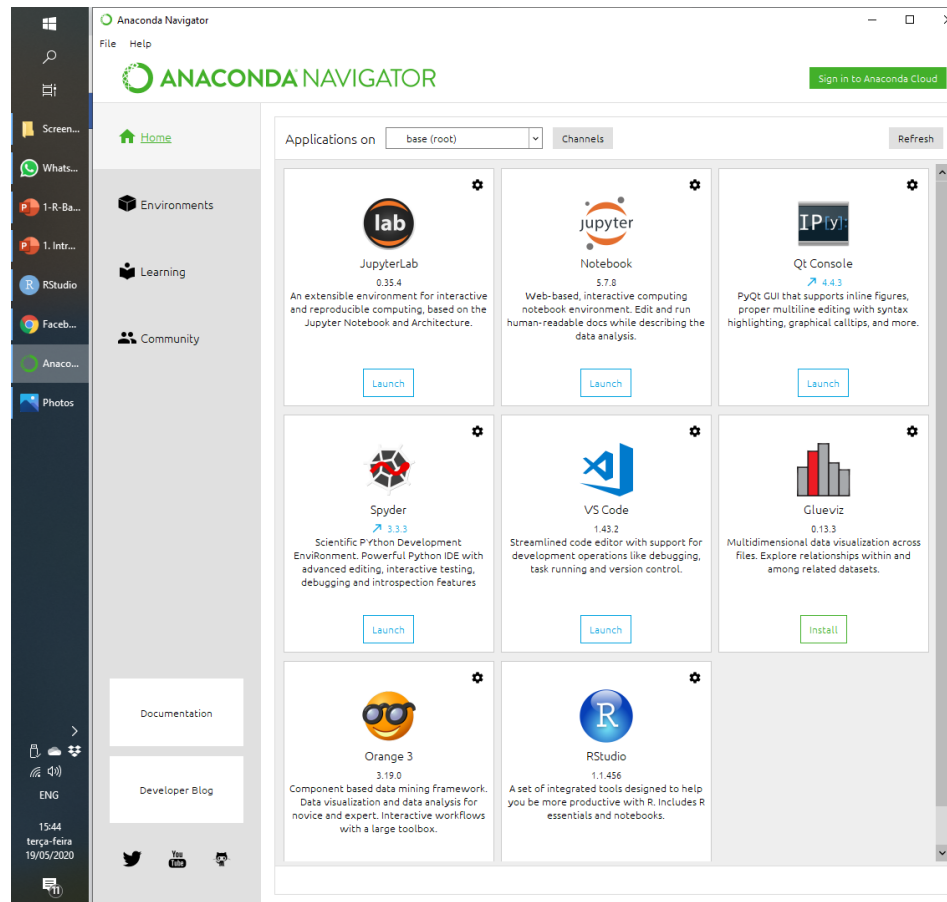
Jupyter notebooks

- Active document (document + code)
- Good tool for reproducibility
- Available at Anaconda

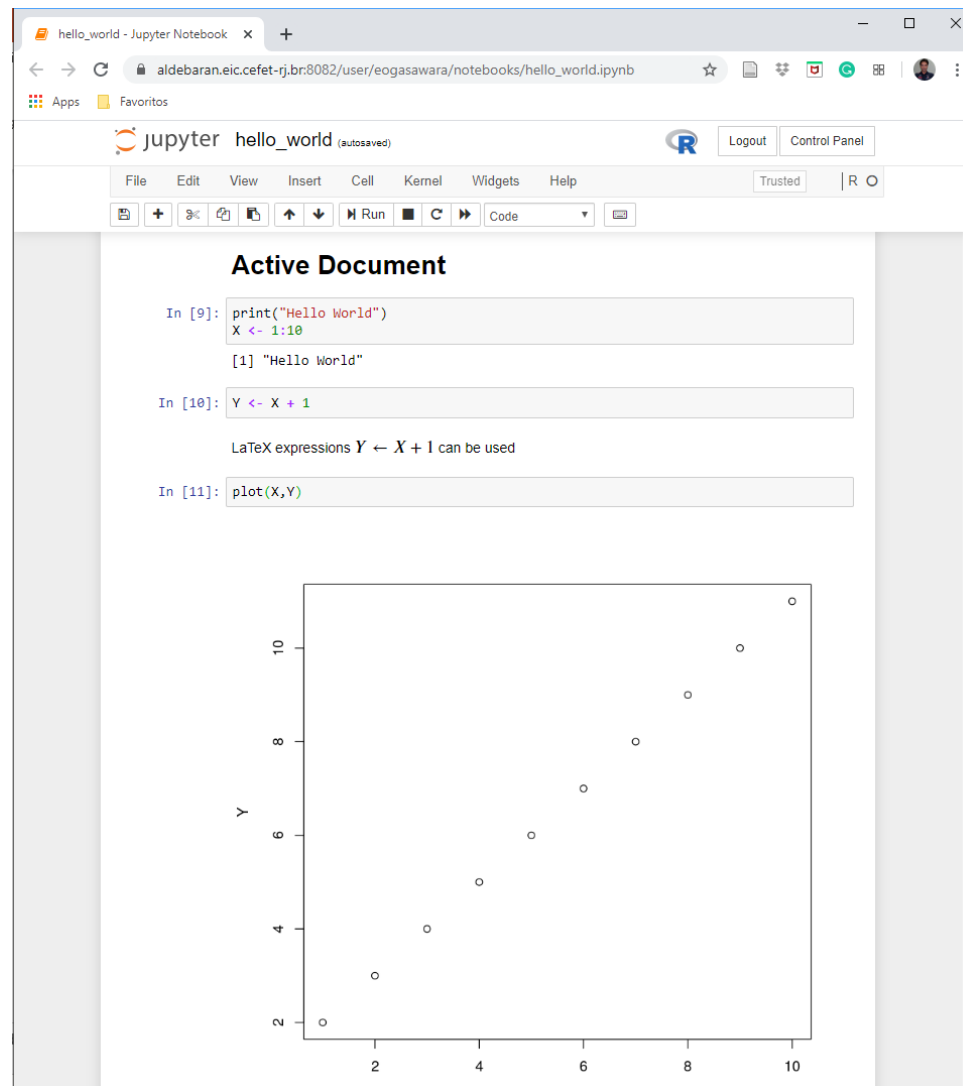


Anaconda

- Data Science Toolkit
 - Python + R + Jupyter + Orange + ...



Jupyter "Hello world"



The screenshot displays a Jupyter Notebook interface in a web browser. The browser's address bar shows the URL: `aldebaran.eic.cefet-rj.br:8082/user/eogasawara/notebooks/hello_world.ipynb`. The notebook's title bar reads "jupyter hello_world (autosaved)". The interface includes a menu bar with options: File, Edit, View, Insert, Cell, Kernel, Widgets, and Help. Below the menu is a toolbar with icons for file operations and a "Run" button. The main content area, titled "Active Document", contains three input cells:

- In [9]:** `print("Hello World")`
`X <- 1:10`
Output: `[1] "Hello World"`
- In [10]:** `Y <- X + 1`
Text below: LaTeX expressions $Y \leftarrow X + 1$ can be used
- In [11]:** `plot(X,Y)`

The output of the third cell is a scatter plot with the x-axis labeled 'X' and the y-axis labeled 'Y'. Both axes range from 0 to 10 with major ticks every 2 units. The plot shows ten data points, each represented by a small open circle, forming a straight line that passes through the origin and the point (10, 10). The points are located at (1, 2), (2, 3), (3, 4), (4, 5), (5, 6), (6, 7), (7, 8), (8, 9), (9, 10), and (10, 11).

Practicing

- Take some time to practice the examples
 - <https://nbviewer.jupyter.org/github/eogasawara/mylibrary/blob/master/Introduction.ipynb>
- Look at how to prepare nice graphics using ggplot2
 - <https://nbviewer.jupyter.org/github/eogasawara/mylibrary/blob/master/Graphic.ipynb>

Main References

