Introduction to R <u>Data S</u>cience Skills Day 2022

Anjali Silva, PhD

Summer Undergraduate Data Science Research Program University of Toronto 03 June 2022









Welcome!

- Instructor: Anjali Silva, PhD
 - Researcher and Lecturer, Department of Cell & Systems Biology, U of T
 - Data Analyst, University of Toronto Libraries
 - Pronouns: she/her
 - Name Phonetic: Un-j-li Sil-va

Course Description

- Introduction to R Data Science Skills Day
 - The vast amount of data produced by evolving information technology requires tools and skills. Among the many tools, R is a free, open-source language for data sciences. R is a programming language that can aid in the process of data analysis. This course is a beginner level, introductory course for R for data analysis. We will learn about R, RStudio (the environment use to work in R), including installation, and apply R for beginner-level data modeling and visualization. By the end of the course, you'll have a introduction to the flexibility of R, different functionalities, and understand how to apply it for basic data exploration.
 - Friday 10:00 am 4 pm EST; online synchronous.

Material

- Instructor Slides:
 - https://github.com/anjalisilva/DSI_IntroductionToR
 - SlideIntroR2022.pdf
- Instructor R Script:
 - https://github.com/anjalisilva/DSI_IntroductionToR
 - Script.R

Course Objectives

Learning Objectives:

- Install R and RStudio
- Navigate the RStudio environment
- Discover how to use RStudio to apply R to your analysis.
- Importing data from a spreadsheet
- View attributes of a dataset
- Understand differences in varying data types and structure
- Write and test functions
- Generate simple visualizations
- Be aware of sources for getting help in R
- Be aware of sources for expanding skills in R

Course Expectations

- Be respectful.
- One speaker at a time.
- Keep yourself on mute, unless you need to speak or ask a question.
- You may save your questions to 'Any questions?' section.
- If you have a question, use raise hand feature. First say your name, then ask the question.
- If you have a question, you may type it to chat as well.



Course Expectations



Figure: Zoom 'Reactions' that you may use.

Navigating an Online Code-Along Workshop

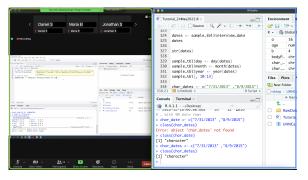


Figure: It is recommended that windows are resized so that both the user RStudio window and Instructor Zoom window (with RStudio) is visible at the same time. User may collapse panels of their RStudio not in current use.

Outline

Time	Торіс
10.00 -10.20 am	Introduction
10.20 - 10.50 am	Setup and RStudio
10.50 - 11.00 am	Short Break
11.00 - 12.15 pm	Basics + Vectors
12.15 - 1.00 pm	Lunch
1.00 - 2.00 pm	Matrices, Lists, Data Frames
2.00 - 2.10 pm	Short Break
2.10 - 3.10 pm	Data Import/Export; Functions
3.10-3.20 pm	Short Break
3.20 - 4.00 pm	Graphics; Next Steps

Any questions?

R and RStudio

Data Science Tools By Popularity

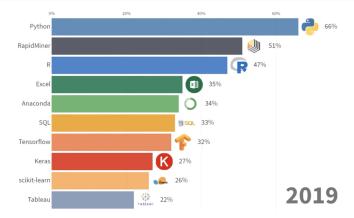


Figure: KDnuggets Survey of Machine Learning Software that asked respondents which data science tools they had used for projects within the past year. The x-axis shows the proportion of users who used a particular data science tool within the past year. Figure from https://www.kdnuggets.com/2020/06/data-science-tools-popularity-animated.html

Data Science Tools By Popularity

Top skills: Data Scientist

Sept 2018 to Sept 2019

Rank	Skill	Percent of jobs
1	python	79%
2	machine learning	72%
3	r	64%
4	sql	53%
5	hadoop	29%
6	spark	28%
7	java	25%
8	sas	24%
9	tableau	21%
10	deep learning	20%

Source: Indeed





Table titled "Top skills: Data Scientist." Indeed ranked the top skills in data scientist job postings from September 2018 to September 2019, comparing the percent of jobs for each skill. Results vary. Caption added oost-publication.



What is R?



- A language and environment for statistical computing and graphics.
- R was initially written by Ross Ihaka and Robert Gentleman.
- Since mid-1997, the R Core Team modify the R source.
- R runs on a wide variety of UNIX platforms, Windows and MacOS.



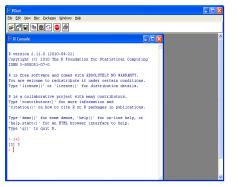
- R is a scripting language, thus an interpreter executes commands one line at a time.
- A Free software under the terms of the GNU General Public License.

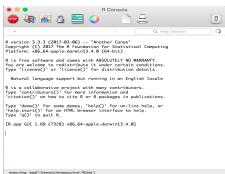
- R home page: https://www.R-project.org/
- How can R be obtained?
 - Via CRAN, the "Comprehensive R Archive Network".
 - https://cran.r-project.org/

Intro



- How can R be installed?
 - Unix
 - https://cran.r-project.org/doc/FAQ/R-FAQ.html# How-can-R-be-installed-_0028Unix_002dlike_0029
 - Windows
 - https://cran.r-project.org/bin/windows/base/
 - Mac
 - https://cran.r-project.org/bin/macosx/





- R can be used interactively or non-interactively.
 - Interactively, with or without an integrated development environment (IDE): RStudio.
 - Non-interactively via scripts.
- R is designed with interactive data exploration in mind.
- A version of R is released each year. Current release is 4.2.0.

• Ihaka,R. and Gentleman,R. (1996) R: a language for data analysis and graphics. *J. Comput. Graph. Statist.*, 5, 299–314.

RStudio



- RStudio is the Graphical User Interface for R.
- Not only for R, but can also use Python.
- An integrated development environment with:
 - A console
 - Syntax-highlighting editor for code execution
 - Tools for plotting, viewing history, debugging and workspace management

Why learn R?

- Not only is R free, but it is also open-source and cross-platform.
- R code is great for reproducibility.
- R is interdisciplinary and extensible.
- R works on data of all shapes and sizes.
- R produces high-quality graphics.
- R has a large and welcoming community.

Intro

- Online documentation for functions and variables in R exists.
- Obtained by typing help(FunctionName) or ?FunctionName at the R prompt, where FunctionName is name of function.
- E.g., if 'sum' is the function then:
 - > help(sum)
 - > ?sum

RStudio

 RStudio contains many features that make the development process easier and faster.

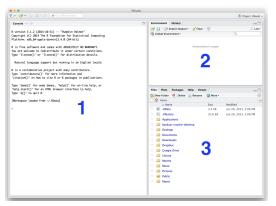


Figure: Anatomy of default RStudio. 1. This is the Console. 2. Environment and History. 3. Files, Plots, Packages, Help and Viewer. If a script is opened up, it will appear on top of Console.

Options To Work With R

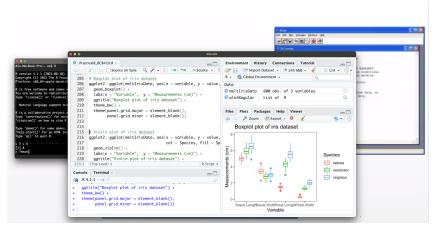


Figure: Options to work with R.

Options To Work With R



Any questions?

Practical - Setup

- This lesson assumes you have (current) versions of the following installed on your computer:
 - The R software itself, and
 - RStudio Desktop
- Any questions with R or RStudio setup?
- If you have downloaded R and RStudio:
 - Read about RStudio http://swcarpentry.github.io/ r-novice-inflammation/09-supp-intro-rstudio/index.html
 - Watch the video about data science tools. Pay attention to R https://www.youtube.com/watch?v=pKPaHH7hnv8&t=99s

Practical - Explore RStudio

RStudio

By now, you should have RStudio installed.



- There are two main ways of interacting with R:
 - Using the console
 - By using script files
- ullet Click on 'Tools' o 'Keyboard Shortcuts Help' for shortcuts.

Interacting with R

- Console:
 - Type commands directly into the console and press 'Enter' to execute.
- Script:
 - Put cursor at the end of the line to execute OR highlight the section.
 - Press 'Ctrl' + 'Enter' on Windows, Mac OR 'Cmd' + 'Return' on Mac.
- Clear console with 'Ctrl' + 'L'.
- If R is still waiting for you to enter more text, the console will show a + prompt.

R Project

- Good to keep data, analyses, and text in a single folder.
- RStudio interface for this is Projects.
 - ullet File o New project; choose New directory o New project
- Enter a name for this new folder ("directory") and choose a convenient location for it. This will be your working directory.
- On Desktop, save as 'DSI_IntroR'
- Click on 'Create' project.
- Create a new file where we will type our scripts.
 - Go to File → New File → R script. Click the save icon on your toolbar and save your script as "script.R".



Location

Current location:

```
getwd() # current location of the file, if saved
```

Set working directory:

```
By typing the path:
setwd("/Users/<Name>/Desktop")
```

Or (recommended method):

'Session' \rightarrow 'Set Working Directory' \rightarrow 'Choose Directory...'

Any questions?

Some Basics

Organize Working Directory

Structure of the working directory is very important.

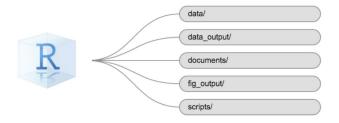


Figure: Examples of suggested directories within working directory or R Project. Figure from: https://datacarpentry.org/r-socialsci/00-intro/index.html

- Limit yourself to 80 characters per line.
- Use comments. Don't describe what the code does, but explain why you wrote it that way.
- Use only <- for assignment, not =.
- Never reassign reserved words.
- You may read more:
 - https://google.github.io/styleguide/Rguide.html
 - http://steipe.biochemistry.utoronto.ca/abc/index.php/ RPR-Coding_style

R Features

- In R, the indexing begins from 1.
- R is case sensitive ("X" is not the same as "x").
- R uses dynamic variable typing, so variables can be used over and over again.

Assignment and Commenting

- The ← symbol is the assignment operator.
- To assign a value to a variable or object called 'test1' test1 <- 123 test1
- Comment using # character
 test1 <- 123 # This is a comment
 test1 # This is called auto-printing

R Built-in Functions

- There are many built-in functions. You will learn these as you go.
- Functions combine a sequence of expressions that are executed to achieve a goal.
- The "argument" of the function is provided inside the brackets.
- The "return value" of the function is the value provided back.

• We will cover some basic functions:

x # auto-printing

print(x) # explicit printing using print() function
typeof(x) # "double" obtained using typeof() function
length(x) # 1 obtained using length() function

- Packages are collections of R functions, data, and compiled code.
- Libraries are directories in R where the packages are stored.
- Built-in functions are part of R standard or base packages and do not need to be downloaded.

```
library(help = "base")
library(help = "stats")
```

Function print() is part of base R package.

R Built-in Functions

- Standard or base R package contains basic functions for R to function as a language: arithmetic, input/output, basic graphics, etc.
- Some functions are not built-in. To get these, need to download pacakges.
- We will cover downloading of packages later.

R Version

To obtain session information: sessionInfo()

Version information:

R. Version()

Show objects in workspacels()

R Help Function

Getting help:

```
?"<-" # help on assignment operator
help("<-") # help on assignment operator
?typeof # help on typeof function
?class # help on class function
?print # help on print function</pre>
```

Any questions?

R Data Types

Numeric: floating types (double precision).

Basics

- Logicals: booleans = TRUE/FALSE or T/F.
- Character strings.
- Examples:

```
xValue <- 100
xValue
yVariable <- FALSE
yVariable
zVariable <- "hello"
zVariable
```



R Class

- Numbers in R are usually treated as numeric objects (i.e., double precision real numbers).
- To explicitly assign an integer, need to specify the L suffix.

```
x <- 1L
x
class(x) # "integer"</pre>
```

Complex class:

• Inf represents infinity:

NaN represents an undefined value/missing value:

Basics

```
NaN # not a number
0 / 0 # NaN
```

• c() function concatenating elements together:

- Character strings are collections of characters.
- Provided as values in single or double quotes.

```
xVariable <- 'hello'
class(xVariable) # "character"
zVariable <- "hello"
class(zVariable) # "character"
```

• "paste" converts inputs to strings, concatenate and return:

```
paste(xVariable)
```

Character Strings

• "cat" concatenates and prints the arguments to the screen:

```
cat("\n", xVariable, zVariable) # "\n" adds new line
```

"print" prints the argument: print(c(zVariable, xVariable)) Missing values are denoted by NA (Not Available) or NaN (Not a Number).

```
x <- c(1, 3, NA, 4, 5)
class(x) # "numeric"

y <- c(1, 3, NaN, 4, 5)
class(y) # "numeric"

# is.na() is used to test objects if they are NA
# is.nan() is used to test for NaN

is.na(x) # FALSE FALSE TRUE FALSE FALSE
is.nan(x) # FALSE FALSE FALSE FALSE</pre>
```

Question: What is the difference between NA and NaN in R?

Any questions?

Tips for Solving Issues

- Copy and paste the entire exact error message into Google.
 - Someone else may have gotten this same error and has asked a question.
- Copy and paste the entire error message into Google, followed by 'r'.
- Google the name of the function with term 'tutorial r' to see tutorials.
- If struggling with code for a plot, Google 'r plot plotname', then click on Images.
- If errors with reading files, ensure path is correct. Check using getwd().



- RStudio Community: https://community.rstudio.com/.
- Carpentries: https://datacarpentry.org/r-socialsci/.
- StackOverflow: https://stackoverflow.com/.
- R-help mailing list: https://stat.ethz.ch/mailman/listinfo/r-help.

Vectors

- Vector is a basic data structure in R.
- An R vector can only contain objects of the same class.
- There are multipel ways to create a vector:

```
y <- 1L:5L # vector using : operator
is.vector(y)</pre>
```

 $y \leftarrow c(1, 2, 3, 4, 5)$ # vector using c() function is.vector(y)

y <- seq(1, 5, by = 1) # vector using seq() function is.vector(y)

Vectors

• There are multipel ways to create a vector:

```
# vector using paste() function
y <- paste("A", 1:5, sep = "")
is.vector(y)

# vector using rep() function
y <- rep(letters[1:5], 3)
is.vector(y)</pre>
```

There are multipel ways to create a vector:

```
?vector # vector using vector() function
# vector(mode, length )
```

The atomic modes are "logical", "integer", "numeric" (synonym "double"), "complex", "character" and "raw":

```
# Initialize vector of certain length
x <- vector(mode = "numeric", length = 5)
x # 0 0 0 0 0
# Initialize vector of certain length
x <- vector(mode = "character", length = 10)
x # "" "" "" "" "" "" "" "" ""
```

 If mixing objects of two different classes in a vector, every element in the vector is forced to be same class.

```
y <- c(2.2, "a")
class(y) # "character"
y</pre>
```

 Sometimes objects can be coerced from one class to another using the as.* functions:

```
x # 1 2 3 4 5 6 7 8 9 10
class(x) # "integer"

z <- as.character(x)
z # "1" "2" "3" "4" "5" "6" "7" "8" "9" "10"
class(z) # "character"</pre>
```

x < -11.:101.

 Sometimes objects can be coerced from one class to another using the as.* functions:

```
w <- c("a", "b", "c")
w
class(w)
q <- as.numeric(w)
q # NA NA NA</pre>
```

Question:

Vector index in R starts from ____.

• To access the contents of the vector use []:

```
x <- 20:30 # vector
x # 20 21 22 23 24 25 26 27 28 29 30
length(x) # 11
x[1] # 20
x[15] # NA
x[c(1, 2, 4)] # 20 21 23</pre>
```

To remove elements:

$$x[c(-2, -4)]$$

Question:

Can you mix positive and negative integers when accessing elements of a vector?

 Can you mix positive and negative integers when accessing elements of a vector?

$$x[c(2, -4)] # ?$$

• There are several ways to modify vectors:

Any questions?

Matrices & Lists

- Matrices are constructed column-wise.
- Entries can be thought of starting in the "upper left" corner and running down the columns.
- Matrices must have every element be the same class (e.g., all integers).

```
?matrix
matrixOne <- matrix(data = 1:5, nrow = 2, ncol = 3)</pre>
matrixOne <- matrix(data = 1:6, nrow = 2, ncol = 3)</pre>
matrixOne
dim(matrixOne) # dimension 2 3
nrow(matrixOne) # 2
ncol(matrixOne) # 3
attributes(matrixOne)
```

Column-binding or row-binding can be done by cbind() and rbind() functions.

```
a <- 1:4
b <- 5:8
c <- cbind(a, b)
c
dim(c) # 4 2
d <- rbind(a, b)
d
dim(d) # 2 4</pre>
```

Lists

 A list is represented as a vector but can contain objects of different classes.

```
listOne <- list(16, "abc", TRUE, 5 + 4i)
listOne
length(listOne) # 4
typeof(listOne) # "list"
class(listOne) # "list"

# access the contents of the list
listOne[[1]] # 16
listOne[[2]] # "abc"
listOne[[4]] # "5+4i"</pre>
```

• Empty lists can be created using vector() function:

```
listTwo <- vector(mode = "list", length = 5)
listTwo
length(listTwo) # 5</pre>
```

Lists can have names:

```
listDestinations1 <- list(1, 2, 3)
listDestinations1
names(listDestinations1) # NULL
names(listDestinations1) <- c("Canada", "Alaska",
"England")
listDestinations1</pre>
```

Any questions?