# Introduction to R, RStudio, and R Markdown

Joseph Rudolf
API-R 2022

1

# Part I

2

**R**

Programming language for data analysis

**RStudio**

Interactive development environment (IDE)

**R Markdown**

Computational document format

3

# Getting Started with RStudio

4

## RStudio: On the Web and In Your Home

**RStudio Server**
Hosted on a server
(in the cloud)

**RStudio Desktop**
Installed locally on
your computer

**Note: Use Rstudio Server only for this course. Do not upload protected health information to the cloud!**

5

5

# Your Turn #1

Go to https://api-r.cloud in your browser and log in using the username and password provided in the course email.

Click "thumbs up" in zoom once you see the RStudio panes.

If you can't access the site click "thumbs down" and we will set you up in a backup configuration shortly.

6

# If You <u>Can't</u> Access api-r.cloud site

1.

# rstudio.cloud

3.
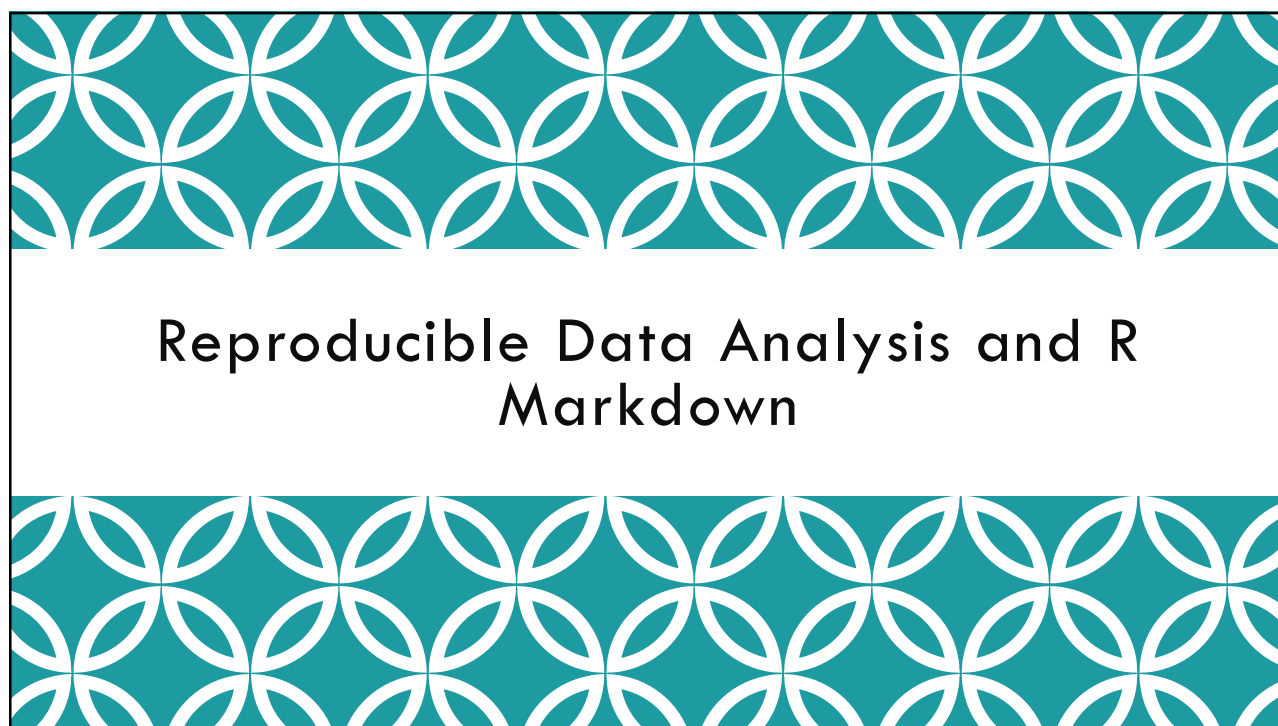


2.



4.



7



8

# Reproducible Data Analysis and R Markdown

---

## The Duke Cancer Scandal

❖ Chemo sensitivity from microarrays

❖ Errors first, then cover-up

❖ Clinical trials based on flawed models

❖ Papers retracted, lawsuits settled

Duke

MD Anderson

```
"1881_at"      "1882_g_at"
"31321_at"     "31322_at"
"31725_s_at"   "31726_at"
"32307_r_at"   "32308_r_at"
...
```

**Off-by-one indexing error**

11

"Common problems are simple…

Off-by-one **indexing error**

Sensitive / resistant **label reversal**

**Confounding** in experimental design

Inclusion of data from **non-reported sources**

**Wrong figure** shown

… and simple problems are common."

12

# Point-and-click is not reproducible



Computer code can precisely document each step of the analysis

13

# Why <u>YOU</u> should analyze your data reproducibly



"Can we redo the analysis with this month's data?"

"Why do the data in Table 1 not seem to agree with Figure 2?"

"Why did I decide to omit these six samples from my analysis?"

YOUR CLOSEST COLLABORATOR IS **YOU** FROM 6 MONTHS AGO

14

15

# Anatomy of an R Markdown Document

Header

Text
(with marks)

Code chunk

```
1  ---
2  title: 'My Markdown Document'
3  output: html_document
4  ---
5
6  # One Hashtag = Large Header
7
8  ## Two Hashtags = Smaller Header
9
10 Here is some text.
11
12 * It's easy to make a list
13 * Here is how you style text *cursive* or **bold**
14
15
16 ```{r}
17 x <- rnorm(100)
18 summary(x)
19 ```
20
```

16

17



18

# Your Turn #2

Open a sample R Markdown document (File -> New File -> R Markdown).

Review the format of the document:  header, text, code chunks

Execute the individual code chunks by selecting the Run Current Chunk arrow.

Knit the document to HTML (Preview or Knit Button -> Knit to HTML).  You may be prompted to save your R Markdown first.  In this case select a name for your document and click save.  Review the knitted document.

03:00

19

# Part II

20

21

# The Basics of Coding: Calculation

- R is a calculator!



press play button to execute code

answer returned here

22

# The Basics of Coding: Functions

- Code that extends our reach beyond the basic operators

```
1
2 ```{r}
3
4  abs(-77)
5
6  ```
```
[1] 77

function
(does stuff)

argument
(input)

**abs**(*-77* )

23

# Putting Functions to Work

- We can use functions to do more than simple math, we can make things!

- We can create a series of integers (a vector) using the seq() function

```
1
2 ```{r}
3
4  seq(from=5, to=150, by=10)
5
6  ```
```
[1]    5  15  25  35  45  55  65  75  85  95 105 115 125 135 145

24

# The Basics of Coding: Objects

• Objects are the container for your output

**object**
(stores output)

**function**
(does stuff)

**arguments**
(input)

sequence_of_10s <- **seq**(from=5, to=150, by=10 )

25

# Checking the contents of an object

• The environment tab shows us the objects we have created.

| Environment | History | Connections |
|---|---|---|

Import Dataset ▾     List ▾

Global Environment ▾

Values

sequence_of_10s     num [1:15] 5 15 25 35 45 55 65 75 85 95 ...

26

## Bending objects to your will

- Once we have created an object we can start to interact with it.

- This includes passing our objects to other functions... Whoa!

```r

min(sequence_of_10s)

```

```
[1] 5
```

```r

max(sequence_of_10s)

```

```
[1] 145
```

# Your Turn #3

I've written some code to create a sequence from 0 to 500 in increments of 25 called sequence_of_25s. Ultimately I want to calculate the median value of this sequence. Unfortunately I've made some mistakes in my code and I am hoping you can help me find them.

```r
sequence_of_25s -< seq(from=0 to=50, by=25)

```

```r
median(sequence of_25s]

```

`01:00`

# Importing Data

29

# The Data Analysis Pipeline



From *R for Data Science* (https://r4ds.had.co.nz/introduction.html)

30

30

31

# Tidyverse: R Packages for Data Science

- A consistent way to organize data

- Human readable, concise, consistent code

- Build pipelines from atomic data analysis steps



32

# Installing and loading R packages

tidyverse



```
function1()
function2()
function3()
function4()
```

**install.packages**("tidyverse")

Downloads files to computer
**1 x per computer**

**library**("tidyverse")

Loads package
**1 x per R Session**

33

# Dataframes: Beyond the Vector

- Dataframe is the term for a table

- Dataframes are composed:
  Columns (Variables)
  Rows (Observations)

- Dataframes are objects and can be acted on like other objects



csv    xls    db

34

# read_csv()

*data_frame* <- **read_csv**(*file_name*)

35

function
(does stuff)

*data_frame* <- **read_csv**(*file_name*)

36

37



38

# Your Turn #4

In the MISC pane, select the folder:
"exercises"

Select the R Markdown file:
"01 - Importing and Exploring Data.Rmd"

In the Editor pane, follow the instructions to complete the exercise.

`05:00`

41

# Recap



**Programming Language**    **IDE**    **Document Format**

**Packages** extend the functionality of R.  They need to be installed once per computer and loaded each session.

**Functions** do stuff.  They accept **Arguments** to define parameters.  We can store the output of functions in **Objects** using the assignment operator ( <- ).

**Importing Data** is the first step data analysis pipeline.  read_csv() is a function from the tidyverse that we can use for importing data.

42

# What else?

43



44

45

# Databases

| | | | | | |
|---|---|---|---|---|---|
| | Microsoft SQL Serve | | Amazon Redshift | | Other Databases |
| | MonetDB | | Apache Hive | | PostgreSQL |
| | MongoDB | | Apache Impala | | SQLite |
| | MySQL | | Athena | | Salesforce |
| | Netezza | | Cassandra | | Teradata |
| | Oracle | | Google BigQuery | | |

https://db.rstudio.com/databases/

46

47

# R Interface to Python



```{python}
import pandas
covid_testing.info()
```

48