

COURSE 07: Data Analysis with R Programming

Page No.

Date

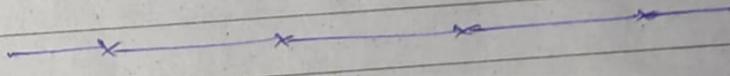
07 | 03 | 23

Journey So Far

Instructor : Carsee

Research Manager within People Operations

- > Use structured thinking to define a problem and ask the eight questions.
- > Work with spreadsheets, databases, and tools like SQL to organize and transform data
- > Clean your data to make sure it has integrity before you analyze it
- > Create impactful data visualization to illustrate key points.
- > Craft a compelling story to communicate insights to stakeholders.



ASK

Prepare

Process

Analyze

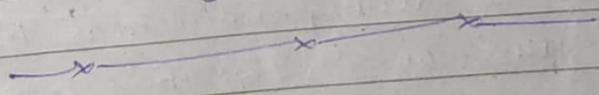
Share

Act

* Computer Programming:

Giving instructions to a computer to perform an action
or set of actions.

R is useful for organizing, cleaning & Analyzing



This course :-

- > intro to programming languages
- > explore main features and functions
- > basic programming concepts in R
- > How to work with data in R
- > Clean, transform, normalize, report data in R

Fun with R

Palmer Penguins Dataset →

Size measurements for three penguin species that live on the Palmer Archipelago in Antarctica



Mostly used for exploration & Teaching

Relax & watch the Power of R

- * Getting Started with R :

Don't be afraid of errors!

Week 1

2/3

08/03/23

Programming languages

The words & symbols we use to write instructions for computers to follow.

Coding is writing instructions to the computer in the syntax of a specific programming language

Ex. R, Python, JS, SAS, Scala, Julia

- * Advantages of using PL for DA :

① Clarify the steps of your Analysis

② Save time

③ Reproduce and Share work with colleagues, etc

Reproduce = Script

Page No.	
Date	

R: A PL frequently used for statistical Analysis, visualization and other data analysis.

R is based on S → 1970 John Chambers

for internal use at Bell Labs

1990s Ross Ihaka & Robert Gentleman developed R ~~\$81.9m~~

at University of Auckland, New Zealand

* Why use R → > Accessible

> Data-Centric

> Open Source

> Community

Check out: R for Data Science Learning Community

R-Studio Community

* When to use R → > Reproducing your Analysis

> Processing lots of data

> Creating data visualizations

↓ R is on Another level.

In the console:-

~~use~~ print("Hello, world!") → [1] Hello, world!

55 + 45 → [1] 100

69 == 96 → [1] FALSE

Lubridate is an R package that makes it easier to work with dates and times.

Intro to RStudio

IDE: A software app that brings together all the tools you may want to use in a single place.

R console, R Editor, tools for managing data & creating visuals.

R is the engine & RStudio is the Dashboard

* Packages are units of Reproducible R-code. Members of R-community create packages to keep track of R functions that they write and reuse.

* Lubridate is a part of tidyverse. The tidyverse is a collection of packages in R with a common design philosophy for data manipulation, exploration and visualization.

* For a lot of analysis tidyverse is essential

In the console type →

install.packages("tidyverse")

Some processes can take a while, the Red icon in the upper right of console indicates that.

Also

when the cursor appears in the console

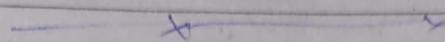
Takeaway \Rightarrow Lubricate is part of tidyverse; tidyverse \rightarrow Essential
install.packages("tidyverse")
library(tidyverse)
library(lubridate)

Page No. _____
Date _____

[R \rightarrow Case sensitive]

To load tidyverse package \Rightarrow library(tidyverse)
 \Rightarrow library(lubridate)

Done



R studio has 4 main windows called Panes \Rightarrow

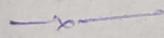
4th is hidden by default so you'll see 3 at first.
 \downarrow

It's the R-script : To open:

File > New File > R script

Drag border to adjust sizes

View > panes > \equiv for more customization



2 ways to code [Console] [Source]
[R-script] [Editor]

Console \Rightarrow install.packages("palmerpenguins")

\Rightarrow library("palmerpenguins")

\Rightarrow summary(penguins)

works together

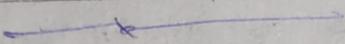
Source editor : View(penguins) \rightarrow shows tabular form

Environment Pane:

- > All the data you currently have uploaded
- > OR if you import using spreadsheet, it'll be here.
- > Can view each object by clicking.

All the previous commands are available in History tab.

Double clicking the line copies it to console

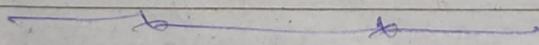


* Lower Right Pane: Files, Plots, Packages, Help, Viewer

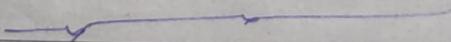
Files: Access to file directory and shows content of current working folder

Plots: Appear here.

Scatter plot
`ggplot(data = penguins, aes(x = flipper_length_mm,
y = body_mass_g) + geom_point(aes(color = species)))`



Help ✓✓



This week

- > Fundamentals of Programming using R in RStudio
- > Coding in RStudio
- > Syntax for performing calculations
- > Pipes make a sequence of code easier to work with and read.
- > R packages
- > Reusable functions & more
- > Tidyverse & packages (ggplot2 for viz.)

— — — — —

* Fundamentals > Basic concepts of R > Functions

> Comments

* Function (R):

A body of reusable code used to perform specific tasks in R.

> Variables

> Data Types

> vectors

> Pipes

* Argument (R)

Info that a function in R needs to run.

> print("Coding in R")

[1] "Coding in R"

> ?print() → Returns a page in help window

* Variable (R) ↗

A representation of a value in R that can be stored for use later during programming.

Variables also called as objects.

↓

Should start with a letter and contain numbers and underscores.

other assignment operator work too, but stick to just one type in your code

Page No.	
Date	

* # Comment in R

first-variable \leftarrow "This is my variable"

when we hit run, Environment Pane updates.

Second-var \leftarrow 12.5

[Run]

(logical, date, date-time, etc.)



* Vector (R) :

A group of data elements of same type stored in a sequence in R

> Can be made using combine function.

Just a c followed by no. ↓

c(a, b, c, d, ..., z)

> Vect_1 \leftarrow c(13, 48, 12.5, 17, 28.5)

> Vect_1

[1] 13.0 48.0 12.5 17.0 28.5



* Pipe (R) \Rightarrow

A tool in R for expressing a sequence of multiple operations, represented with "%>%"

④ Apply output of one function \Rightarrow to another function

This filters & sorts data ↴

Ex: Tooth growth % > % filter (dose = 0.5) % > % arrange (len)



* There are 2 types of vectors:

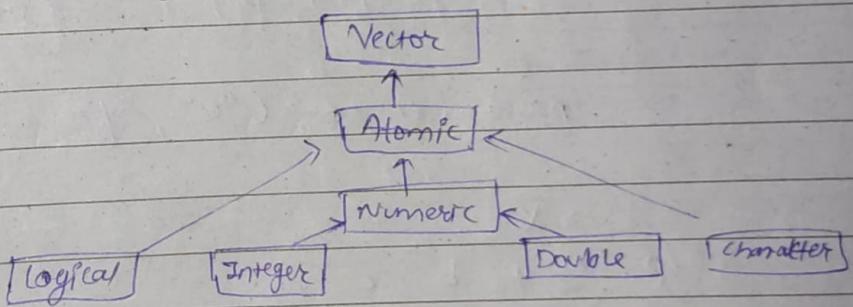
- o Atomic Vectors
- o Lists

Vector is a group of data elements of same type stored in a sequence.

There are 6 primary types of Atomic vectors

4 main Types	Logical	True/False	TRUE
	Integer	Positive and negative whole values	3
	Double	Decimal values	120.175
	Character	String / character values	"Coding"

Hierarchy ↴



Creating vector: Using c() function

Note: → When creating integer vector append L after each.

nos ← c(1L, 5L, 7L)

Also, c ← c(105.2, 7.9, 69.25)

c ← ("Sora", "Zoro", "Choro")

c ← (TRUE, FALSE, TRUE)

Length(vector) → determine length of vector.
↓
Var name

Page No.	
Date	

type of (c("a", "b"))
> [1] "character"

cool ← c(1, 3)
type of (cool)
> [1] "integer"

length (cool)
> [1] 2

Boolean Checking:

is.logical(), is.double(), is.integer(), is.character()
↓
is.vector()

Returns TRUE/FALSE

* Naming Vector ← names(x) ↳ vector name

x ← (1, 3, 5) { stored as double & not int}

names(x) ← ("a", "b", "c")

> x
0/p > a b c
> 1 2 3

Ex

remain ← c(11, 12, 11, 13)

suits ← c("spades", "hearts", "diamonds", "clubs")

names(remain) ← suits

> remain

#	Spades	hearts	diamonds	clubs
#	11	12	11	13

*Note → Atomic vector contains elements of same type

If other-way around use lists.

Page No.	
Date	

Also

`sremain ← c (spades = 11, hearts = 12,
diamonds = 11, clubs = 13)`

While simultaneously naming, it optional to keep
the names in quotes

`sremain ← c ("spades" = 11, "hearts" = 12,
"diamonds" = 11, "clubs" = 13)`

> str(sremain)

Named num [1:4] 11 12 11 13

- attr(f, "names") = chr[1:4] "spades" "hearts"
"diamonds" "clubs"

> my_apples ← 5

> my_oranges ← "six"

> is.vector(my_apples) # TRUE

> is.vector(my_oranges) # TRUE

> length(suits) # 4

single valued vectors

vectors allows only homogeneous types so it performs
any coercion if necessary,

converting different types to one ex. int to str, etc.

> class(suits)

character.

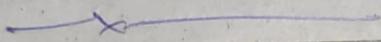
* Lists:

Elements can be of any type:
dates, dataframe, vectors, matrices, etc.

To create use list() similar to c()

Ex. `list("a", 1L, 1.5, TRUE)`

valid Ex. `list(list(list(1, 3, 5)))`



* Structure: str()

To determine size & type.

`str(list("a", 1L, 1.5, TRUE))`

- > List of 4
- > \$: chr "a"
- > \$ int 1
- > \$ num 1.5
- > \$ logi TRUE

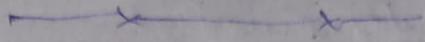
It describes the variable called.

`z <- list(list(1, 3, 5))`

`str(z)`

- > List of 2
- > \$ List of 1
- > .. \$ List of 3
- > num 1
- > num 3
- > num 5

Indentation suggest that it's a list of first of (PST).



* Naming lists:

List ('Chicago') = 1, 'New York' = 2, 'Los Angeles' = 3)

\$ Chicago

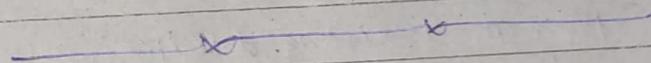
[1] 4

↳ 'New York'

[1] 2

\$ 'Los Angeles'

[1] 3



R - creates dates in standard: YYYY-MM-DD

Page No.

Date

14/03/23

Dates & times

- > install.packages ("tidyverse")
- > load(tidyverse)
- > load(lubridate)

Types → Date : "2016-08-16"
Date-time : "2011:19 UTC"
Date Time : "2018-09-17 18:15:58 UTC"

today() → gives current date

now() → gives current date time

* 3-way will be creating date-time formats:

> from a string

> from an individual date

> from an existing date/time objects

* Converting from string:

ymd ("2021-01-20") → [1] "2021-01-20"

Converts into
YYYY-MM-PP
ONLY

mdy ("January 20th, 2021") → [1] "2021-01-20"

dmy ("20-Jan-2021") → [1] "2021-01-20"

These three functions can also take quoted numbers and convert them into dates.

ymd (20210120)

mdy (01202021)

dmy (20012021)

→ [1] "2021-01-20"

* Date-time to Date \Rightarrow
as_date(now()) $> [1] \text{ 2021-01-20}$

2021		
Date		

* Date-time \Rightarrow

ymd() & its variations creates dates. To create date time add underscore (-) and add one or more of letters h,m,s

e.g. ymd_hms

ex. ymd_hms ("2021-01-20 20:11:59")

$> [1] \text{ "2021-01-20 20:11:59 UTC"}$

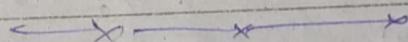
mdy_hm ("01/20/2021 08:01")

$> [1] \text{ "2021-01-20 08:01:00 UTC"}$

dm_y-h ("12/15/2022 20:01")

$> [1] \text{ "2022-12-15 20:00:00 UTC"}$

Every
combo



* Data Frames

data.frame(x = c(1, 2, 3), y = c(1.5, 5.5, 7.5))

$> \begin{array}{ccc} x & y \\ 1 & 1 & 1.5 \\ 2 & 2 & 5.5 \\ 3 & 3 & 7.5 \end{array}$

* Files

dir.create("dest_folder") # To create dir
file.create("new-file.txt")
file.create("new.docx") } self explanatory
file.create("new.csv")

If file.create returns TRUE else FALSE

file.copy("new.csv", "dest-fol")

unlink("som.txt") # To delete

*matrices :

- > objects in which elements are arranged in 2-D rectangular format
- > Only homogenous element type
- > Create using matrix()

Basic syntax →

`matrix(data, nrow, ncol, byrow, dimnames)`

\downarrow or \downarrow \downarrow
 $=$ $=$ Bool

Ex: `P <- matrix(c(3:14), nrow=4, byrow=TRUE)`

`print(P)`

	[,1]	[,2]	[,3]
#> [1,]	3	4	5
[2,]	6	7	8
[3,]	9	10	11
[4,]	12	13	14

access element at (1,3) & (4,2)

`print(P[1,3])` → 5

`print(P[4,2])` → 13

`Q <- matrix(c(14:28), nrow=4, byrow=TRUE)`

`print(P + Q)`

`print(P - Q)`

`print(P / Q)`

`print(P * Q)`

BODMAS()

Page No.

Date

14/03/23

* Logical Operators →

- 2 → AND
- | → OR
- ! → NOT

> if () { }
 > else () { }
 > else if () { }

Ex. datacamp.views ← c(100, 20, 5, 200, 60, 88, 190, 33, 290, 64)

youtube.views ← c(10, 20, 1000, 64, 100, 9, 19, 3, 90, 4)

⑪ datacamp.views > youtube.views

1. TRUE
2. FALSE
3. FALSE
- ⋮
10. TRUE

views ← matrix(c(datacamp.views, youtube.views),
 nrow=2, byrow=TRUE)

views

	100	20	5	200	60	88	190	33	290	64
10	20	1000	64	100	9	19	3	90	4	

views == 20

	F	T	F	F	F	F	F	F	F	R
F	F	T	F	F	F	F	F	F	R	

Packages → Units of reproducible R-code

Includes →

- > Re-usable R functions
- > Documentation about the functions
- > Sample datasets
- > Tests for checking your code

By default, R includes a set of packages called Base R.

Installed packages

↓
shows a list of installed package

priority → base → installed & loaded

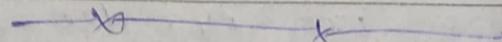
recommended → installed but not loaded

To load → library(name)

most commonly used source →

CRAN (Comprehensive R Archive Network)

An online archive with R packages, source code, manuals & R documentation.



Tidyverse (R):

A system of packages in R with a common design philosophy for data manipulation, exploration & visualization.

Using Tidyverse can help you work your way through pretty much the entire data analysis process.

Date			

Why was CRAN created?

- > Since packages are made by users, people need a reliable way to check and validate submitted code
- > CRAN makes sure any R content open to public meets required quality standards.

Github also is large source of R packages.

→ install package "tidyverse"

Options aren't always necessary.

→ library("tidyverse")

→ 8 other also imported & some conflicts

 Conflicts happen when packages have functions with same name as other functions

→ Function related to last package loaded is used

8 core tidyverse packages →

- > ggplot2
- > tibble
- > tidyverse
- > readr
- > purrr
- > dplyr
- > stringr
- > forecast

Tidyverse update(s): To keep track of updated

`tidyverse_update()` → keep track of updates

`update.packages()` → update all packages

`Install-package ("name")` → update single

- * Vignette: Documentation that acts as a guide to an R package.

> shares details about problems that the package is designed to solve.

② browseVignettes("name")

4 most essential packages for DA &

* {> ggplot2 → Data viz
> dplyr → Data manipulation
> tidyverse → Cleaning
> readr → Importing (read_csv, etc)

> tibble → Dataframe

> pure + functions & vectors

> strange \rightarrow function related to strings

> forecast → some common problem with factory

Factors (R)

Store categorical data in R where data values are limited and usually based on a finite group like country or year.

—→

Working with Pipes

* Pipe:

A tool in R for expressing a sequence of multiple operations, represented with "%>%"

↓
It takes the output of one statement and makes it the input of next statement.

↓
We can use pipe instead of typing out functions inside function. (nesting)

Example:

① Call up data (and then)

② Group the data (and then)

③ Summarize the grouped data using a mean function.

* To use already installed dataset we use data()

data("ToothGrowth")

View(ToothGrowth)

↓
Say we need to filter and sort this data for analysis

→ Nesting commands / sequence of data frames.

Filter

install.packages("dplyr") dplyr

library(dplyr)

filtered_tg ← filter(ToothGrowth, dose == 0.5)
View(filtered_tg)

Filter & Sort

View (%) To view

```
data("ToothGrowth")  
View(ToothGrowth)
```

Page No.
CTRL + SHIFT + M
Type

```
filtered_tg ← filter(ToothGrowth, dose == 0.5)  
View(filtered_tg)
```

```
arrange(filtered_tg, len)
```

→ Output in console

Ascending sorted by len

So the tooth growth data is filtered ($dose == 0.5$)
& sorted by len.

↓

Another way ↗

```
arrange(filter(ToothGrowth, dose == 0.5), len)
```

Pipe

```
filtered_toothgrowth ← ToothGrowth %>%
```

```
filter(dose == 0.5) %>%
```

```
group_by(supp) %>%
```

```
summarize(mean(len) = mean(len, na.rm = T), .group = "drop")
```

- If pipe operator is detected, then automatically
indentation is corrected.
- Last line doesn't have a pipe

- Pipe op is added at the end
- Check code after adding pipe
- Reused to check what went wrong.

Connor : Marketing Analytics Manager!

- > Overall readability is an important aspect.
- > Using Comments!
- > Documentation will explain in depth exactly what your code is doing, why it was built, what is the purpose for it and any limitations.
- > Building it for scalability as well as making it dynamic.
Q: Will the code be used in future?
if yes, make it in such a way that it will be able to scale.

Dynamic: Not hard coding YUP!