

Introduction to Data Science (11372 & G 11516)  
Semester 1 2021


UNIVERSITY OF  
CANBERRA

# INTRODUCTION TO DATA SCIENCE

## Lecture 5

Dr. Ibrahim Radwan

DISTINCTIVE BY DESIGN

1

---

---

---

---


---

---

---

---

# OUTLINE


UNIVERSITY OF  
CANBERRA

- Data Science, a Practical View
- Data Sources and Data Formats
- Tibbles vs. Data Frames
- Data Import with 'readr' Package
- Variable Specifications and Parsing

©Dr. Ibrahim Radwan – University of Canberra

2

---

---

---

---


---

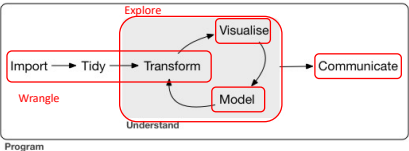
---

---

---

# DATA SCIENCE; A PRACTICAL VIEW


UNIVERSITY OF  
CANBERRA



### Program Steps

- 1- Reading Data
- 2- Data Wrangling
- 3- Data Exploratory
- 4- Modelling
- 5- Result Communication

©Dr. Ibrahim Radwan – University of Canberra

3

---

---

---

---

---

---

---

---

## DATA SOURCES



Data scientist deals with a diverse set of data sources (e.g. sensors, Internet, etc.)



Internet

Sensors

Manually created or collected

©Dr. Ibrahim Radwan – University of Canberra

<https://www.kdnuggets.com/2017/12/big-data-free-sources.html>

4

## DATA SOURCES (2)



- DataHub
  - <https://datahub.io/>
- OpenML
  - <https://www.openml.org>
- Australia Open Data
  - <https://data.gov.au/>
- GitHub repositories, manually collected work
  - <https://github.com/awesomedata/awesome-public-datasets>

©Dr. Ibrahim Radwan – University of Canberra

5

## DATA FORMATS



Data scientist also deals with different types of content such as:



Texts



Images



Sounds

This means that the data scientist spends lot of time on pre-processing these data into some data structure that is manageable by the data analysis techniques.

©Dr. Ibrahim Radwan – University of Canberra

6

## DATA FORMATS (2)



- ▶ The most important and the commonly used data structure is the **2-D data table**, which is known as data frame in R
- ▶ Characteristics of the data tables:
  - Each row represents an entity (e.g. product, person, etc.) and the columns represent the properties (e.g. name, age, temperature, etc.) we have measured for the entities.
  - The entities are frequently referred to as objects, tuples, records, examples or feature-vectors.
  - The properties are often called features, attributes, variables, dimensions or fields.

©Dr. Ibrahim Radwan – University of Canberra

7

---

---

---

---

---

---

---

---

## DATA FORMATS (3)



- ▶ Characteristics of the data tables (continued):
  - Rows of a dataset can be either:
    1. independent from each others; or
    2. There is some dependency among them.
      - \* For examples, dependencies in the form of time order (e.g. measurements of a set of variables in successive time steps).
  - Columns of a dataset can be either:
    1. Quantitative, or
    2. Categorical.
  - Columns also can be independent or correlated

©Dr. Ibrahim Radwan – University of Canberra

8

---

---

---

---

---

---

---

---

## DATA FORMATS (4)



- ▶ Characteristics of the data tables (continued):
  - A finer categorization for the types of the columns can be:
    1. **Interval**: quantitative variables like for instance dates.
    2. **Ratio**: quantitative variables like for instance height of a person or price of a product.
    3. **Nominal**: these are categorical variables whose values are some sort of labels without any ordering among them (e.g. colors).
    4. **Ordinal**: again categorical variables but this time with some implicit ordering among their finite set of values (e.g. small, medium and large).

©Dr. Ibrahim Radwan – University of Canberra

9

---

---

---

---

---

---

---

---

## DATA FORMATS (5)



- Most of the data scientists are dealing with the following data sources:
  - Text files, separated by some delimiters (e.g. spaces, commas (i.e. csv), etc.)
  - Spreadsheets (eg, .xls, .xlsx, etc.)
  - Databases (MySQL, SQL, Oracle, etc.)
  - Other software-specific formats
- In IDS unit, we will mostly deal with the plain text, **rectangular** data files such as .txt, .csv, etc.
- To handle these data, we need to read them in **data frames**

©Dr. Ibrahim Radwan – University of Canberra

10

---

---

---

---

---

---

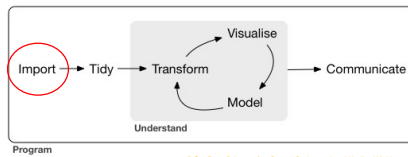
---

---

## READING DATA



### Program Steps



- 1- Reading Data
- 2- Data Wrangling
- 3- Data Exploratory
- 4- Modelling
- 5- Result Communication

©Dr. Ibrahim Radwan – University of Canberra

11

---

---

---

---

---

---

---

---

## DATA FRAMES, A RECAP



ID, Name, Age  
23424, Ana, 45  
11234, Charles, 23  
77654, Susanne, 76

data.csv

How to read this using  
the base functions in R?

Data Frame				
	X	Y	Z	
1	...	...	...	Observations How many?
2	...	...	...	
3	...	...	...	
4	...	...	...	
5	...	...	...	
6	...	...	...	
	Variables Types Names			

©Dr. Ibrahim Radwan – University of Canberra

12

---

---

---

---

---

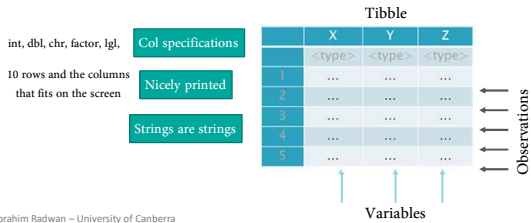
---

---

---

## TIBBLES – A CUSTOMIZED DATA FRAMES UNIVERSITY OF CANBERRA

- Tibbles are enhanced data frame objects where some of the features have been added and some others were dropped to make the life a bit easier when dealing with big datasets.



13

## TIBBLES UNIVERSITY OF CANBERRA

- Tibbles are part of the **'tidyverse'** library that has been developed by *Hadley Wickham*.
- The **'tidyverse'** library provides an effective way to read, visualize and to clean the data. It has packages such as:
  - 'readr' to import data
  - 'tidyr' and 'dplyr' to wrangle the data
  - 'ggplot2' to visualize the data

©Dr. Ibrahim Radwan – University of Canberra

14

## TIBBLES VS. DATA FRAMES UNIVERSITY OF CANBERRA

```

1 # To use tibble, we first need to load it
2 # It is part of the tidyverse package
3 # If it's not installed, install it
4 install.packages("tidyverse")
5
6 library("tidyverse")
7 # creating a tibble from existing data frame
8 my_data <- iris
9 my_data # how is it printed?
10
11 # convert it to tibble
12 my_data2 <- as_tibble(my_data)
13 my_data2 # and now
14
15 #####
16 # another way to create a tibble
17
18 # recycling feature is allowed when it is just one element in a column
19
20 # also with a transposed tibble, 'tribble'
21
22
23
24
25
26
27
28
29
30
31
32

```

©Dr. Ibrahim Radwan – University of Canberra

15

## TIBBLES VS. DATA FRAMES (2)



```

34 # tibbles vs. dataframes
35 # i- printing (first 10 rows, all cols that fits with the screen, type of column)
36 tbl_3 <- tibble(
37   a = lubridate::now() + runif(1e3) * 86400,
38   b = lubridate::today() + runif(1e3) * 30,
39   c = 1:1e3,
40   d = runif(1e3),
41   e = sample(letters, 1e3, replace = TRUE)
42 )
43
44 df_3 <- data.frame(
45   a = lubridate::now() + runif(1e3) * 86400,
46   b = lubridate::today() + runif(1e3) * 30,
47   c = 1:1e3,
48   d = runif(1e3),
49   e = sample(letters, 1e3, replace = TRUE)
50 )
51
52 # If you want to print more rows, n = 7
53 print(tbl_3, n=20, width=inf)
54 # also you can use the view
55 view(tbl_3)
56
57 # 2- subsetting
58 df4 <- data.frame(x = 1:3, y = 3:1)
59 class(df4)
60
61 tbl4 <- tibble(x = 1:3, y = 3:1)
62 class(tbl4)
63
64 class(tbl4)
65 # Tibbles are also stricter with $. Tibbles never do partial matching.
66 df5 <- data.frame(abc = 1)
67 class(df5)
68
69 tbl5 <- tibble(abc = 1)

```

©Dr. Ibrahim Radwan – University of Canberra

16

## DATA IMPORT



- ▶ Reading flat files such as CSV or text files
- ▶ Suppose we have the following CSV file:
  - ▶ We will use the **readr** library to read the csv files
  - ▶ Then, we can use the `read_csv("filename.csv")` function to read the contents.
  - ▶ The returned object is a tibble.
- ▶ It reads the data with 10x faster than the base R functions.

ID	Name	Age
23424	Ana	45
11234	Charles	23
77654	Susanne	76

```

> read_csv("data_sample.csv")
#> data frame with 3 columns and 3 rows
#> # A tibble: 3 x 3
#>   ID Name   Age
#>   <dbl> <chr> <dbl>
#> 1 23424 Ana    45
#> 2 11234 Charles 23
#> 3 77654 Susanne 76

```

©Dr. Ibrahim Radwan – University of Canberra

17

## DATA IMPORT (2)



```

read_(file, col_names = TRUE, col_types = NULL, locale = default_locale(), na = c("", "NA"),
      quoted_na = TRUE, comment = "", trim_ws = TRUE, skip = 0, n_max = Inf, guess_max = min(1000,
      n_max), progress = interactive())

```

Comma Delimited Files  
read\_csv("file.csv")

Tab Delimited Files  
read\_tsv("file.tsv")  
Also read\_table()

Semi-colon Delimited Files  
read\_csv2("file2.csv")

To save data into csv or txt file

Files with Any Delimiter  
read\_delim("file.txt", delim = "|")

```

write_csv(x, path, na = "NA", append = FALSE,
          col_names = lappend)
write_delim(x, path, delim = "|", na = "NA",
            append = FALSE, col_names = lappend)

```

©Dr. Ibrahim Radwan – University of Canberra

18

## DATA IMPORT (3)



```

1 # Import data into R
2
3 # It is part of the "tidyverse" package
4 if (!("tidyverse" %in% rownames(installed.packages()))){
5   install.packages("tidyverse")
6 }
7
8 library("tidyverse")
9
10 heights <- read_csv("data/heights.csv")
11
12 # supply inline csv file
13 read_csv("a,b,c
14 1,2,3
15 4,5,6")
16
17 # skip first two lines, i.e. meta data the beginning of a file
18 read_csv("The first line of metadata
19 The second line of metadata
20 x,y,z
21 1,2,3", skip = 2)
22 # or skip when comment
23 read_csv("# A comment I want to skip
24 x,y,z
25 1,2,3", comment = "#")
26
27 # by default the first row is for column, we can skip this rule
28 read_csv("1,2,3\n4,5,6", col_names = FALSE) # \n is convenient shortcut for new line
29 # pass column names
30 read_csv("1,2,3\n4,5,6", col_names = c("x", "y", "z"))
31
32 # deal with na values
33 read_csv("a,b,c\n1,2,,", na = "",)

```

©Dr. Ibrahim Radwan – University of Canberra

19

## DATA IMPORT – COL SPECIFICATIONS



- The functions of 'readr' package **guess** the types of each column and convert types when appropriate (but will NOT convert strings to factors automatically), a message shows the type of each column in the result.

```

## Parsed with column specification:
##   col1
##   age = col_integer(), age is an integer
##   sex = col_character(), sex is a character
##   earn = col_double(), earn is a double (numeric)
##

```

The solution is to use  
problems() function

For each col\_\* function, there is a parse\_\* function

col\_guess function uses the first 1000 rows from each column to guess the column type.

What about if:

- The first 1000 rows are just special case and the rest are different?
- The first 1000 rows are 'na' ?

©Dr. Ibrahim Radwan – University of Canberra

20

## DATA IMPORT – COL SPECIFICATIONS (2)



```

1 # col specifications
2 # logicals
3 lg1 <- parse_logical(c("TRUE", "FALSE", "NA"))
4 str(lg1)
5 # integers
6 intr <- parse_integer(c("1", "2", "3"))
7 str(intr)
8 # date
9 dt <- parse_date(c("2010-01-01", "1979-10-14"))
10 str(dt)
11 # which strings should be specified as missing
12 intr_missing <- parse_integer(c("1", "231", ".", "456"), na = ".")
13 str(intr_missing)
14 # sometimes the parsing fails, NAs will be used when failure is occurring
15 x <- parse_integer(c("123", "345", "abc", "123.45", "221"))
16 x
17 # here you can use problems
18 problems(x)
19
20 # parse numbers with . or , or $ or %
21 db1 <- parse_double("1.24")
22 str(db1)
23 num <- parse_number("$123,456,789")
24 num
25 num <- parse_number("$123,456,789")
26 num
27 num <- parse_number("$123.456.789", locale=locale(grouping_mark = '.'))
28 num
29 num <- parse_number("123'456'789", locale = locale(grouping_mark = "'"))
30 num

```

©Dr. Ibrahim Radwan – University of Canberra

21

## DATA IMPORT – COL SPECIFICATIONS (3) UNIVERSITY OF CANBERRA

- Using the `problems()` function with the `read_*` functions from the `'readr'` package to check the problems in parsing and guessing the column types.

```
problems(tibble_object)
```

©Dr. Ibrahim Radwan – University of Canberra

22

## DATA IMPORT – COL SPECIFICATIONS (4) UNIVERSITY OF CANBERRA

```
1 # base function to read data.frame
2 df <- read.csv(readr_example("challenge.csv"))
3
4 df
5 view(df)
6 df[1000,]
7 str(df)
8 # problems(df) # no problems as everything are chars
9 #####
10 # readr functions to read csv
11 tbl
12 view(tbl)
13 tbl[1000,]
14 str(tbl)
15 #####
16 # use problems to check the problems in the tbl parsing
17 problems(tbl)
18
19 # how to fix the problems
20 # first let us fix the problem by hard-specifying the column types
21 tbl2 <- read_csv(
22   readr_example("challenge.csv"),
23   col_types = cols(
24     x = col_double(),
25     y = col_character()
26   )
27 )
28
29 # there are dates stored with this column
30 # change the character to date
31 tbl3 <- read_csv(
32   readr_example("challenge.csv"),
33   col_types = cols(
34     x = col_double(),
35     y = col_date()
36   )
37 )
38
```

©Dr. Ibrahim Radwan – University of Canberra

23

## KEY TAKEAWAYS



- As a data scientist, you will deal with different data formats and data sources
- Reading the rectangular data into tibbles are much efficient than using the `data.frame`
- Tibble is a customized and enhanced version of the basic data frame, where the features that resist the test of the time have been kept and the frustrating ones have been dropped
- Using the functions inside `readr` package to read and to write flat data are much faster and efficient than using the basic functions in R

©Dr. Ibrahim Radwan – University of Canberra

24



## RECOMMENDED READING



- You are recommended to read chapters 10 & 11 from the "*R for Data Science*" book:
  - <https://r4ds.had.co.nz/tibbles.html>
  - <https://r4ds.had.co.nz/data-import.html>

©Dr. Ibrahim Radwan – University of Canberra

---

---

---

---

---

---

---