Introduction to Data Science (11372 & G 11516)

Semester 1 2021



INTRODUCTION TO DATA SCIENCE

Lecture 5

Dr. Ibrahim Radwan

DISTINCTIVE BY DESIGN

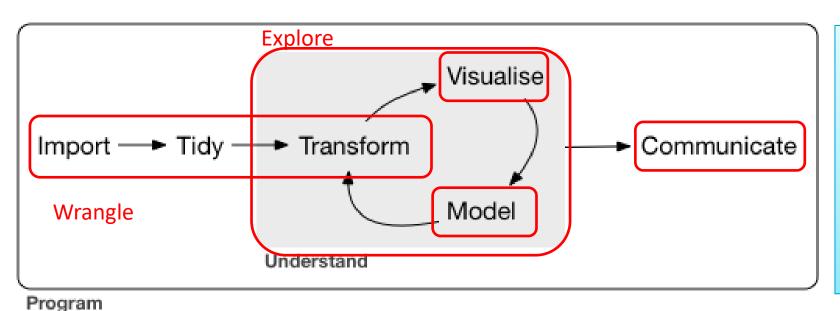
OUTLINE



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

DATA SCIENCE; A PRACTICAL VIEW





Program Steps

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

R for Data Science, by Garrett Grolemund and Hadley Wickham

DATA SOURCES



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



Internet

Sensors

Manually created or collected

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

DATA FORMATS



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

amorphice or vinegar and oil fontend of sale alty toppings such as bacon bits, chaeser, as at olives.

and beverages such as tomato juice, vegetabget sims. Choose fresh fruits and fruit juices, of food restaurants. Order a plain sandwich mustard, catsup, cheese or special sauces, benato and onion instead. Order french friest.





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.

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.

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

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).

DATA FORMATS (5)

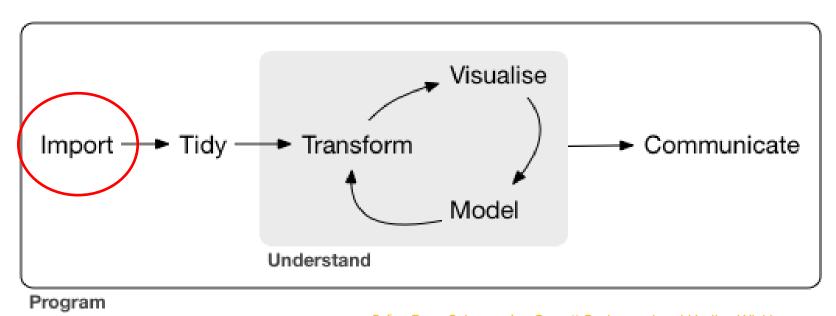


- Most of the data scientists are dealing with the following data sources:
 - 1. Text files, separated by some delimiters (e.g. spaces, commas (i.e. csv), etc.)
 - 2. Spreadsheets (eg, .xls, .xlsx, etc.)
 - 3. Databases (MySQL, SQL, Oracle, etc.)
 - 4. Other software-specific formats
- In IDS unit, we 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 <u>data frames</u>

READING DATA



Program Steps



R for Data Science, by Garrett Grolemund and Hadley Wickham

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

DATA FRAMES, A RECAP



ID, Name, Age

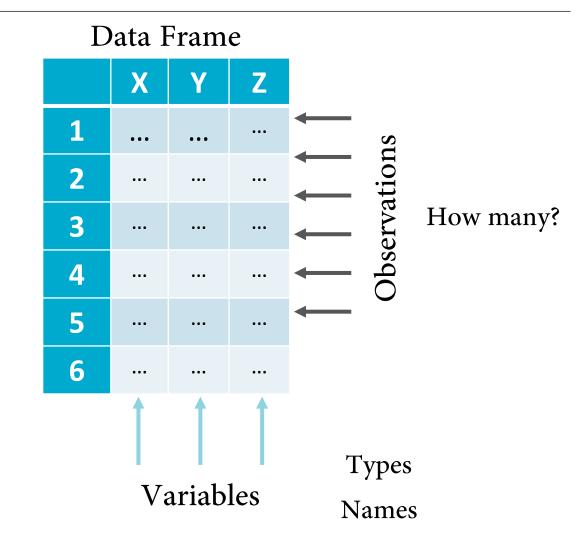
23424, Ana, 45

11234, Charles, 23

77654, Susanne, 76

data.csv

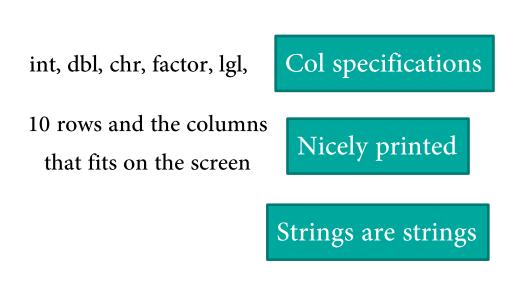
How to read this using the base functions in R?

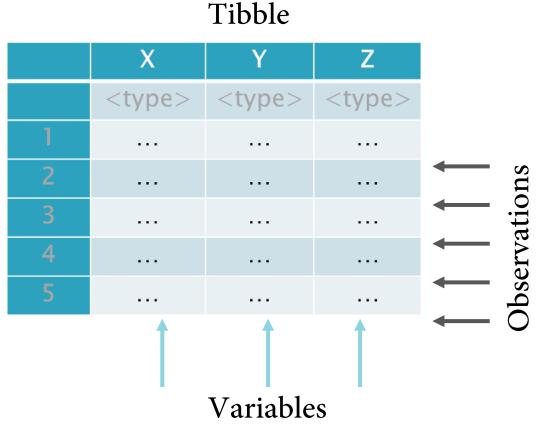


TIBBLES - A CUSTOMIZED DATA FRAMES A CANBE



 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.





TIBBLES



- 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

TIBBLES VS. DATA FRAMES



```
1 # To use tibble, we first need to load it
   # It is part of the `tidyverse` package
 3 * if (!("tidyverse" %in% rownames(installed.packages()))){
     install.packages("tidyverse")
 5
  library("tidyverse")
   # creating a tibble from existing data frame
  my_data <- iris
   my_data # how is it printed?
11
12
13
  # convert it to tibble
  my_data2 <- as_tibble(my_data)</pre>
15
   my_data2 # and now
16
17
   # another way to create a tibble
19 tbl_1 <- tibble(</pre>
20
  x=1:5.
21
  y=1, # recycling feature is allowed when it is just one element in a column
22
     z=x\wedge 2+y
23
24
  tbl 1
25
  # also with a transposed tibble, `tribble`
27 tbl_2 <- tribble(</pre>
28
    ~x, ~y, ~z,
    "a", 2, 3.6, "b", 1, 8.5
29
30
31
32 tb1_2
```

TIBBLES VS. DATA FRAMES (2)



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

© Dr. Ibrahim Radwan – University of Canberra

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.

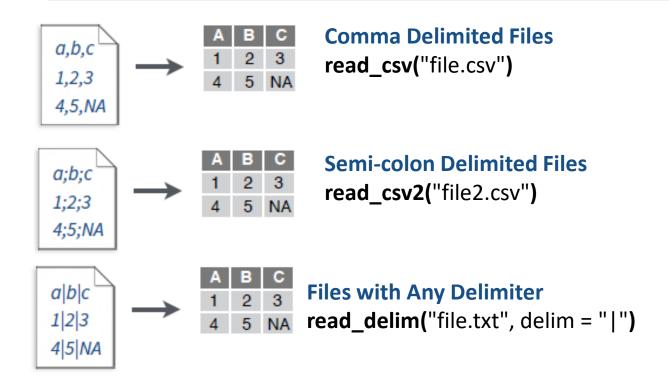
```
ID, Name, Age
23424, Ana, 45
11234, Charles, 23
77654, Susanne, 76
```

```
> read_csv("data_sample.csv")
Parsed with column specification:
cols(
  ID = col_double(),
  Name = col_character(),
  Age = col_double()
 A tibble: 3 \times 3
     ID Name
                   Age
  <db1> <chr>
  23424 Ana
                    45
  11234 Charles
                    23
 77654 Susanne
                    76
>
```

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())





To save data into csv or txt file

Comma delimited file write_csv(x, path, na = "NA", append = FALSE, col_names = !append) File with arbitrary delimiter write_delim(x, path, delim = " ", na = "NA", append = FALSE, col_names = !append)

DATA IMPORT (3)



```
# Import data into R
   # It is part of the `tidyverse` package
4 - if (!("tidyverse" %in% rownames(installed.packages()))){
     install.packages("tidyverse")
 6
   library("tidyverse")
 9
  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
1,2,3", skip = 2)
22 # or skip when comment
23 read_csv("# A comment I want to skip
24
    x,y,z
     1,2,3", comment = "#")
25
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
  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

DATA IMPORT – COL SPECIFICATIONS



• The functions of 'readr' package <u>guess</u> 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:
## cols(
## age = col_integer(),
## sex = col_character(),
## earn = col_double()
## )

sex is a
character
```

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'?

DATA IMPORT – COL SPECIFICATIONS (2) TO CANBERRA

```
1 # col specifications
 2 # logicals
  lgl <- parse_logical(c("TRUE", "FALSE", "NA"))</pre>
4 str(lql)
 5 # integers
6 intgr <- parse_integer(c("1", "2", "3"))</pre>
7 str(intgr)
8 # date
   dt <- parse_date(c("2010-01-01", "1979-10-14"))
10 str(dt)
# which strings should be specified as missing
intgr_missing <- parse_integer(c("1", "231", ".", "456"), na = ".")</pre>
13 str(intgr_missing)
14 # sometimes the parsing fails. NAs will be used when failure is occuring
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
    dbl <- parse_double("1.24")
    str(dbl)
   num <- parse_number("$123,456,789")
24 num
25  num <- parse_number("$123,456,789")</pre>
26 num
    num <- parse_number("$123.456.789", locale=locale(grouping_mark = '.'))</pre>
27
28 num
   num <- parse_number("123'456'789", locale = locale(grouping_mark = "'"))</pre>
30
   num
```

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)

DATA IMPORT – COL SPECIFICATIONS (4) A CANBERRA



```
1 # base function to read data.frame
 2 df <- read.csv(readr_example("challenge.csv"))</pre>
 3 df
 4 view(df)
 5 df[1001,]
 6 str(df)
 7 problems(df) # no problems as everything are chars
 9 # readr functions to read csv
10 tbl <- read_csv(readr_example("challenge.csv"))</pre>
11 tbl
12 view(tbl)
13 tbl[1001,]
14 str(tbl)
16 # use problems to check the problems in the tbl parsing
   problems(tbl)
17
18
19 # how to fix the problems
20 # first let us fix the problem by hard-specifying the column types
21 tbl2 <- read_csv(</pre>
22
     readr_example("challenge.csv"),
23
     col_types = cols(
24
       x = col_double(),
25
       y = col_character()
26
27
28 tbl2
29 tail(tbl2) # there are dates stored with this column
30 # change the character to date
31 tbl3 <- read_csv(</pre>
     readr_example("challenge.csv"),
32
     col_types = cols(
33
       x = col_double(),
34
       y = col_date()
35
36
37
```

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

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