

DSA2101

Essential Data Analytics Tools: Data Visualization

Yuting Huang

Weeks 3 – 4: Importing Data to R

Contents

In the next two weeks, we will learn how to import data into R.

1. CSV files
2. Excel Files
3. R data files
4. JSON Files
5. Data from the Web

Recap: Navigation

An important pre-requisite to loading data into R is that we are able to point to the location at which the data files are stored.

1. Where am I?
2. Where are my data?

Working directory

The first question addresses the notion of our current **working directory**.

- ▶ Typically, it is the location of our current R script.
- ▶ We can use the function `getwd()` to obtain the current working directory.

```
getwd()
```

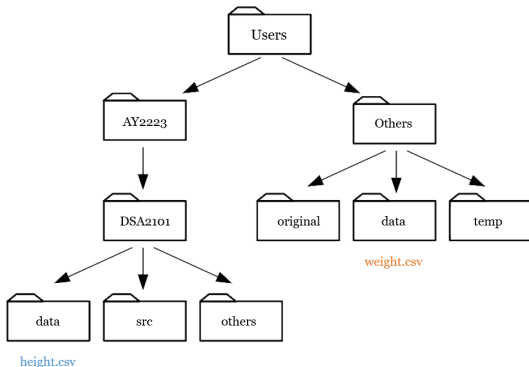
The function returns the **absolute path** of the current working directory.

File path

The second question implies that data are not necessarily stored at the location of our current working directory.

- ▶ Absolute path: the exact address of a file on our computer.
- ▶ **Relative path:** the address of a file relative to our current working directory.
 - ▶ Access files directly in the current working path.
 - ▶ Use two dots `..` to denote “one level up in the directory hierarchy”.

Using relative path in all code you write. This allows you to share your scripts and data files easily with others.



Let's say `getwd()` gives us `C:/Users/AY2223/DSA2101/src`

- ▶ To access **height** data: `../data/height.csv`
- ▶ To access **weight** data: `../../../Others/data/weight.csv`

File path

We strongly recommend the following practice:

- ▶ Create a folder **DSA2101** for our class and store all the code, data, and markdown files inside.
- ▶ Within DSA2101, create a folder called **src** to store all your source codes and **Rmd** files.
- ▶ Within DSA2101, create another folder called **data** to store all your data files.
- ▶ The **src** and **data** folders should be at the same level.

Use **relative path** in all code you write.

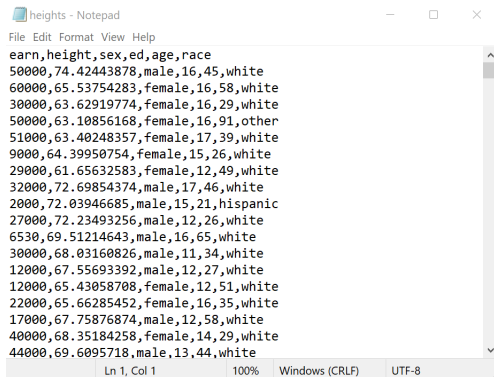
CSV files

CSV stands for Comma-separated values.

- ▶ These files are in fact just text files, with
 - ▶ an optional header, listing the column names.
 - ▶ each observation separated by commas within each row
- ▶ CSV is the easiest format to read into R.

What does a CSV file look like?

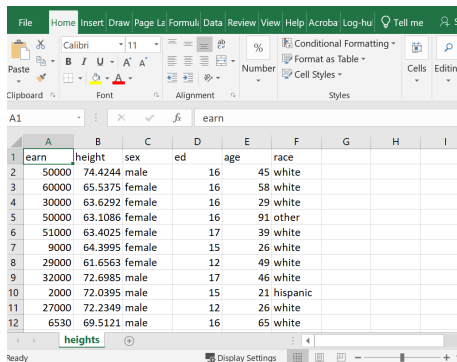
A .csv file, opened in a text editor:



```
heights - Notepad
File Edit Format View Help
earn,height,sex,ed,age,race
50000,74.42443878,male,16,45,white
60000,65.53754283,female,16,58,white
30000,63.62919774,female,16,29,white
50000,63.10856168,female,16,91,other
51000,63.40248357,female,17,39,white
9000,64.39950754,female,15,26,white
29000,61.65632583,female,12,49,white
32000,72.69854374,male,17,46,white
2000,72.03946685,male,15,21,hispanic
27000,72.23493256,male,12,26,white
6530,69.51214643,male,16,65,white
30000,68.03160826,male,11,34,white
12000,67.55693392,male,12,27,white
12000,65.43058708,female,12,51,white
22000,65.66285452,female,16,35,white
17000,67.75876874,male,12,58,white
40000,68.35184258,female,14,29,white
44000,69.6095718,male,13,44,white
Ln 1, Col 1 100% Windows (CRLF) UTF-8
```

What does a CSV file look like?

Here is the same file opened in Excel:



The screenshot shows the Microsoft Excel interface with the 'File' tab selected. The ribbon includes 'Home', 'Insert', 'Draw', 'Page Layout', 'Formulas', 'Data', 'Review', 'View', 'Help', 'Acroba', 'Log-hu', and 'Tell me'. The 'Home' ribbon is active, showing options for Clipboard, Font, Alignment, Number, Styles, Cells, and Editing. The formula bar shows 'A1' and the formula 'earn'. The spreadsheet contains the following data:

	A	B	C	D	E	F	G	H	I
1	earn	height	sex	ed	age	race			
2	50000	74.4244	male		16	45 white			
3	60000	65.5375	female		16	58 white			
4	30000	63.6292	female		16	29 white			
5	50000	63.1086	female		16	91 other			
6	51000	63.4025	female		17	39 white			
7	9000	64.3995	female		15	26 white			
8	29000	61.6563	female		12	49 white			
9	32000	72.6985	male		17	46 white			
10	2000	72.0395	male		15	21 hispanic			
11	27000	72.2349	male		12	26 white			
12	6530	69.5121	male		16	65 white			

Read a CSV file into R

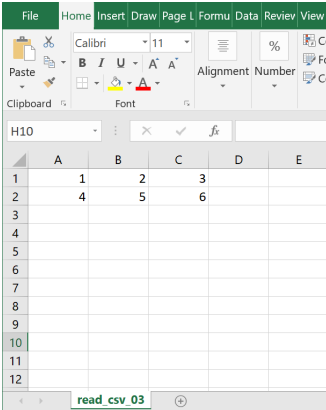
The command to read a CSV file into R is `read.csv()`

The main arguments to this function are:

- ▶ `file`: the file name.
- ▶ `header`: absence / presence of a header row.
- ▶ `skip`: number of lines at the beginning to skip.
- ▶ `col.names`: the names to identify columns in the table.
- ▶ `stringsAsFactors`: whether to convert character vectors to factors.
- ▶ `na.strings`: specify a character vector to be interpreted as NA values.

Example: A simple CSV file

- ▶ Take a first look at the data.
- ▶ 2 rows \times 3 columns.
- ▶ The data set has no header.



The screenshot shows the Microsoft Excel interface. The ribbon at the top includes File, Home, Insert, Draw, Page Layout, Formulas, Data, Review, and View. The Home ribbon is active, showing the Clipboard, Font, and Paragraph groups. The font is set to Calibri, size 11. The active cell is H10. The spreadsheet contains data in columns A, B, and C, with rows 1 and 2 populated. The data is as follows:

	A	B	C	D	E
1	1	2	3		
2	4	5	6		
3					
4					
5					
6					
7					
8					
9					
10					
11					
12					

The status bar at the bottom shows the file name 'read_csv_03' and a plus sign icon.

Example: A simple CSV file

```
df = read.csv("../data/read_csv_03.csv", header = FALSE,  
              col.names = c("a", "b", "c"))  
df
```

```
##    a b c  
## 1 1 2 3  
## 2 4 5 6
```

- ▶ This file does not contain a header row, thus `header = FALSE`
- ▶ We can name the column as `a`, `b`, `c`. If we do not supply column names, R will name the columns by itself.

Example: Education, Height, and Income

`heights.csv` contains information on 1192 individuals.

- ▶ Take a look at the data, you will find that it contains 6 columns and 1 header.
- ▶ Hence, we read in the data in the following way:

```
heights = read.csv("../data/heights.csv",  
                    header = TRUE, stringsAsFactors = TRUE)  
dim(heights)
```

```
## [1] 1192    6
```

- ▶ The function `dim()` (stands for **dimensions**) tells us that the data frame has 1192 rows and 6 columns.

Data checks

1. What type has each column been read in as?

```
str(heights)
```

```
## 'data.frame':    1192 obs. of  6 variables:
## $ earn   : num  50000 60000 30000 50000 51000 9000 29000 32000 2000 27000 ..
## $ height: num  74.4 65.5 63.6 63.1 63.4 ...
## $ sex    : Factor w/ 2 levels "female","male": 2 1 1 1 1 1 1 2 2 2 ...
## $ ed     : int   16 16 16 16 17 15 12 17 15 12 ...
## $ age    : int   45 58 29 91 39 26 49 46 21 26 ...
## $ race   : Factor w/ 4 levels "black","hispanic",...: 4 4 4 3 4 4 4 4 2 4 ...
```

- ▶ The function `str()` (stands for **structure**) reveals information about the columns, giving the names of the columns and a peek into the contents of each.
- ▶ We can see that the data types make sense.

Data checks

2. `race` is a categorical variable (a **factor** class in R). What are the different races that have been read in?

```
levels(heights$race)
```

```
## [1] "black"      "hispanic" "other"     "white"
```

- ▶ The function `levels()` returns the level of a factor variable.
- ▶ Recall that the dollar sign `$` extracts variable from a data frame.

Data checks

3. Are there any missing values in the data?

```
sum(is.na(heights))
```

```
## [1] 0
```

- ▶ Use `is.na()` to check missing entries in the entire data set.
- ▶ If there are missing values, we would also like to know which variable contains missing value.

```
apply(heights, 2, function(x) sum(is.na(x)))
```

```
##   earn height   sex   ed   age   race  
##     0       0     0     0     0     0
```

Summary statistics

We can compute summary statistics for `earn`:

```
summary(heights$earn)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##       200   10000   20000   23155   30000  200000
```

Group statistics with `tapply()`:

```
tapply(heights$earn, heights$sex, mean)
```

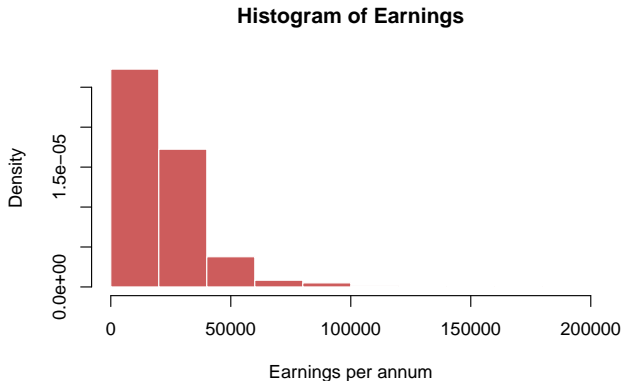
```
##   female    male
## 18280.20 29786.13
```

Histogram

Let us plot a histogram of income earned by individuals.

- ▶ A histogram divides the range of quantitative values into bins, then counts the number of values that fall into each bin.
- ▶ By default, the height of each bar represents frequencies.
- ▶ `freq = FALSE` alters a histogram such that the height represents the probability densities (that is, the histogram has a total area of one).

```
hist(heights$earn, freq = FALSE,  
      main = "Histogram of Earnings", xlab = "Earnings per annum",  
      col = "indianred", border = "white")
```



- The distribution of income is right-skewed, as expected.

Histogram (revised code)

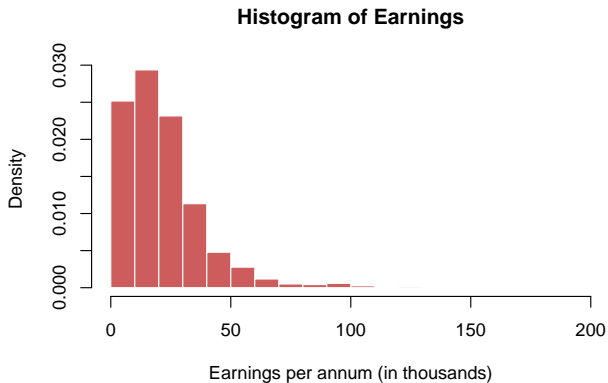
Our presentation of the histogram can be improved:

1. Disable scientific notations.
2. The bins correspond to intervals of width 20,000. Perhaps we would like bins of width 10,000 instead.

```
options(scipen = 999)
hist(heights$earn/1000, freq = FALSE,
      main = "Histogram of Earnings", xlab = "Earnings per annum (in thousands)",
      col = "indianred", border = "white", breaks = seq(0, 200, by = 10))
```

- ▶ `heights$earn/1000` divides earnings by a thousand. Now the earnings value ranges from 0 to 200.
- ▶ `breaks = seq(0, 200, by = 10)` sets the range of the x-axis from 0 to 200, and split it into bins with width 10.

Histogram (revised code)



The income distribution

Who are those high-earning individuals – earn more than 100K a year?

```
# install.packages("tidyverse")  
library(tidyverse)  
filter(heights, earn > 1e5)
```

```
##      earn  height    sex ed age  race  
## 1 125000 74.34062  male 18  45 white  
## 2 170000 71.01003  male 18  45 white  
## 3 175000 70.58955  male 16  48 white  
## 4 148000 66.74020  male 18  38 white  
## 5 110000 65.96504  male 18  37 white  
## 6 105000 74.58005  male 12  49 white  
## 7 123000 61.42908 female 14  58 white  
## 8 200000 69.66276  male 18  34 white  
## 9 110000 66.31203 female 18  48 other
```

The income distribution

The code on the previous slide uses the `tidyverse` syntax.

- ▶ It is a excellent tool for cleaning data.
- ▶ We shall study it very soon in Week 5.
- ▶ For now, only need to understand that it **filters out** irrelevant rows from the `heights` data frame, keeping only those who earned more than 10^5 per year.

Tncome distribution

From the new histogram, it is easier to tell that more than 50% of the individuals earned less than 20K per year. This calculation comes from

$$0.025 \times 10 + 0.03 \times 10$$

- ▶ Also, more than 90% of individuals earned less than 50K per year.
- ▶ From inspecting the outliers,
 - ▶ Only two females earned more than 100K per year, comapred to seven males.
 - ▶ None of those earning more than 100K were black or hispanic.
 - ▶ The highest earner is also the youngest guy in the group.

Re-cap

- ▶ Remember that you should inspect your data before and after you read them in.
- ▶ Try to think of as many ways in which it could have gone wrong and check them.
- ▶ As we covered here, you should at least consider the following:
 - ▶ Correct number of rows and columns
 - ▶ Column variables read in with the correct class type
 - ▶ Missing values

Excel files

To read data from `xls` and `xlsx` files, we need the `readxl` package.

```
# install.packages("readxl")  
library(readxl)
```

- ▶ The `read_xlsx()` function automatically detects the rectangle region that contains non-empty cells in the Excel spreadsheet.
- ▶ Nonetheless, ensure that you open up your file in Excel first, to see what it contains and how you can provide further contextual information for the function to use.

Excel example

Let us see a simple example.

```
read_excel("../data/read_excel_01.xlsx")
```

```
## # A tibble: 7 x 5
##   'Table 1' ...2 ...3 ...4 ...5
##   <lgl>      <lgl> <chr> <dbl> <chr>
## 1 NA      NA    <NA>    NA <NA>
## 2 NA      NA    <NA>    NA <NA>
## 3 NA      NA    <NA>    NA <NA>
## 4 NA      NA    <NA>    NA <NA>
## 5 NA      NA    a        1 m
## 6 NA      NA    b        2 m
## 7 NA      NA    c        3 m
```

In this case, `read_excel()` needs a little help as the data seems to be “floating” in the center of the worksheet.

Excel example

```
read_excel("../data/read_excel_01.xlsx", skip = 5)
```

```
## # A tibble: 2 x 3
##   a          '1' m
##   <chr> <dbl> <chr>
## 1 b          2 m
## 2 c          3 m
```

- ▶ The `skip` argument tells R to skip a certain number of rows.
- ▶ Looks like the function is reading the first row as the header. We can disable it by specifying `col_names = FALSE`.
- ▶ Notice that `read_excel()` uses a `col_names` argument, instead of `header`.

Excel example

Another way is to specify the data range exactly.

```
read_excel("../data/read_excel_01.xlsx", range = "C6:E8", col_names = FALSE)
```

```
## # A tibble: 3 x 3
##   ...1    ...2 ...3
##   <chr> <dbl> <chr>
## 1 a          1 m
## 2 b          2 m
## 3 c          3 m
```

- In case you were wondering, a **tibble** is an improved version of a data frame. We shall learn more about it in Week 5.

Example: UNESCAP data

The excel file `UNESCAP_population_2010_2015.xlsx` contains population counts for the Asia-Pacific countries.

- ▶ The counts are broken down by age group and gender.
- ▶ In the file, data for each age group and gender are stored in different spreadsheets.
- ▶ Suppose we want to read in data for **females aged 0–14 years** and combine them into one data set.

UNESCAP data on population

Read in data for female aged 0–4 years old.

```
female_0_4 = read_excel("../data/UNESCAP_population_2010_2015.xlsx", sheet = 3)
head(female_0_4)
```

```
## # A tibble: 6 x 7
##   e_fname      Y2010 Y2011 Y2012 Y2013 Y2014 Y2015
##   <chr>      <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
## 1 Afghanistan 2447  2459  2454  2438  2422  2412
## 2 Armenia      92    94    97    99   101   101
## 3 Australia    710   731   740   743   745   752
## 4 Azerbaijan   333   348   370   394   413   425
## 5 Bangladesh  7725  7622  7565  7540  7525  7503
## 6 Bhutan        35    35    35    34    33    32
```


RDS file

RDS is R's own data file format. To read it in, we invoke the `readRDS()` function.

Last time, we used a RDS file that contains a nested list.

```
hawkers = readRDS("../data/hawker_ctr_raw.rds")
str(hawkers[[1]][[2]], max.level = 1)
```

```
## List of 12
## $ ADDRESSBUILDINGNAME : chr ""
## $ ADDRESSFLOORNUMBER : chr ""
## $ ADDRESSPOSTALCODE   : chr "141001"
## $ ADDRESSSTREETNAME   : chr "Commonwealth Drive"
## $ ADDRESSUNITNUMBER   : chr ""
## $ DESCRIPTION          : chr "HUP Standard Upgrading"
## $ HYPERLINK            : chr ""
## $ NAME                 : chr "Blks 1A/ 2A/ 3A Commonwealth Drive"
## $ PHOTOURL             : chr ""
## $ ADDRESSBLOCKHOUSENUMBER: chr "1A/2A/3A"
## $ XY                   : chr "24055.5,31341.24"
## $ ICON_NAME            : chr "HC icons_Opt 8.jpg"
```

Retrieving street names

Remove the first sublist in hawkers

```
hawkers_116 = hawkers[[1]][-1]
```

Name	Type	Value
hawkers	list [1]	List of length 1
SrcResults	list [117]	List of length 117
[[1]]	list [1]	List of length 1
[[2]]	list [12]	List of length 12
ADDRESSBUILDING...	character [1]	"
ADDRESSFLOORN...	character [1]	"
ADDRESSPOSTALC...	character [1]	"141001"
ADDRESSSTREETN...	character [1]	"Commonwealth Drive"
ADDRESSUNITNUM...	character [1]	"
DESCRIPTION	character [1]	"HUP Standard Upgrading"
HYPERLINK	character [1]	"
NAME	character [1]	"Blks 1A/ 2A/ 3A Commonwealth Drive"
PHOTOURL	character [1]	"
ADDRESSBLOCKHO...	character [1]	"1A/2A/3A"
XY	character [1]	"24055.5,31341.24"
ICON_NAME	character [1]	"HC icons_Opt 8.jpg"
[[3]]	list [12]	List of length 12

Name	Type	Value
hawkers_116	list [116]	List of length 116
[[1]]	list [12]	List of length 12
ADDRESSBUILDING...	character [1]	"
ADDRESSFLOORN...	character [1]	"
ADDRESSPOSTALCODE	character [1]	"141001"
ADDRESSSTREETNAME	character [1]	"Commonwealth Drive"
ADDRESSUNITNUM...	character [1]	"
DESCRIPTION	character [1]	"HUP Standard Upgrading"
HYPERLINK	character [1]	"
NAME	character [1]	"Blks 1A/ 2A/ 3A Commonwealth Drive"
PHOTOURL	character [1]	"
ADDRESSBLOCKHOU...	character [1]	"1A/2A/3A"
XY	character [1]	"24055.5,31341.24"
ICON_NAME	character [1]	"HC icons_Opt 8.jpg"
[[2]]	list [12]	List of length 12
[[3]]	list [12]	List of length 12

Retrieving street names

The new object `hawkers_116` contains 116 lists, each has 12 components.

- ▶ Retrieve the street names of the first component with the following

```
hawkers_116[[1]]$ADDRESSSTREETNAME
```

```
## [1] "Commonwealth Drive"
```

- ▶ The following code produces the same output.

```
hawkers_116[[1]][[4]]
```

Retrieving street names

To retrieve all street names, use `sapply()` with an anonymous function to store them in a vector.

```
street_name = sapply(hawkers_116, function(x) x$ADDRESSSTREETNAME)
head(street_name, n = 10)
```

```
## [1] "Commonwealth Drive"      "Marsiling Lane"          "Boon Lay Place"
## [4] "Havelock Road"           "Circuit Road"            "Whampoa Drive"
## [7] "Upper Bukit Timah Road"  "Smith Street"            "Kensington Park Road"
## [10] "Yishun Ring Road"
```

Converting to a data frame

Using the same trick on different components in the object, we can store variables in different vectors, and then combine them as a new data frame.

```
postal_code = sapply(hawkers_116, function(x) x$ADDRESSPOSTALCODE)
name = sapply(hawkers_116, function(x) x$NAME)
coordinates = sapply(hawkers_116, function(x) x$XY)

hawkers_df = data.frame(postal_code, name, coordinates)
head(hawkers_df, n = 4)
```

##	postal_code	name	coordinates
## 1	141001	Blks 1A/ 2A/ 3A Commonwealth Drive	24055.5,31341.24
## 2	730020	Blks 20/21 Marsiling Lane	21755.23,47282.71
## 3	641221	Blks 221A/B Boon Lay Place	14587.57,36373.7899
## 4	161022	Blks 22A/B Havelock Road	27589.1399,30043.3

JavaScript Object Notation (JSON)

JSON (JavaScript Object Notation) is a standard **text-based format** for storing structured data.

- ▶ On the internet, it is a very popular format for data interchange.
- ▶ The full description of the format can be found at <http://www.json.org/>
- ▶ The syntax is easy for humans to read and write, and for computers to parse and generate.

We shall work with the `jsonlite` package.

```
# install.packages("jsonlite")  
library(jsonlite)
```

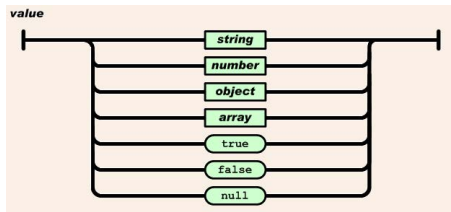
JSON description

- ▶ JSON is built on two structures:
 - ▶ An **object** is an unordered collection of name/value pairs.
 - ▶ An **array** is an ordered list of values.
- ▶ By repeatedly stacking these structures on top of one another, we will be able to store quite complex data structures.

```
object
  {}
  { members }
members
  pair
  pair , members
pair
  string : value
array
  []
  [ elements ]
elements
  value
  value , elements
value
  string
  number
  object
  array
  true
  false
  null
```

JSON value

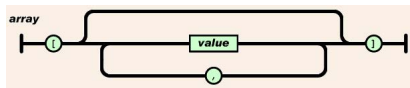
A **value** can be a string (in double quotes), a number, an object, an array, or a true or false or null.



JSON array

An **array** is an ordered collection of values.

- ▶ Surrounded with square brackets, starts with [and ends with]
- ▶ Values are separated by a comma ,



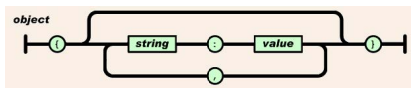
Example:

- ▶ [12, 3, 7] is an JSON array with three elements, all are numbers.
- ▶ ["Hello", 3, 7] is also valid.

JSON object

An **object** is an unordered set of name/value pairs.

- ▶ Surrounded with curly braces, starts with { and ends with }
- ▶ Each name is followed by a colon : and the name/value pairs are separated by a comma ,



Example:

- ▶ {"fruit": "Apple"} is a valid JSON object.
- ▶ {"fruit": "Apple", "price": 2.03} is also valid.
 - ▶ Two name/value pairs. The names are “fruit” and “price”.

Read JSON objects in R

The `fromJSON()` function in the `jsonlite` package allows us to read JSON from files, the web, or straight from the console.

1. In the following example, `fromJSON()` detects that the values are homogeneous and so reads them into a numeric vector.

```
txt = "[12, 3, 7]"  
fromJSON(txt)
```

```
## [1] 12 3 7
```

Read JSON objects in R

2. In this case, the values are not all of the same type. So the function reads them in as a character vector.

```
txt2 = '[12, "a", 7]'  
fromJSON(txt2)
```

```
## [1] "12" "a"  "7"
```

3. The missing value is coded as NA within R.

```
txt3 = '[12, null, 7]'  
fromJSON(txt3)
```

```
## [1] 12 NA  7
```

Read a single JSON object from a file

- ▶ JSON stores everything as a text.
- ▶ If we are sure that the `txt` file only contains one JSON object, we can use the command `fromJSON()`.

```
fromJSON("../data/read_json_01.txt")
```

```
## $fruit  
## [1] "Apple"  
##  
## $price  
## [1] 2.03  
##  
## $shelf  
## [1] "lower" "middle"
```

Read multiple JSON objects from a file

- If the file has multiple JSON objects, we need to first read each line into R using `readLines()`, and then apply `fromJSON()` to each of them.

```
all_lines = readLines("../data/read_json_02.txt")
json_list = lapply(all_lines, fromJSON)
str(json_list)
```

```
## List of 3
## $ :List of 3
## ..$ fruit: chr "Apple"
## ..$ price: num 2.03
## ..$ shelf: chr [1:2] "lower" "middle"
## $ :List of 3
## ..$ fruit: chr "Orange"
## ..$ price: num 1.03
## ..$ shelf: chr [1:2] "middle" "upper"
## $ :List of 3
## ..$ fruit: chr "Watermelon"
## ..$ price: num 0.99
## ..$ shelf: chr "lower"
```

Convert to a data frame

The next step is to convert it into a data frame.

- ▶ Notice that watermelons can only be stored on the lower shelf, but the other two fruits can be stored in two possible shelves.
- ▶ How should the data frame look like?

fruit	price	shelf
Apple	2.03	lower, middle
Orange	1.03	middle, upper
Watermelon	0.99	lower



OR

fruit	price	lower	middle	upper
Apple	2.03	1	1	0
Orange	1.03	0	1	1
Watermelon	0.99	1	0	0



Convert to a data frame

Let us first write a function (`convert_2_df`) that takes one component at a time and then converts it to a data frame.

```
convert_2_df = function(x) {  
  
  lower  = ifelse("lower" %in% x$shelf, 1, 0)  
  middle = ifelse("middle" %in% x$shelf, 1, 0)  
  upper  = ifelse("upper" %in% x$shelf, 1, 0)  
  
  data.frame(fruit = x$fruit, price = x$price, lower, middle, upper)  
}
```


Convert to a data frame

Apply this new function `convert_2_df` to the list `json_list` to obtain a **list** of three data frames.

```
df_row = lapply(json_list, convert_2_df)
df_row
```

```
## [[1]]
##   fruit price lower middle upper
## 1 Apple  2.03      1       1     0
##
## [[2]]
##   fruit price lower middle upper
## 1 Orange  1.03      0       1     1
##
## [[3]]
##   fruit price lower middle upper
## 1 Watermelon 0.99      1       0     0
```

Convert to a data frame

We then combine these individual rows into one single data frame using `rbind()`.

```
df_fruit = rbind(df_row[[1]], df_row[[2]], df_row[[3]])  
df_fruit
```

```
##           fruit price lower middle upper  
## 1      Apple  2.03     1      1      0  
## 2     Orange  1.03     0      1      1  
## 3 Watermelon  0.99     1      0      0
```

Data from the web

We can read data files directly from a website to R.

TidyTuesday is a weekly social data project in R born out of the *R for Data Science* textbook and its online learning community.

- ▶ It posts raw data set(s) and a related article every week.
- ▶ Emphasizes on the understanding of how to summarize and arrange data to make meaningful visuals in the **tidyverse** ecosystem.

Full list of data sets can be found on

<https://github.com/rfordatascience/tidytuesday>

TidyTuesday data

Let's explore the data set posted on April 20, 2021.

- ▶ A data set on TV shows and movies available on Netflix.
- ▶ You can find an overview of the data at:

<https://github.com/rfordatascience/tidytuesday/blob/master/data/2021/2021-04-20/readme.md>



Follow the instruction to get the data.

- Read in the data with the `tidytuesdayR` package:

```
# install.packages("tidytuesdayR")
tuesdata = tidytuesdayR::tt_load("2021-04-20")
```

```
##
## Downloading file 1 of 1: 'netflix_titles.csv'
```

```
netflix = tuesdata$netflix_titles
head(netflix, n = 4)
```

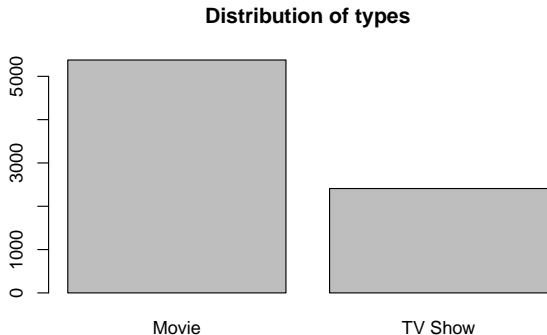
```
## # A tibble: 4 x 12
##   show_id type    title director    cast    country date_~1 relea~2 rating dur
##   <chr>   <chr>   <chr> <chr>      <chr> <chr>   <chr>      <dbl> <chr> <chr>
## 1 s1      TV Show 3%    <NA>      João~ Brazil August~    2020 TV-MA 4 S
## 2 s2      Movie  7:19  Jorge Mich~ Demi~ Mexico Decemb~    2016 TV-MA 93
## 3 s3      Movie  23:59 Gilbert Ch~ Tedd~ Singap~ Decemb~    2011 R      78
## 4 s4      Movie   9     Shane Acker Elij~ United~ Novemb~    2009 PG-13 80
## # ... with 2 more variables: listed_in <chr>, description <chr>, and
## # abbreviated variable names 1: date_added, 2: release_year, 3: duration
## # i Use 'colnames()' to see all variable names
```

```
summary(netflix)
```

```
##      show_id          type          title          director
## Length:7787      Length:7787      Length:7787      Length:7787
## Class :character  Class :character  Class :character  Class :character
## Mode  :character  Mode  :character  Mode  :character  Mode  :character
##
##
##
##      cast          country          date_added      release_year
## Length:7787      Length:7787      Length:7787      Min.   :1925
## Class :character  Class :character  Class :character  1st Qu.:2013
## Mode  :character  Mode  :character  Mode  :character  Median :2017
##                                     Mean   :2014
##                                     3rd Qu.:2018
##                                     Max.   :2021
##
##      rating          duration          listed_in      description
## Length:7787      Length:7787      Length:7787      Length:7787
## Class :character  Class :character  Class :character  Class :character
## Mode  :character  Mode  :character  Mode  :character  Mode  :character
##
##
##
```

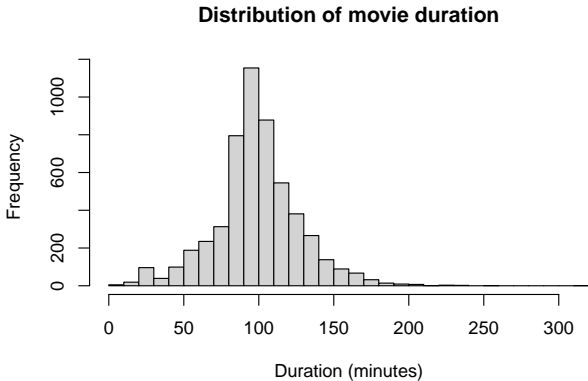
A bar plot on the types of Netflix titles.

```
netflix$type = as.factor(netflix$type)  
barplot(table(netflix$type), main = "Distribution of types")
```



A histogram of movie duration.

```
movie = netflix[netflix$type == "Movie", ]  
movie$minute = as.numeric(gsub(" min", "", movie$duration))  
hist(movie$minute, breaks = 30,  
      xlab = "Duration (minutes)", main = "Distribution of movie duration")
```



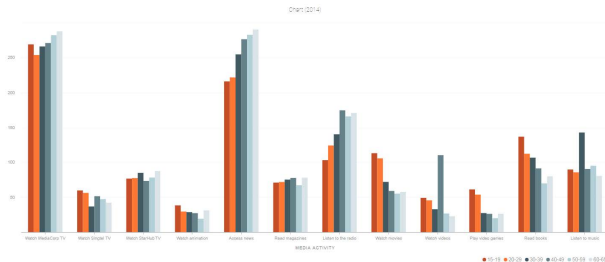


- ▶ data.gov.sg was launched in 2011 as Singapore's one-stop open data portal.
- ▶ Data sets from 70 government agencies, in the field of economy, education, environment, finance, health, infrastructure, etc.
- ▶ From the website, data sets can be downloaded in `csv` format. It is also possible to download the data using a script. The data would then be returned as a `JSON` object.

Media usage data from IMDA

Every year, the Infocomm Media Development Authority (IMDA) commissions a Media Consumer Experience Study.

- ▶ The data set describes the percentage of consumers who have ever used a traditional media device (e.g., TV, newspaper) for media activities.



Download IMDA data

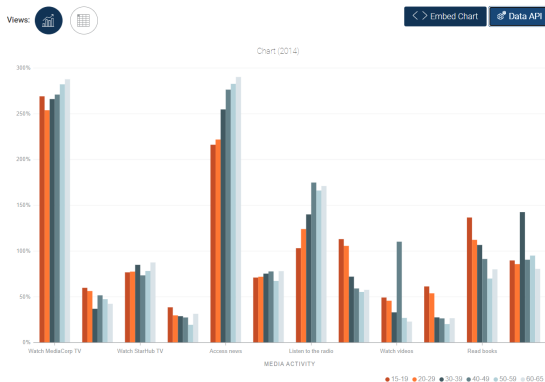
Now we demonstrate how to download the data through R.

- ▶ Instructions on downloading the data can be found in the **Developers** sub-page on the website.
- ▶ Essentially what is needed is to identify the **resource id** for this data set, and then tag it onto a template URL.
- ▶ The **Data API** link on the data set page shows the resource id for this data.
- ▶ However, there is a limit on the number of records that can be retrieved per query. Thus it is necessary to run a loop until all records have been retrieved.

- Visit the data web page:

https://data.gov.sg/dataset/usage-of-traditional-media-devices?view_id=c1b506c9-bf0e-4601-8ef1-f9ec9c533c61&resource_id=50c2820a-a18f-4514-9bb0-ff43048ddff5

- Then click **Data API** on the upper right corner.



Download IMDA data

- It will show you the resource id for this data set.

View Embedding Instructions



Endpoints »

The Data API can be accessed via the following actions of the CKAN action API.

Query	<code>https://data.gov.sg/api/action/datastore_search</code>
--------------	--

Querying »

Query example (first 5 results)

```
https://data.gov.sg/api/action/datastore_search?resource_id=50c2820a-a18f-4514-9bb0-ff43048ddff5&limit=5
```

Query example (results containing 'jones')

```
https://data.gov.sg/api/action/datastore_search?resource_id=50c2820a-a18f-4514-9bb0-ff43048ddff5&q=jones
```

Example: Javascript »

Example: Python »

Download IMDA data

```
root_url = "https://data.gov.sg"
url1 = paste(root_url,
              "/api/action/datastore_search?",
              "resource_id=50c2820a-a18f-4514-9bb0-ff43048ddff5", sep = "")
media_json1 = fromJSON(url1)
str(media_json1, max.level = 2)
```

```
## List of 3
## $ help : chr "https://data.gov.sg/api/3/action/help_show?name=datastore_s
## $ success: logi TRUE
## $ result :List of 5
## ..$ resource_id: chr "50c2820a-a18f-4514-9bb0-ff43048ddff5"
## ..$ fields : 'data.frame': 6 obs. of 2 variables:
## ..$ records : 'data.frame': 100 obs. of 6 variables:
## ..$ _links :List of 2
## ..$ total : int 210
```

Download IMDA data

The previous slide tells us that we have managed to retrieve a data frame with 100 rows.

- ▶ However, this component tells us that the final data set should contain 210 rows.

```
media_json1$result$total
```

```
## [1] 210
```

Download IMDA data

- ▶ The following two links inform us that we need to submit another query.

```
media_json1$result$`_links`
```

```
## $start
```

```
## [1] "/api/action/datastore_search?resource_id=50c2820a-a18f-4514-9bb0-ff4304
```

```
##
```

```
## $`next`
```

```
## [1] "/api/action/datastore_search?offset=100&resource_id=50c2820a-a18f-4514-
```

- ▶ Notice that in the second component of `_links`, it tells us to offset the first 100 rows in the next query.

Download IMDA data

- ▶ Continue to submit queries **until** the requisite number of rows are obtained.

```
media_data = media_json1$result$records
total_records = media_json1$result$total

while(nrow(media_data) < total_records) {

  url1 = paste(root_url,
               media_json1$result$`_links`[[2]],
               sep = "")
  media_json1 = fromJSON(url1)
  media_data = rbind(media_data, media_json1$result$records)
}
```

Download IMDA data

To confirm that we have the data now:

```
dim(media_data)
```

```
## [1] 210    6
```

```
str(media_data)
```

```
## 'data.frame':    210 obs. of  6 variables:
## $ ever_used      : chr  "97.1" "32.9" "39.5" "30.9" ...
## $ age            : chr  "15-19" "15-19" "15-19" "15-19" ...
## $ sample_size    : chr  "161" "161" "161" "161" ...
## $ year           : chr  "2013" "2013" "2013" "2013" ...
## $ _id            : int   1  2  3  4  5  6  7  8  9 10 ...
## $ media_activity: chr  "Watch MediaCorp TV" "Watch Singtel TV" "Watch StarH
```

Plotting IMDA data

Let us make a bar chart for the 20-29 years age group.

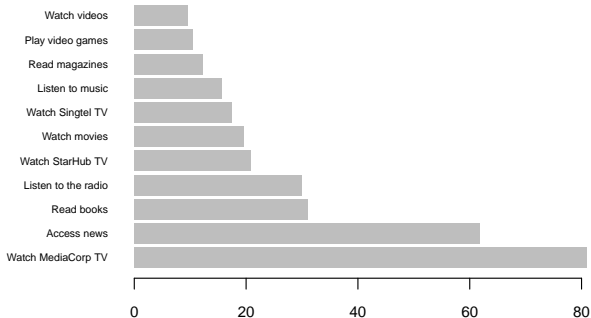
- ▶ The following code will become comprehensible after the next lecture. For now, I only need you to understand its purpose:
 - ▶ Filter and keep the groups we want.
 - ▶ Convert the variable `ever_used` from character to numeric.
 - ▶ Sort the data set by descending `pct`.

```
library(tidyverse)
young = filter(media_data, age == "20-29", year == 2015) %>%
  mutate(pct = as.numeric(ever_used)) %>%
  arrange(desc(pct))
head(young, n = 2)
```

##	ever_used	age	sample_size	year	_id	media_activity	pct
## 1	81	20-29	395	2015	156	Watch MediaCorp TV	81.0
## 2	61.8	20-29	395	2015	161	Access news	61.8

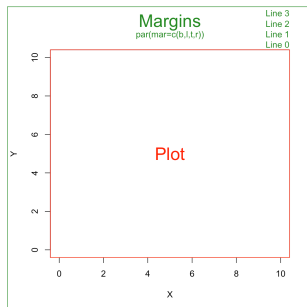
- Alter the arguments below and study their effects on the plot.

```
par(mar = c(5, 7, 2, 2))  
barplot(young$pct,  
        names.arg = young$media_activity,  
        horiz = TRUE, las = 1, cex.names = 0.6, cex.axis = 0.8, border = NA)
```



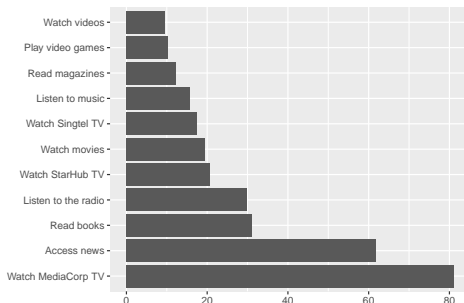
Adjusting the margins of the plot

- ▶ `par(mar = c(5, 7, 4, 2))`
on the previous slide specifies
the margins on the four sides
of the plot
- ▶ The default is `c(5, 4, 4, 2)`



The `ggplot()` way

```
ggplot(data = young, aes(x = reorder(media_activity, -pct), y = pct)) +  
  geom_bar(stat = "identity") +  
  coord_flip() + labs(x = "", y = "")
```



- Later in the semester, we shall learn about graphing with `ggplot()`.

Summary

We learn about importing data from different formats and sources:

1. CSV file using `read.csv()`
2. Excel file with `read_excel()` from the `readxl` package
3. R data file with `readRDS()`
4. JSON file with `fromJSON()` from the `jsonlite` package
5. Data from the web

Also a few more ways to clean data and make visualizations.

Summary

- ▶ Importing data becomes complicated when data is not stored in a friendly format.
- ▶ When reading data from the web, we need to have some creativity to identify patterns or keywords that can be used in a loop.
 - ▶ The paths and patterns are unlikely to be the same every time, but the experience you gather will help you along.