

## Module 2: Reporting, Data Wrangling and Graphing

Siyue Yang

04/26/2022

# Outline

We will review R, Rstudio, and Syntax of R together.

- ▶ LaTeX/Markdown
- ▶ Tidy data, processing (tidyverse)
- ▶ Graphing (ggplot2)

# LaTeX and Markdown

LaTeX is useful for documents with mathematical formulas.

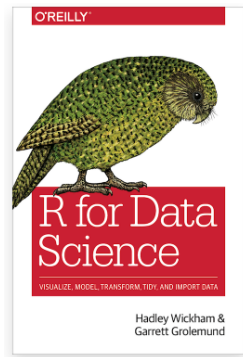
- ▶ Overleaf - an online, collaborative LaTeX editor
- ▶ LaTeX mathematical symbols
- ▶ Inline equation e.g. (`\alpha`) returns  $\alpha$
- ▶ Equation e.g. (`e = mc^2`) returns

$$e = mc^2$$

Markdown is appealing for formatting, e.g. headings, bold text, text with codes, ...

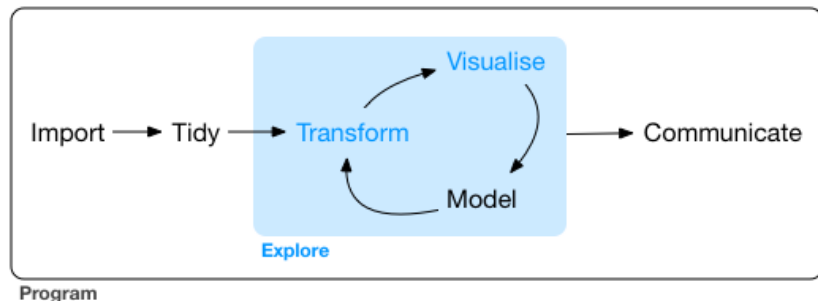
# Resources

“R for Data Science: Import, Tidy, Transform, Visualize, and Model Data” by Hadley Wickham.



# Let's code!

Data science project workflow:



## Data import

```
df <- read.table("mtcars.txt", header = TRUE)
head(df) # Show the first 6 rows.
```

##	Cntry	lper100k	weight	length
## 1	US	19.8	2178	5.92
## 2	Japan	9.9	1026	4.32
## 3	US	10.8	1188	4.27
## 4	US	12.5	1444	5.11
## 5	US	12.5	1485	5.03
## 6	US	12.5	1485	5.03

## Other options

CSV files.

- ▶ `read.csv()` in the base `r`.
- ▶ `read.csv()` in “`readr`” package (much faster).
- ▶ `fread()` in “`data.table`” package (much more faster).

Rdata.

- ▶ `load()` in the base `r`.

# Tidy data

The goal is to clean the dataset so it is much easier to use.

Specifically,

- ▶ Each variable must have its own column.
- ▶ Each observation must have its own row.
- ▶ Each value must have its own cell.

We will focus on the functions from “tidyverse” package.

```
library(tidyverse)
```



## Tidy data 1: pivoting

For a dataset having column names are not names of variables, but values of a variable, e.g.

```
table4a
```

```
## # A tibble: 3 x 3
##   country      '1999' '2000'
## * <chr>      <int>  <int>
## 1 Afghanistan    745    2666
## 2 Brazil        37737   80488
## 3 China          212258  213766
```

- ▶ Need to change 1999, 2000 to a column named as “year”.
- ▶ Need to change the values of 1999, 2000 as “cases”.

We can use `pivot_longer()` from the “tidyverse” package.

## Pivot longer

```
table4a %>%  
  pivot_longer(c(`1999`, `2000`),  
               names_to = "year", values_to = "cases")
```

```
## # A tibble: 6 x 3  
##   country      year  cases  
##   <chr>      <chr> <int>  
## 1 Afghanistan 1999     745  
## 2 Afghanistan 2000    2666  
## 3 Brazil      1999   37737  
## 4 Brazil      2000   80488  
## 5 China       1999  212258  
## 6 China       2000  213766
```

## Another example

```
table2 %>% head(5)
```

```
## # A tibble: 5 x 4
##   country      year type      count
##   <chr>      <int> <chr>    <int>
## 1 Afghanistan  1999 cases      745
## 2 Afghanistan  1999 population 19987071
## 3 Afghanistan  2000 cases      2666
## 4 Afghanistan  2000 population 20595360
## 5 Brazil       1999 cases      37737
```

- ▶ case and population are two variables and should be converted into columns.

We can use `pivot_wider()`.

## Pivot wider

```
table2 %>%  
  pivot_wider(names_from = type, values_from = count)
```

```
## # A tibble: 6 x 4  
##   country      year  cases population  
##   <chr>      <int> <int>      <int>  
## 1 Afghanistan 1999     745  19987071  
## 2 Afghanistan 2000    2666  20595360  
## 3 Brazil      1999   37737  172006362  
## 4 Brazil      2000   80488  174504898  
## 5 China       1999  212258  1272915272  
## 6 China       2000  213766  1280428583
```

# Transform data

Use the “pipes” from the “tidyverse” package, a powerful tool for clearly expressing a sequence of multiple operations, with the combination of the following functions:

- ▶ `select()`
- ▶ `filter()`
- ▶ `arrange()`
- ▶ `mutate()`
- ▶ `summarise()`
- ▶ `group_by()`

## Dataset - Diamonds

A dataset containing the prices and other attributes of almost 54,000 diamonds.

```
head(diamonds)
```

```
## # A tibble: 6 x 10
```

##	carat	cut	color	clarity	depth	table	price	x
##	<dbl>	<ord>	<ord>	<ord>	<dbl>	<dbl>	<int>	<dbl>
## 1	0.23	Ideal	E	SI2	61.5	55	326	3.95
## 2	0.21	Premium	E	SI1	59.8	61	326	3.89
## 3	0.23	Good	E	VS1	56.9	65	327	4.05
## 4	0.29	Premium	I	VS2	62.4	58	334	4.2
## 5	0.31	Good	J	SI2	63.3	58	335	4.34
## 6	0.24	Very Good	J	VVS2	62.8	57	336	3.94

# Select

Use `select()` to get a column, e.g. “color”

```
diamonds %>%  
  select(color)
```

In the base R, this is equivalent to `select()`

```
diamonds$color
```

Why bother use `select()` ?

- ▶ Because we can do a sequence of operations (later).

# Select

Use `select()` to remove a column, e.g. “color”

```
diamonds %>%  
  select(-color)
```

Need to assign the change to the original dataset, otherwise, the deletion won't affect the dataset:

```
diagmonds <- diamonds %>%  
  select(-color)
```



## Filter

Use `filter()` to filter by some condition, e.g. filter all price  $> 335$

```
diamonds %>%  
  filter(price > 335)
```

Filters with multiple conditions.

```
diamonds %>%  
  filter(price > 335 & depth < 64)
```

```
diamonds %>%  
  filter(cut == "Very Good" | cut == "Fair")
```

## Filter after select

```
diamonds %>%  
  select(price) %>%  
  filter(price > 335)
```

This is an example of “a sequence of operations”.

# Arrange

Use `arrange()` to order data.

```
diamonds %>%  
  arrange(price)
```

Arrange descending order, e.g. from the cheapest!

```
diamonds %>%  
  arrange(-price)
```

Arrange by multiple conditions.

```
diamonds %>%  
  arrange(price, cut)
```

## Filter, select, arrange

```
diamonds %>%  
  filter(table < 340) %>%  
  select(carat, cut, price) %>%  
  filter(price ) %>%  
  arrange(price, cut)
```



## Group by and Summarise

Use `group_by` and `summarise` to group variables:

```
diamonds %>%  
  group_by(cut) %>%  
  summarise(n = n())
```

```
## # A tibble: 5 x 2  
##   cut          n  
##   <ord>      <int>  
## 1 Fair        1610  
## 2 Good        4906  
## 3 Very Good  12082  
## 4 Premium    13791  
## 5 Ideal      21551
```

## More examples

```
diamonds %>%  
  group_by(cut) %>%  
  summarise(n = n(), price_avg = mean(price))
```

```
## # A tibble: 5 x 3  
##   cut          n price_avg  
##   <ord>      <int>     <dbl>  
## 1 Fair       1610     4359.  
## 2 Good       4906     3929.  
## 3 Very Good 12082     3982.  
## 4 Premium   13791     4584.  
## 5 Ideal     21551     3458.
```

# Proportions

```
diamonds %>%  
  group_by(cut) %>%  
  summarise(n = n(), price_avg = mean(price)) %>%  
  ungroup() %>%  
  mutate(prop = n/sum(n))
```

```
## # A tibble: 5 x 4  
##   cut          n price_avg  prop  
##   <ord>      <int>    <dbl> <dbl>  
## 1 Fair       1610    4359. 0.0298  
## 2 Good       4906    3929. 0.0910  
## 3 Very Good 12082    3982. 0.224  
## 4 Premium   13791    4584. 0.256  
## 5 Ideal     21551    3458. 0.400
```



## With percentage

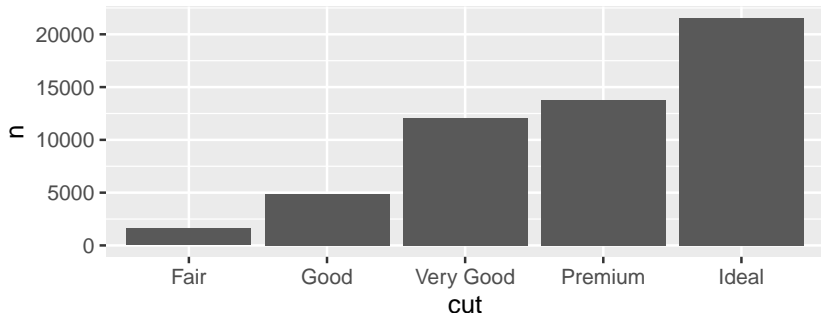
Use `scales::percent()` to add %.

```
diamonds %>%  
  group_by(cut) %>%  
  summarise(n = n(), price_avg = mean(price)) %>%  
  ungroup() %>%  
  mutate(prop = scales::percent(n/sum(n)))
```

```
## # A tibble: 5 x 4  
##   cut          n price_avg prop  
##   <ord>      <int>      <dbl> <chr>  
## 1 Fair      1610      4359. 3.0%  
## 2 Good      4906      3929. 9.1%  
## 3 Very Good 12082      3982. 22.4%  
## 4 Premium   13791      4584. 25.6%  
## 5 Ideal     21551      3458. 40.0%
```

## Graphing after transformation

```
diamonds %>%  
  group_by(cut) %>%  
  summarise(n = n(), price_avg = mean(price)) %>%  
  ggplot() +  
  geom_bar(aes(x = cut, y = n), stat = "identity")
```



# ggplot

Here we used functions from “ggplot2” package. Same pattern as “tidyverse”, but using “+” to connect.

How to write?

- ▶ Specify the data using `ggplot(data = diamonds)`
- ▶ Specify the x-/y-axis, `ggplot(data = diamonds, mapping = aes(x = cut))`
- ▶ Specify the types of plots with `geom`, e.g. `+ geom_bar()`

```
ggplot(data = diamonds) +  
  geom_bar(mapping = aes(x = cut))
```

## More plots

- ▶ `geom_histogram()`, `geom_density()`, `geom_line()`, `geom_point()`
- ▶ `geom_facet()` generates subplots
- ▶ color package
  - ▶ “RColorBrewer”
  - ▶ “ggsci”