# DATA WRANGLING WITH R

## BY DR. ALVIN ANG

## TABLE OF CONTENTS

## A.  SET YOUR WORKING DIRECTORY TO YOUR DOWNLOADS FOLDER



1.  CHECK YOUR CURRENT WORKING DIRECTORY



```
#---------------------------------------------------------------
# Data Wrangling using Core R by Dr. Alvin Ang
#---------------------------------------------------------------
#1. Check Current Working Directory

getwd()
```



```
> getwd()
[1] "C:/Users/User/Downloads"
```

## B. WRANGLING WEATHER.CSV

### 1. IMPORT CSV

https://www.alvinang.sg/s/weather.csv

```
#---------------------------------------------------------------
#2. Wrangling Weather.csv

#2a. Import CSV
#https://www.alvinang.sg/s/weather.csv
weather = read.csv('weather.csv',header = TRUE)
```

2. SLICE OUT COLUMN USING SUBSET

```
#2b.  Slice Out Column using Subset
#Slice out the "Ozone" column
weather.mayOzone <-
  subset(weather, select=Ozone, subset = Month==5)
```

3. CHECK WHICH ROWS HAVE NAS

```
#2c. Check which Rows have NAs
m = !is.na(weather.mayOzone)
```



Showing 1 to 10 of 31 entries, 1 total columns

## 4.   COMPUTE THE AVERAGE OZONE LEVEL IN THE MONTH OF MAY

```
#2d.   Compute the average Ozone level in the month of May
mean(weather.mayOzone[m])
```

5. FILTER OUT ALL NAS IN THE MONTH OF MAY

```
#2e.   Filter Out all NAs in the month of May
a = weather.mayOzone[m]
```

```
Values
  a                    int [1:26] 41 36 12 18 28 23 19 8 7 16 ...
```

6. OUTPUT AS CSV

```
#2f. Output as CSV
write.csv(a,'may_weather_data.csv')
```

| may_weather_data.csv - LibreOffice Calc |

| | A | B | C | D |
|---|---|---|---|---|
| 1 | | x | | |
| 2 | 1 | 41 | | |
| 3 | 2 | 36 | | |
| 4 | 3 | 12 | | |
| 5 | 4 | 18 | | |
| 6 | 5 | 28 | | |
| 7 | 6 | 23 | | |
| 8 | 7 | 19 | | |
| 9 | 8 | 8 | | |
| 10 | 9 | 7 | | |
| 11 | 10 | 16 | | |
| 12 | 11 | 11 | | |
| 13 | 12 | 14 | | |
| 14 | 13 | 18 | | |
| 15 | 14 | 14 | | |
| 16 | 15 | 34 | | |

may_weather_data

**C. WRANGLING MTCARS**

1. SLICING OUT MPG / AM / WT COLUMNS

```
#--------------------------------------------------
#3. Wrangling Mtcars

#3a. Slicing Out mpg / am / wt columns
b = mtcars[c('mpg','am','wt')]
```

| | mpg | am | wt |
|---|---|---|---|
| Mazda RX4 | 21.0 | 1 | 2.620 |
| Mazda RX4 Wag | 21.0 | 1 | 2.875 |
| Datsun 710 | 22.8 | 1 | 2.320 |
| Hornet 4 Drive | 21.4 | 0 | 3.215 |
| Hornet Sportabout | 18.7 | 0 | 3.440 |
| Valiant | 18.1 | 0 | 3.460 |
| Duster 360 | 14.3 | 0 | 3.570 |
| Merc 240D | 24.4 | 0 | 3.190 |
| Merc 230 | 22.8 | 0 | 3.150 |
| Merc 280 | 19.2 | 0 | 3.440 |

Showing 1 to 11 of 32 entries, 3 total columns

**Data**

| b | 32 obs. of 3 variables |
|---|---|
| m | logi [1:31, 1] TRUE TRUE TRUE TRUE FALSE TRUE... |
| weather | 153 obs. of 6 variables |
| weather.mayOzone | 31 obs. of 1 variable |

**Values**

| a | int [1:26] 41 36 12 18 28 23 19 8 7 16 ... |
|---|---|

2. PREVIEWING HEADS AND TAILS

```
#3b. Viewing Heads and Tails
head(b,7)
tail(b,3)
```

```
38    #3b. Viewing Heads and Tails
39    head(b,7)
40    tail(b,3)
41
42
40:1    # (Untitled)                    R Scri

Console   Terminal ×   Background Jobs ×
  R 4.2.1 · C:/Users/User/Downloads/
> b = mtcars[c('mpg','am','wt')]
> View(b)
> #3b. Viewing Heads and Tails
> head(b,7)
                    mpg am    wt
Mazda RX4          21.0  1 2.620
Mazda RX4 Wag      21.0  1 2.875
Datsun 710         22.8  1 2.320
Hornet 4 Drive     21.4  0 3.215
Hornet Sportabout  18.7  0 3.440
Valiant            18.1  0 3.460
Duster 360         14.3  0 3.570
>
```

```
38    #3b. Viewing Heads and Tails
39    head(b,7)
40    tail(b,3)
41
42
36:1    # (Untitled)                    R Script

Console   Terminal ×   Background Jobs ×
  R 4.2.1 · C:/Users/User/Downloads/
Mazda RX4          21.0  1 2.620
Mazda RX4 Wag      21.0  1 2.875
Datsun 710         22.8  1 2.320
Hornet 4 Drive     21.4  0 3.215
Hornet Sportabout  18.7  0 3.440
Valiant            18.1  0 3.460
Duster 360         14.3  0 3.570
> tail(b,3)
                mpg am    wt
Ferrari Dino   19.7  1 2.77
Maserati Bora  15.0  1 3.57
Volvo 142E     21.4  1 2.78
>
```

3. SLICING OUT MPG / HP COLUMNS

```
#3c. Slicing Out mpg / hp columns
c = subset(mtcars, select=c(mpg,hp))
```

| | mpg | hp |
|---|---|---|
| Mazda RX4 | 21.0 | 110 |
| Mazda RX4 Wag | 21.0 | 110 |
| Datsun 710 | 22.8 | 93 |
| Hornet 4 Drive | 21.4 | 110 |
| Hornet Sportabout | 18.7 | 175 |
| Valiant | 18.1 | 105 |
| Duster 360 | 14.3 | 245 |
| Merc 240D | 24.4 | 62 |
| Merc 230 | 22.8 | 95 |
| Merc 280 | 19.2 | 123 |
| Merc 280C | 17.8 | 123 |

Showing 1 to 12 of 32 entries, 2 total columns

4. OUTPUT AS CSV



```
#3d. Output as CSV
write.csv(c,"mtcarssubset.csv")
```



| | mpg | hp |
|---|---|---|
| Mazda RX4 | 21 | 110 |
| Mazda RX4 Wag | 21 | 110 |
| Datsun 710 | 22.8 | 93 |
| Hornet 4 Drive | 21.4 | 110 |
| Hornet Sportabout | 18.7 | 175 |
| Valiant | 18.1 | 105 |
| Duster 360 | 14.3 | 245 |
| Merc 240D | 24.4 | 62 |
| Merc 230 | 22.8 | 95 |
| Merc 280 | 19.2 | 123 |
| Merc 280C | 17.8 | 123 |
| Merc 450SE | 16.4 | 180 |
| Merc 450SL | 17.3 | 180 |
| Merc 450SLC | 15.2 | 180 |
| Cadillac Fleetwood | 10.4 | 205 |
| Lincoln Continental | 10.4 | 215 |
| Chrysler Imperial | 14.7 | 230 |
| Fiat 128 | 32.4 | 66 |

5. FILTER ALL THE MPG > 15 AND AM = 1

```
#3e. Filter all the mpg>15 and am=1
d = mtcars[mtcars$mpg>15 & mtcars$am==1,]
```

| | mpg | cyl | disp | hp | drat | wt | qsec | vs | am | gear |
|---|---|---|---|---|---|---|---|---|---|---|
| Mazda RX4 | 21.0 | 6 | 160.0 | 110 | 3.90 | 2.620 | 16.46 | 0 | 1 | |
| Mazda RX4 Wag | 21.0 | 6 | 160.0 | 110 | 3.90 | 2.875 | 17.02 | 0 | 1 | |
| Datsun 710 | 22.8 | 4 | 108.0 | 93 | 3.85 | 2.320 | 18.61 | 1 | 1 | |
| Fiat 128 | 32.4 | 4 | 78.7 | 66 | 4.08 | 2.200 | 19.47 | 1 | 1 | |
| Honda Civic | 30.4 | 4 | 75.7 | 52 | 4.93 | 1.615 | 18.52 | 1 | 1 | |
| Toyota Corolla | 33.9 | 4 | 71.1 | 65 | 4.22 | 1.835 | 19.90 | 1 | 1 | |
| Fiat X1-9 | 27.3 | 4 | 79.0 | 66 | 4.08 | 1.935 | 18.90 | 1 | 1 | |
| Porsche 914-2 | 26.0 | 4 | 120.3 | 91 | 4.43 | 2.140 | 16.70 | 0 | 1 | |
| Lotus Europa | 30.4 | 4 | 95.1 | 113 | 3.77 | 1.513 | 16.90 | 1 | 1 | |
| Ford Pantera L | 15.8 | 8 | 351.0 | 264 | 4.22 | 3.170 | 14.50 | 0 | 1 | |

```
#3f. Filter Out only mpg and am columns with am=1 (automatic)
e = subset(mtcars,
           select=c('mpg','am'),
           subset=am==1)
```

| | mpg | am |
|---|---|---|
| Mazda RX4 | 21.0 | 1 |
| Mazda RX4 Wag | 21.0 | 1 |
| Datsun 710 | 22.8 | 1 |
| Fiat 128 | 32.4 | 1 |
| Honda Civic | 30.4 | 1 |
| Toyota Corolla | 33.9 | 1 |
| Fiat X1-9 | 27.3 | 1 |
| Porsche 914-2 | 26.0 | 1 |
| Lotus Europa | 30.4 | 1 |
| Ford Pantera L | 15.8 | 1 |
| Ferrari Dino | 19.7 | 1 |

Showing 1 to 12 of 13 entries, 2 total columns

7. SUMMARY OF MTCARS SUBSET

```
#3g. Summary of mtcars subset
summary(e)
```

```
> summary(e)
       mpg                am
 Min.    :15.00    Min.      :1
 1st Qu.:21.00    1st Qu.:1
 Median :22.80    Median :1
 Mean    :24.39    Mean      :1
 3rd Qu.:30.40    3rd Qu.:1
 Max.    :33.90    Max.      :1
>
```

8. CREATE A TABLE FROM MTCARS AM COLUMNS

```
#3h. Create a Table from Mtcars AM columns
factor = factor(mtcars$am)
table(factor)


#---------------------------------------------------
#THE END
#---------------------------------------------------
```

```
factor                    Factor w/ 2 levels "0","1": 2 2 2 1 1 1 1 1 1 1 ...
```

```
> table(factor)
factor
 0  1
19 13
```

https://www.alvinang.sg/s/Data-Wrangling-with-Tidyverse-by-Dr-Alvin-Ang.R

Tons of great Data Wrangling with R here:

https://www.marsja.se/how-to-rename-column-or-columns-in-r-with-dplyr/

### A. INSTALLING TIDYVERSE PACKAGE

```r
#-----------------------------------------------
#Data Wrangling with Tidyverse by Dr Alvin Ang
#-----------------------------------------------
#1. Install Tidyverse and Load Libraries

install.packages("tidyverse", dependencies = TRUE)

library(tidyverse)
library(tibble)
library(tidyr)
library(dplyr)
library(readxl)
library(ggplot2)
library(lubridate)
```

**B.  READING IN CSV**

File can be found here: https://www.alvinang.sg/s/dengue.csv



- Do you know where you stored the dengue.csv downloaded file?

- Most probably is in your download folder.

- Make sure that you set the working directory to that folder (download folder)… so that it can import in the CSV.

```
#----------------------------------------------------------
#2. Reading in the Dengue.csv
dengue <- read_csv("dengue.csv")

#file is here: https://www.alvinang.sg/s/dengue.csv

#or if you want to read in .xls
# dengue_xls <- read_excel("dengue.xlsx")
```



| | year | eweek | type_dengue | number |
|---|------|-------|-------------|--------|
| 1 | 2014 | 1 | Dengue | 436 |
| 2 | 2014 | 1 | DHF | 1 |
| 3 | 2014 | 2 | Dengue | 479 |
| 4 | 2014 | 2 | DHF | 0 |
| 5 | 2014 | 3 | Dengue | 401 |
| 6 | 2014 | 3 | DHF | 0 |

Showing 1 to 6 of 530 entries, 4 total columns

dengue          530 obs. of 4 variables

```
#---------------------------------------------------
#3. Selecting Columns

# Select 'year' and 'number' columns from dengue.csv
a = dengue %>%
  select(year,number)
```

1. FILTER OUT 'YEAR' == 2018 FROM DENGUE.CSV

```
#----------------------------------------------
#4. Filter data

#Filter out 'year' == 2018 from dengue.csv
b = dengue %>%
  filter(year==2018)
```

**E. FILTER DATA BASED ON MULTIPLE CONDITIONS**

1. FILTER OUT 2017 AND 2018

```
#-------------------------------------------------------------
#5. Filter data based on multiple conditions

#5a. Filter out 2017 and 208
c = dengue %>%
  filter(year==2017 | year==2018 )
```

2. FILTER OUT 2018 AND DENGUE TYPE

```
#5b. Filter out 2018 and 'Dengue' type
d = dengue %>%
    filter(year==2018,type_dengue=='Dengue' )
```

3. ANOTHER WAY TO FILTER OUT 2018 AND DENGUE TYPE

```
#5c. Another way to Filter out 2018 and 'Dengue' type
e = dengue %>%
    filter(year==2018) %>%
    filter(type_dengue=='Dengue')
```
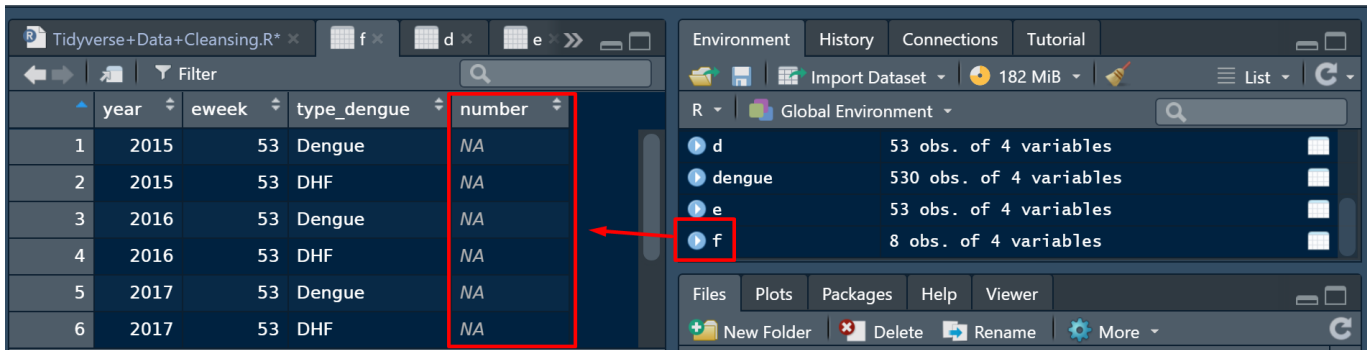
**F. HANDLING MISSING VALUES IN DENGUE.CSV DATASET**

1. SHOW ALL NAS IN "NUMBER" COLUMN

```
#-----------------------------------------------------
#6. Handling Missing values in Dengue.csv Dataset

#6a. Show All NAs in "number" column
f = dengue %>%
  filter(is.na(number))
```



2. ANOTHER WAY OF SHOWING ALL NAS IN ALL COLUMNS

```
#6b. Another way of showing all NAs in all columns
g = dengue %>%
  filter(!complete.cases(.))
```

3. SHOWING ALL NO NAS (FILLED COLUMNS) NOW

```
#6c. Showing All NO NAs (filled columns) now....
h = dengue %>%
    filter(complete.cases(.))
```

```
#----------------------------------------------------------------
#7. Mutate data

#using "eweek" to create a new column called "date"...
i = dengue %>%
  mutate(date = ymd(paste0(year,"-01-01"))+weeks(eweek))
```

## H. FILTER, MUTATE THEN PLOT

### 1. DENGUE.CSV

```
#-------------------------------------------------------------------
#8. Filter, Mutate then Plot

#8a.  Dengue.csv
dengue %>%
  filter(complete.cases(.)) %>%          #remove all NAs

  filter(type_dengue=='Dengue') %>%      #filter out only "Dengue" type

  mutate(date = ymd(paste0(year,"-01-01"))+weeks(eweek)) %>%
  #create a new column called "date"

  select(date,number) %>%
  #selecting out only "date" and "number" columns to plot

  plot()
```

2. ANOTHER EXAMPLE FOR FILTER, MUTATE THEN PLOT
(VACCINATION.XLS)

```r
#---------------------------------------------------------------------
#8b. Another Example for Filter, Mutate then Plot (vaccination.xls)
#https://www.alvinang.sg/s/vaccination.xlsx

vaccination <- read_excel("vaccination.xlsx")

vaccination %>%
  filter(complete.cases(.)) %>%                    #remove all NAs

  filter(vaccination_type=='Poliomyelitis') %>%
  #filter out only 'Poliomyelitis'

  mutate(date = ymd(paste0(year,"-01-01"))) %>%
  #create a new column called "date"

  mutate(doses = no_of_doses_in_thousands) %>%
  #rename the column

  select(date,doses) %>%
  #selecting out only "date" and "doses" columns to plot

  plot()
```

3. FILTER, MUTATE THEN EXPORT TO CSV

```r
# Export to CSV
dengue_filtered <- dengue %>%
  filter(complete.cases(.)) %>%
  filter(type_dengue=='Dengue') %>%
  mutate(date = ymd(paste0(year,"-01-01"))+weeks(eweek)) %>%
  select(date,number)

write_csv(dengue_filtered, path = "dengue_filtered.csv")
```
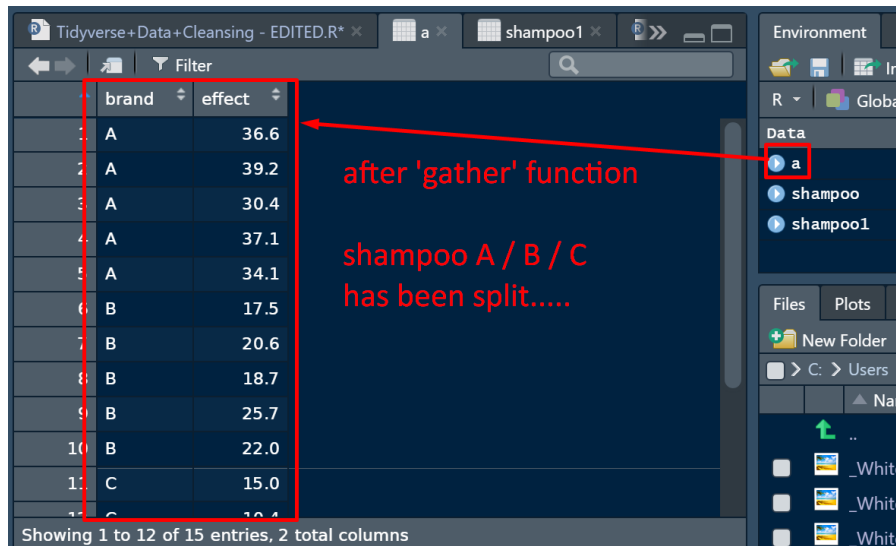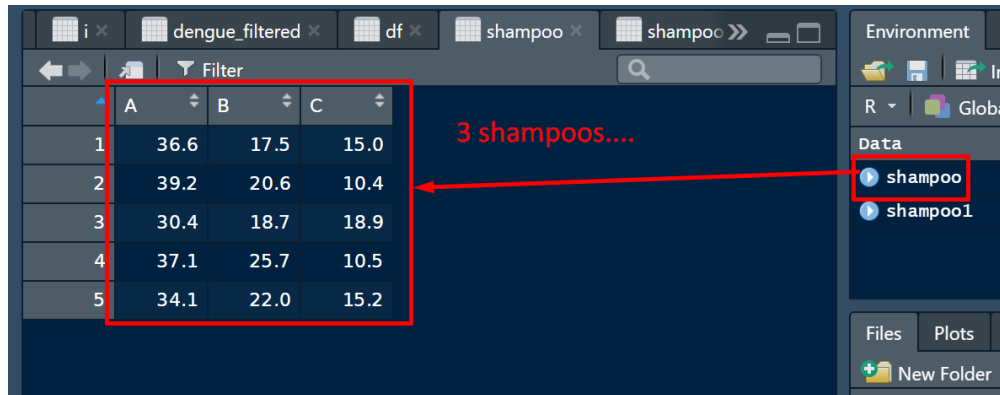


this has now been created (write csv) into the working directory folder.....

4. USING GATHER TO PIVOT DATA

```r
shampoo =
  data.frame('A'=c(36.6,39.2,30.4,37.1,34.1),
             'B' = c(17.5,20.6,18.7,25.7,22.0),
             'C'=c(15.0,10.4,18.9,10.5,15.2))

shampoo1 <- as_tibble(shampoo)

a = shampoo1 %>%
  gather(brand, effect)
```

3 shampoos....

after 'gather' function

shampoo A / B / C
has been split.....

Showing 1 to 12 of 15 entries, 2 total columns

```
#----------------------------------------------------------------
#11. Data Joins

df1 = data_frame(name=c('Ally','Steve','John'),age=c(45,46,47))
df2 = data_frame(name=c('Ally','Belinda','John'),age=c(45,48,47))
```

1. LEFT JOIN



```
#11a. Left Join
left_join(df1,df2,by='name')
```



df1

| | name | age |
|---|---|---|
| 1 | Ally | 45 |
| 2 | Steve | 46 |
| 3 | John | 47 |

df2

| | name | age |
|---|---|---|
| 1 | Ally | 45 |
| 2 | Belinda | 48 |
| 3 | John | 47 |

LEFT JOIN

| | name | age.x | age.y |
|---|---|---|---|
| 1 | Ally | 45 | 45 |
| 2 | Steve | 46 | NA |
| 3 | John | 47 | 47 |

2. RIGHT JOIN

```
#11b. Right Join
right_join(df1,df2,by='name')
```

| | name | age |
|---|---|---|
| 1 | Ally | 45 |
| 2 | Steve | 46 |
| 3 | John | 47 |

df1

| | name | age |
|---|---|---|
| 1 | Ally | 45 |
| 2 | Belinda | 48 |
| 3 | John | 47 |

df2

| | name | age.x | age.y |
|---|---|---|---|
| 1 | Ally | 45 | 45 |
| 2 | John | 47 | 47 |
| 3 | Belinda | NA | 48 |

RIGHT JOIN

3. INNER JOIN

```
#11c. Inner Join
inner_join(df1,df2,by='name')
```



| | name | age |
|---|---|---|
| 1 | Ally | 45 |
| 2 | Steve | 46 |
| 3 | John | 47 |

df1

| | name | age.x | age.y |
|---|---|---|---|
| 1 | Ally | 45 | 45 |
| 2 | John | 47 | 47 |

INNER JOIN

| | name | age |
|---|---|---|
| 1 | Ally | 45 |
| 2 | Belinda | 48 |
| 3 | John | 47 |

df2

4. FULL JOIN

```
#11d. Full Join
full_join(df1,df2,by='name')
```

| | name | age |
|---|---|---|
| 1 | Ally | 45 |
| 2 | Steve | 46 |
| 3 | John | 47 |

df1

| | name | age |
|---|---|---|
| 1 | Ally | 45 |
| 2 | Belinda | 48 |
| 3 | John | 47 |

df2

| | name | age.x | age.y |
|---|---|---|---|
| 1 | Ally | 45 | 45 |
| 2 | Steve | 46 | NA |
| 3 | John | 47 | 47 |
| 4 | Belinda | NA | 48 |

FULL JOIN

## J.    GROUPBY



the
sleep
dataset

```
s1 = sleep %>%
  group_by(group) %>%
  summarize(avg_extra=mean(extra))
```

```
data("starwars", package = "dplyr")
d = starwars
```

| | name | height | mass | hair_color | skin_color |
|---|---|---|---|---|---|
| 1 | Luke Skywalker | 172 | 77.0 | blond | fair |
| 2 | C-3PO | 167 | 75.0 | NA | gold |
| 3 | R2-D2 | 96 | 32.0 | NA | white, blue |
| 4 | Darth Vader | 202 | 136.0 | none | white |
| 5 | Leia Organa | 150 | 49.0 | brown | light |
| 6 | Owen Lars | 178 | 120.0 | brown, grey | light |
| 7 | Beru Whitesun lars | 165 | 75.0 | brown | light |
| 8 | R5-D4 | 97 | 32.0 | NA | white, red |
| 9 | Biggs Darklighter | 183 | 84.0 | black | light |
| 10 | Obi-Wan Kenobi | 182 | 77.0 | auburn, white | fair |

1. REMOVING THE 'HEIGHT' COLUMN

```
#remove the 'height' column
a = select(starwars, -height)
```

height column is now removed

| | name | mass | hair_color | skin_color | eye_color |
|---|---|---|---|---|---|
| 1 | Luke Skywalker | 77.0 | blond | fair | blue |
| 2 | C-3PO | 75.0 | NA | gold | yellow |
| 3 | R2-D2 | 32.0 | NA | white, blue | red |
| 4 | Darth Vader | 136.0 | none | white | yellow |
| 5 | Leia Organa | 49.0 | brown | light | brown |
| 6 | Owen Lars | 120.0 | brown, grey | light | blue |
| 7 | Beru Whitesun lars | 75.0 | brown | light | blue |
| 8 | R5-D4 | 32.0 | NA | white, red | red |
| 9 | Biggs Darklighter | 84.0 | black | light | brown |
| 10 | Obi-Wan Kenobi | 77.0 | auburn, white | fair | blue-g |

Environment

Data
a
d
starwars

Files    Plots
New Folder
C: > Users

2. RENAME THE 'NAME' COLUMN

```
#rename the 'name' column
b = starwars %>%
    rename(BLABLABLA = name)
```

name column has been renamed

| | BLABLABLA | height | mass | hair_color | skin_color |
|---|---|---|---|---|---|
| 1 | Luke Skywalker | 172 | 77.0 | blond | fair |
| 2 | C-3PO | 167 | 75.0 | NA | gold |
| 3 | R2-D2 | 96 | 32.0 | NA | white, blue |
| 4 | Darth Vader | 202 | 136.0 | none | white |
| 5 | Leia Organa | 150 | 49.0 | brown | light white |
| 6 | Owen Lars | 178 | 120.0 | brown, grey | light |
| 7 | Beru Whitesun lars | 165 | 75.0 | brown | light |
| 8 | R5-D4 | 97 | 32.0 | NA | white, red |
| 9 | Biggs Darklighter | 183 | 84.0 | black | light |
| 10 | Obi-Wan Kenobi | 182 | 77.0 | auburn, white | fair |

Environment

Data
- a
- b
- d
- starwars

Files   Plots

## L. DIFFERENCES BETWEEN TIBBLE VS DATAFRAME

1. TIBBLE

```
#-------------------------------------------------
#15. Differences between Tibble vs Dataframe
#15a. Tibble

df <- tibble(
  'male' = c(2.3,3.5,4.6,3.2,2.5),
  'female' = c(1.3,2.6,1.7,1.9,2.1)
)
df
```

2. DATAFRAME

```
#15b. Dataframe
shampoo = data.frame(
   'A'=c(36.6,39.2,30.4,37.1,34.1),
   'B' = c(17.5,20.6,18.7,25.7,22.0),
   'C'=c(15.0,10.4,18.9,10.5,15.2))
```

3. AS TIBBLE

```
#15c. As Tibble
shampoo1 <- as_tibble(shampoo)
```



even after tibbling,
you don't see any
difference....

# Data Frame

# Tibble

```
df1 <- data.frame(
  gender = c("Female", "Female","Male"),
  height = c(152, 171.5, 165),
  weight = c(81,93, 78),
  age =c(42,38,26),
  row.names=c('Ally','Belinda','Alvin')
)
```

```
df2 <- tibble(
  gender = c("Female", "Female","Male"),
  height = c(152, 171.5, 165),
  weight = c(81,93, 78),
  age =c(42,38,26),
  row.names=c('Ally','Belinda','Alvin')
)
```

## Not Much Difference....

| | gender | height | weight | age |
|---|---|---|---|---|
| Ally | Female | 152.0 | 81 | 42 |
| Belinda | Female | 171.5 | 93 | 38 |
| Alvin | Male | 165.0 | 78 | 26 |

| | gender | height | weight | age | row.names |
|---|---|---|---|---|---|
| 1 | Female | 152.0 | 81 | 42 | Ally |
| 2 | Female | 171.5 | 93 | 38 | Belinda |
| 3 | Male | 165.0 | 78 | 26 | Alvin |

- There's not much visible difference between a Data Frame vs Tibble…..

- Except that Tibble adds an extra column….

5. COMPARING STRUCTURE (STR)

# Data Frame

# Tibble



```
> str(df1)
'data.frame':    3 obs. of  4 variables:
 $ gender: chr   "Female" "Female" "Male"
 $ height: num   152 172 165
 $ weight: num   81 93 78
 $ age   : num   42 38 26
```

```
> str(df2)
tibble [3 x 5] (S3: tbl_df/tbl/data.frame)
 $ gender    : chr [1:3] "Female" "Female" "Male"
 $ height    : num [1:3] 152 172 165
 $ weight    : num [1:3] 81 93 78
 $ age       : num [1:3] 42 38 26
 $ row.names : chr [1:3] "Ally" "Belinda" "Alvin"
>
```

The Structure of a Data Frame vs Tibble also don't show much difference....

- Even if you look at the strucure…they display the same things….

# Data Frame

# Tibble

```
> df1$ge
[1] "Female" "Female" "Male"
```

Even though the proper column name
Is called "gender"..... If you use a DataFrame,
You can misspell it as $ge and it will still
Show the column items....

This might cause future errors if you accidentally
Call out the wrong column with similar column
"ge" headings.....

```
> df2$ge
NULL
Warning message:
Unknown or uninitialised column: `ge`.
> df2$gender
[1] "Female" "Female" "Male"
```

However, for Tibble, you are not able to
Display the column items if you misspell
The column name.....it will show an error...

You have to type out the whole "$gender"
To get the items.....

This prevents future errors.....

- But you are not able to use short forms for the column names to call out the items…

7.  COMPARING DISPLAY



- The Penguins dataset can be found here: https://www.alvinang.sg/s/penguins.csv

- For dataframe as shown above, we see that it displays very messily in the console….



- For tibble, we now see that its displayed neatly….in the console….

**ABOUT DR. ALVIN ANG**

Dr. Alvin Ang earned his Ph.D., Masters and Bachelor degrees from NTU, Singapore. He was a previously a Professor, Scientist and Financial Consultant. Currently, he owns multiple self-started businesses and is a Personal/Business Advisor.

More about him at www.AlvinAng.sg