

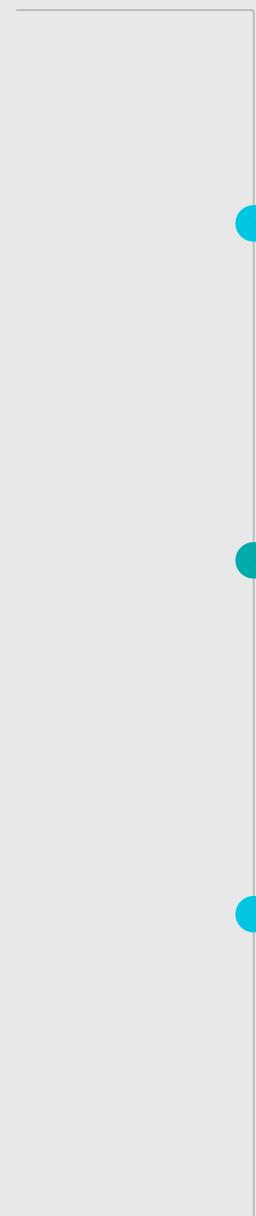
Business Analytics

Session 2

Ing. Juan José Franklin Uraga, PhD.



Session 1



Data Wrangling

Class evidence # 2

Homework review

Statistics Review



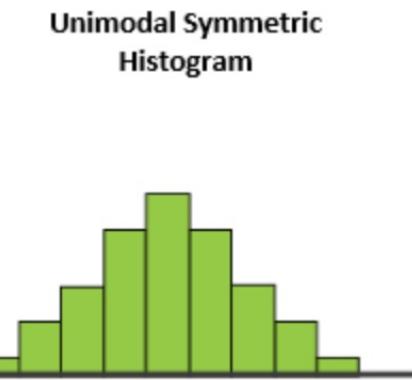
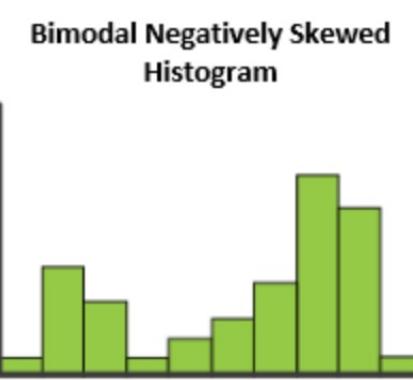
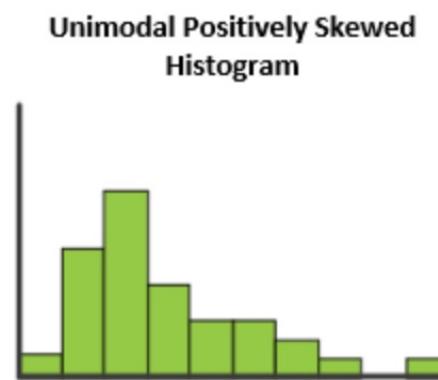
Statistics

- It is the art and science of collecting data, analyzing it, presenting it and interpreting it.



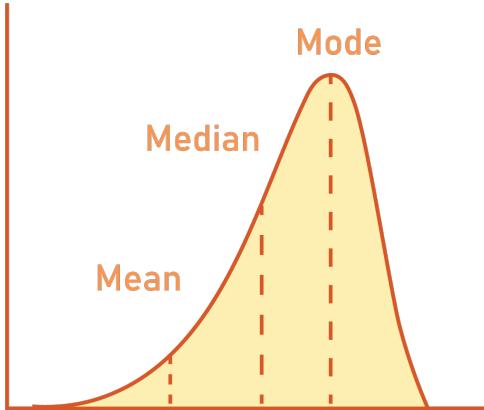
Histogram

- It is used to represent the frequencies of the observed values of a quantitative variable.
- It constitutes a graphical way of visualizing the behavior and trends of the variable, detecting atypical data and forms of a frequency distribution.

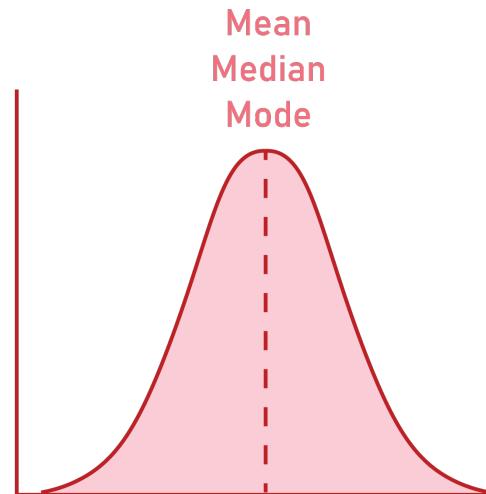


Shapes of numerical distributions

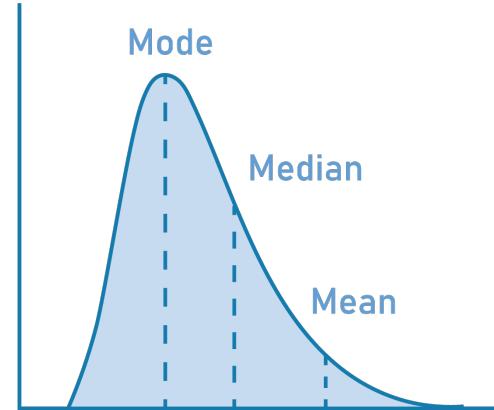
- It describes how data are distributed.



Left skew



Normal distribution
Symmetric



Right skew

Central tendency measures

MEAN OR AVERAGE

the sum of the numbers divided
by the amount of numbers

$$5 + 5 + 5 + 6 + 7 + 7 + 14 \\ 49 / 7 = 7$$

MEDIAN

the number in the middle

5 5 5 6 7 7 14

(numbers must be in ascending order)

MODE

the number that appears
the most

5 5 5 6 7 7 14

Mean

- Arithmetic mean or average.

$$\bar{X} = \frac{\sum_{i=1}^n X_i}{n} = \frac{X_1 + X_2 + \dots + X_n}{n}$$

▪ Each one of the observed values

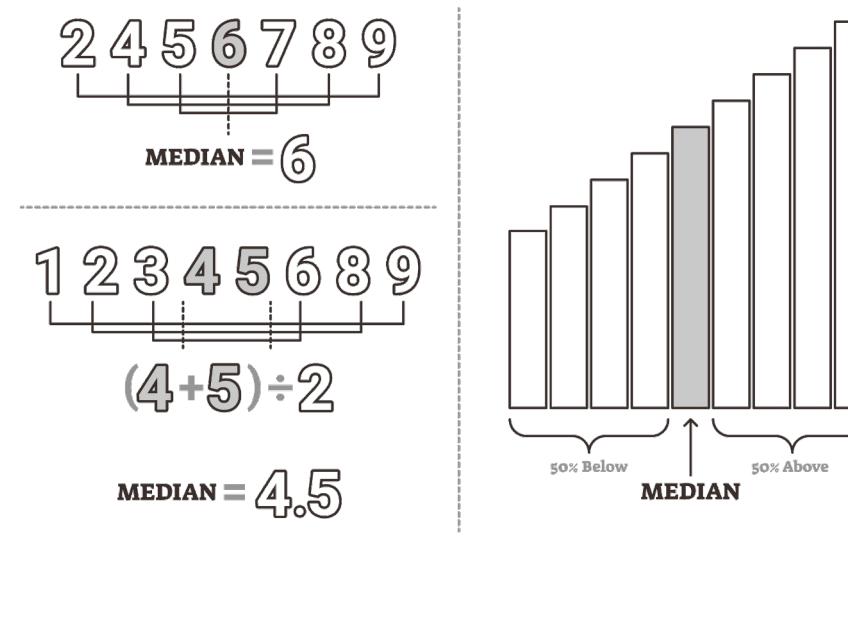
▪ sample size



=PROMEDIO(dataset)

Median

- Finding the median value: position $(n+1)/2$ of the **sorted dataset**.
- If the number of values is odd, the median is the middle number.
- If the number of values is even, the median is the average of the two middle numbers.



=MEDIANA(dataset)

Mean vs Median

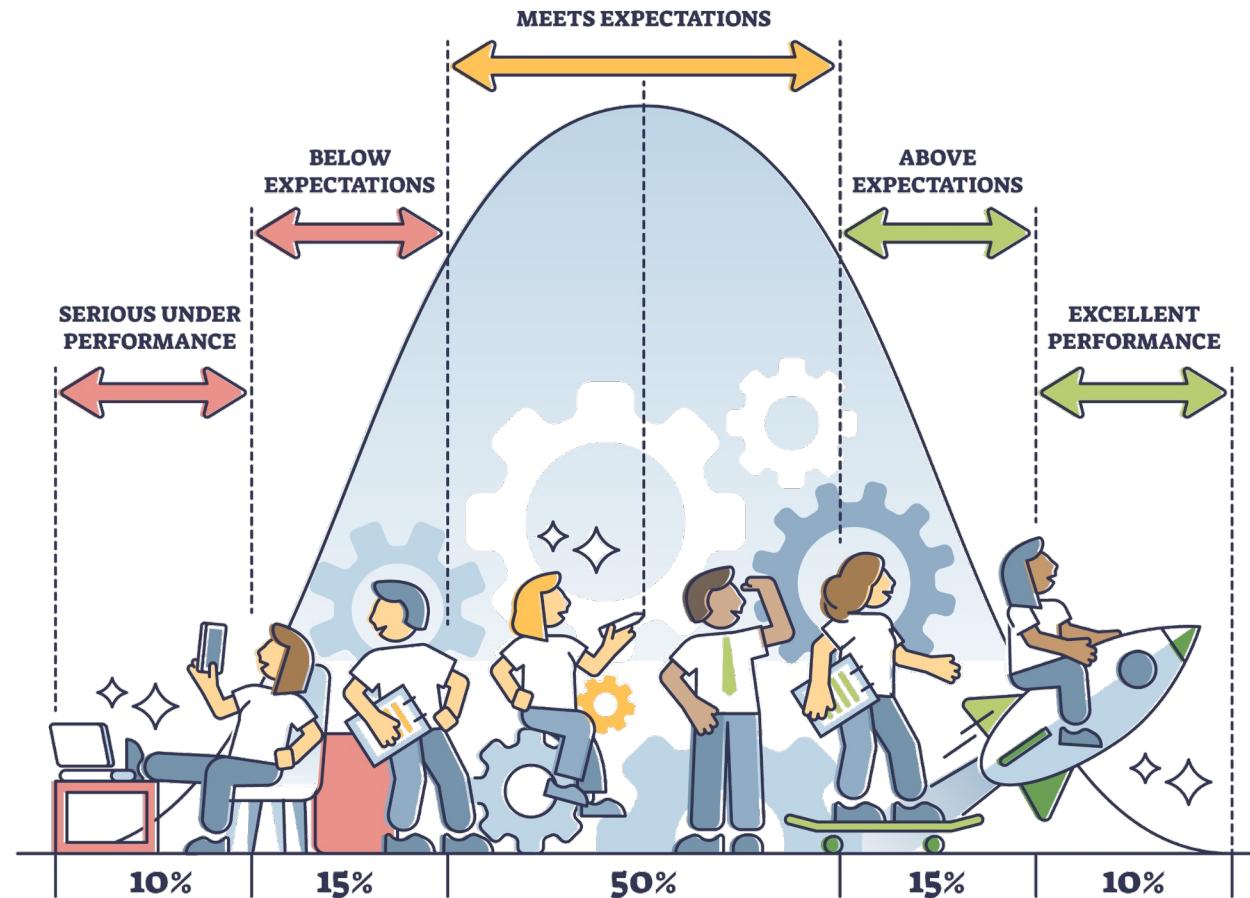
- Month income (USD) is analyzed for a group of friends:
 - \$1000, \$1250, \$1300, \$1500, \$1600, \$2000, \$70,000
 - Mean: \$11,236
 - Median: \$1,500
- Which one is **better** to represent this dataset?
- The **mean is heavily influenced** by extreme values.
- **Median** will be a better option when dataset is **skewed** or under the presence of **outliers**.

Probability distributions



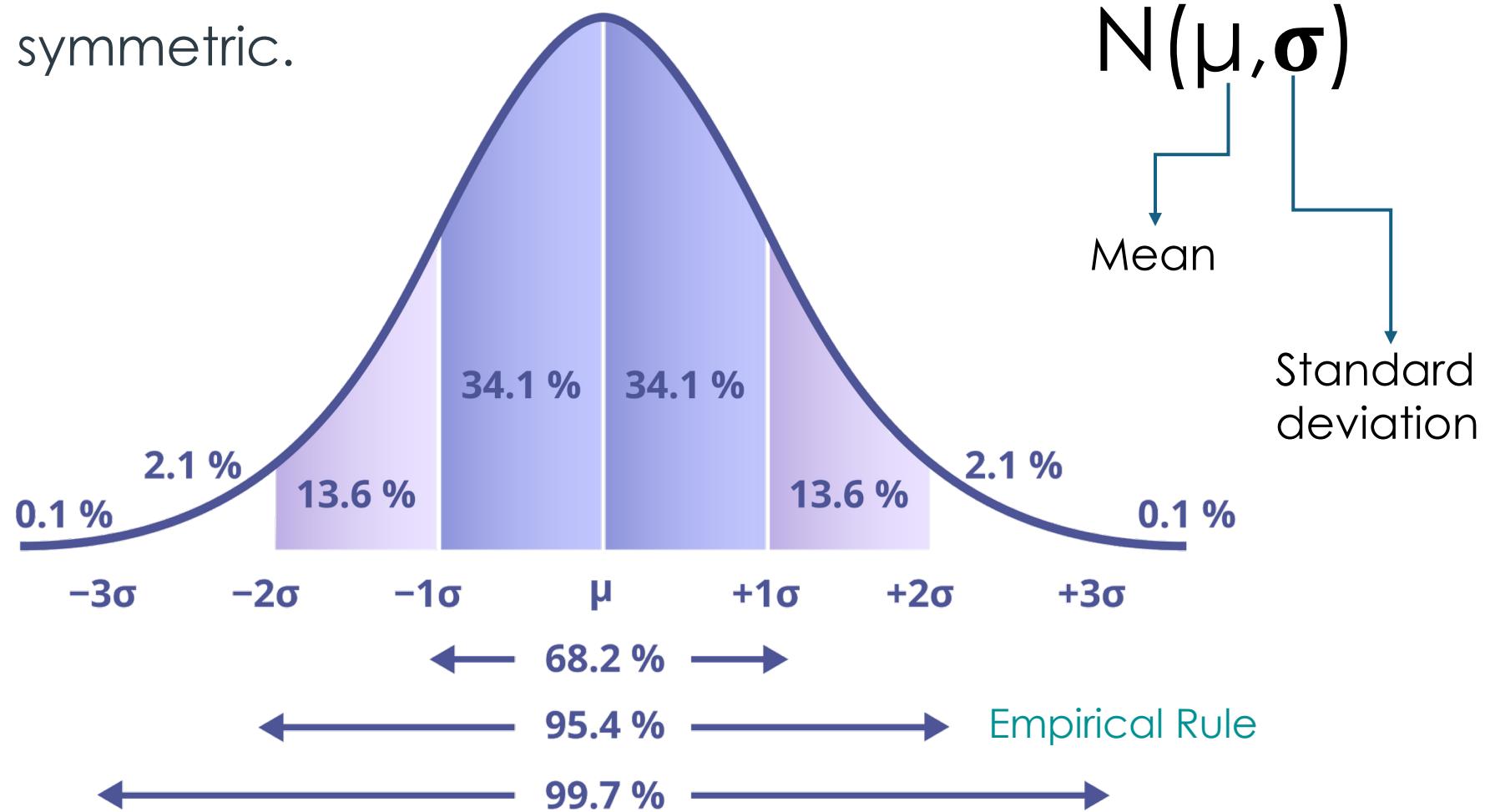
Normal distribution

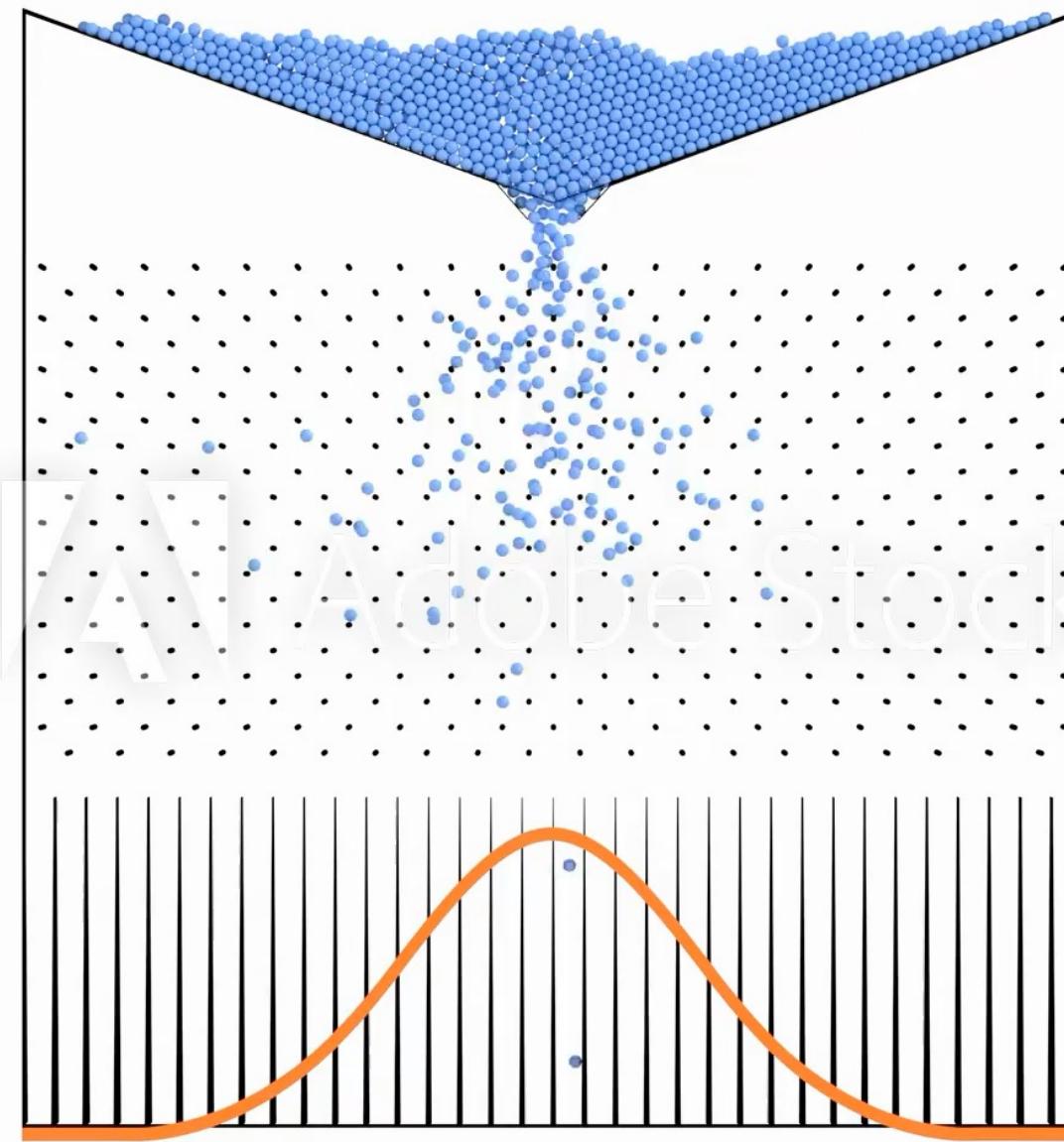
- Many variables in nature are normally distributed.



Normal distribution

- Unimodal and symmetric.
- Bell-shaped.





Chebysev Rule

- Regardless of how the data are distributed, at least $(1 - 1/k^2) \times 100\%$ of the values will fall within k standard deviations of the mean (for $k > 1$)
 - Examples:

At least

within

$$\frac{(1 - 1/1^2) \times 100\% = 0\%}{(1 - 1/2^2) \times 100\% = 75\%} \dots \dots \dots k=1 (\mu \pm 1\sigma)$$

$$(1 - 1/2^2) \times 100\% = 75\% \dots \dots \dots k=2 (\mu \pm 2\sigma)$$

$$(1 - 1/3^2) \times 100\% = 89\% \dots \dots \dots k=3 (\mu \pm 3\sigma)$$

Business Analytics

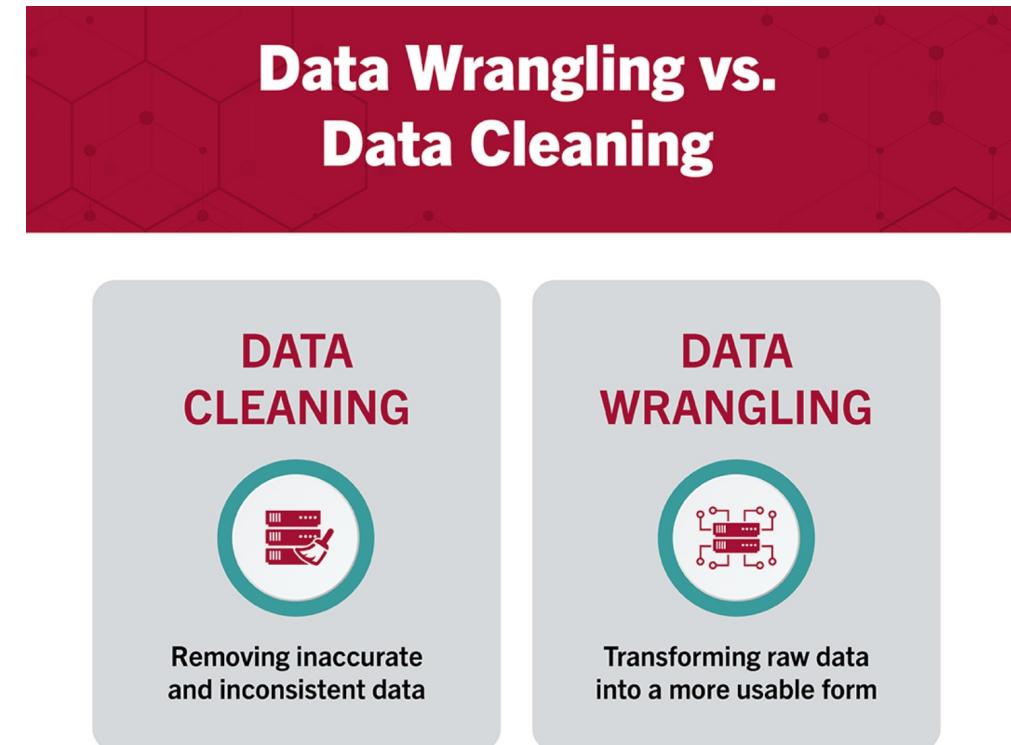
Data wrangling With R Studio

Ing. Juan José Franklin Uraga, PhD.



Data Wrangling

- Data wrangling—also called data cleaning, data remediation, or data munging—refers to a variety of processes designed to transform raw data into more readily used formats.



Objectives for today



- Pipe operator %>% (and then...)
- Filter
- Select
- Group_by
- Count
- Arrange
- Summarise
- Mutate
- Ggplot (histograms with facets)

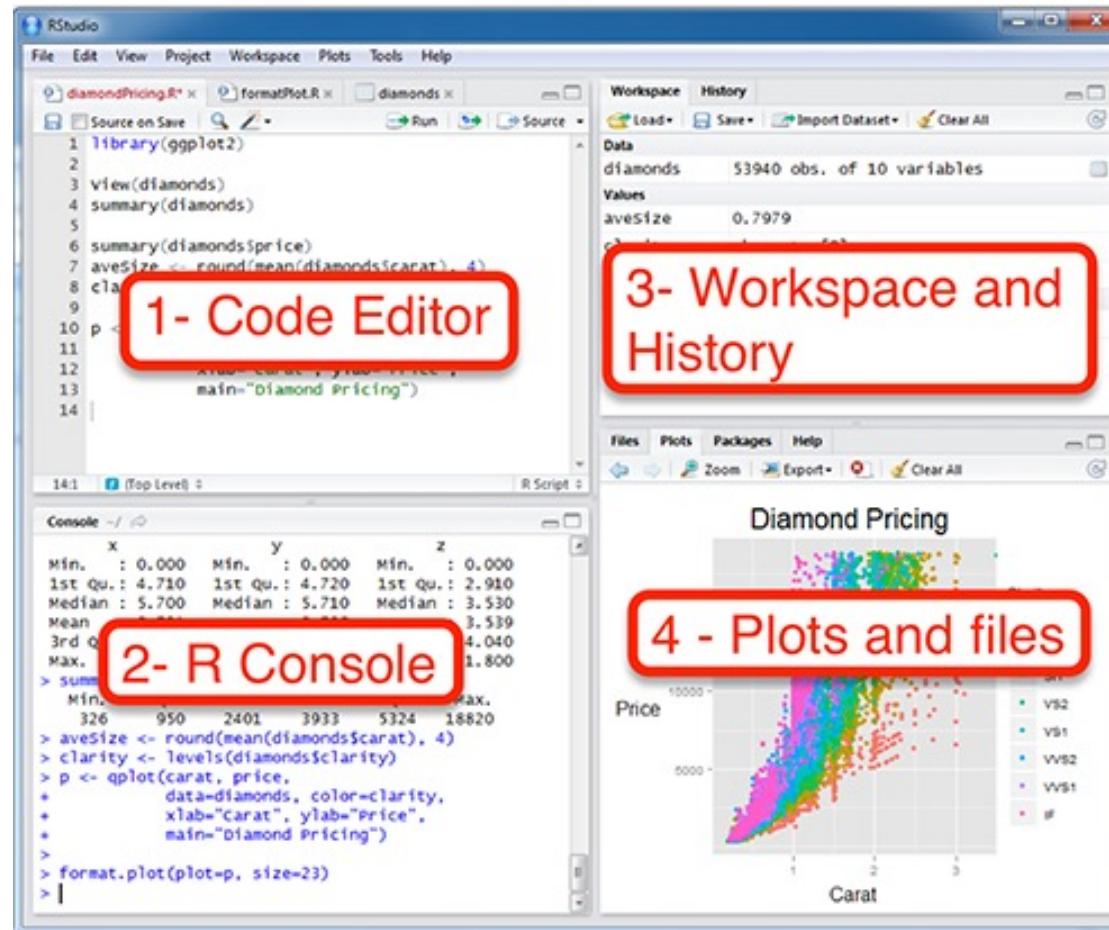
Objectives for today



Start working with R Markdown

- R Markdown is a special type of document that mixes together regular text with R code.
- In an R Markdown document, you can write explanations, make lists, or add images just like in any other document.
- But what's special is that you can also include pieces of R code directly in the document.
- When you "knit" or compile an R Markdown document, it runs the R code and inserts the results back into the document.

The R studio 4 windows



Filter

```
```{r}
library(tidyverse)
starwars %>%
 filter(homeworld=="Tatooine") %>%
 filter(height>170)
```



A tibble: 6 × 14

name	height	mass	hair_color
<chr>	<int>	<dbl>	<chr>
Luke Skywalker	172	77	blond
Darth Vader	202	136	none
Owen Lars	178	120	brown, grey
Biggs Darklighter	183	84	black
Anakin Skywalker	188	84	blond
Clegg Lars	183	NA	brown



# Select

```
```{r}
starwars %>%
  select(height, mass, homeworld, species, ends_with("color"))
```

```

A tibble: 87 × 7

| height | mass  | homeworld | species | hair_color |   |
|--------|-------|-----------|---------|------------|---|
| <int>  | <dbl> | <chr>     | <chr>   | <chr>      | ▶ |
| 172    | 77.0  | Tatooine  | Human   | blond      | ▶ |
| 167    | 75.0  | Tatooine  | Droid   | NA         | ▶ |
| 96     | 32.0  | Naboo     | Droid   | NA         | ▶ |

```
```{r}
starwars1 <- starwars %>%
  select(height, mass, homeworld, species, ends_with("color"))
```

```

<- Creates an object with  
the selection.

# Group\_by

```
```{r}
starwars %>%
  group_by(species) %>%
  count() %>%
  arrange(desc(n))
```
```

A tibble: 38 × 2 Groups: species [38]

| species  | n     |
|----------|-------|
| <chr>    | <int> |
| Human    | 35    |
| Droid    | 6     |
| NA       | 4     |
| Gungan   | 3     |
| Kaminoan | 2     |
| Mirialan | 2     |
| Twi'lek  | 2     |

# Dealing with NA values

```
is.na(starwars$species)
```

```
[1] FALSE
[13] FALSE FALSE FALSE FALSE FALSE TRUE FALSE FALSE FALSE FALSE FALSE FALSE
[25] FALSE
[37] FALSE
[49] FALSE TRUE TRUE
[61] FALSE
[73] FALSE FALSE FALSE FALSE FALSE FALSE FALSE TRUE FALSE FALSE FALSE
[85] FALSE FALSE FALSE
```

```
starwars[is.na(starwars$species),]
```

```
A tibble: 4 × 14
name height mass hair_color skin_color eye_color birth_year sex gender
<chr> <int> <dbl> <chr> <chr> <chr> <dbl> <chr> <chr>
1 Jek Tono... 180 110 brown fair blue NA <NA> <NA>
2 Gregar T... 185 85 black dark brown NA <NA> <NA>
3 Cordé 157 NA brown light brown NA <NA> <NA>
4 Sly Moore 178 48 none pale white NA <NA> <NA>
i 5 more variables: homeworld <chr>, species <chr>, films <list>,
vehicles <list>, starships <list>
```

# Group\_by, not considering NA

```
```{r}
starwars %>%
  drop_na(species) %>%
  group_by(species) %>%
  count() %>%
  arrange(desc(n))
```

```



A tibble: 37 × 2

Groups: species [37]

**species**  
<chr>

**n**  
<int>

Human

35

Droid

6

Gungan

3

Kaminoan

2

Mirialan

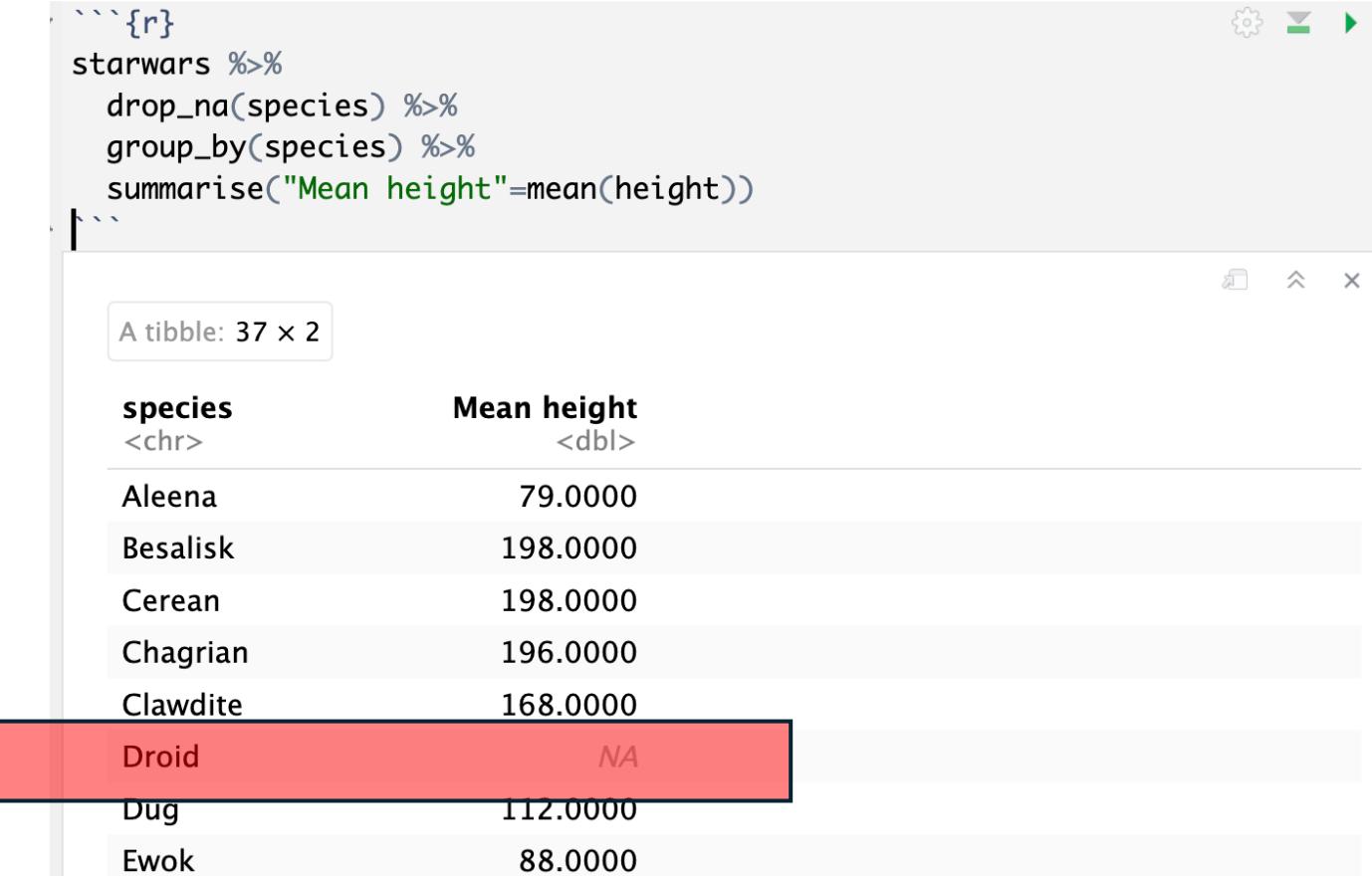
2

Twi'lek

2

# Group\_by & summarise

**Task:** find the mean height by species



```
```{r}
starwars %>%
  drop_na(species) %>%
  group_by(species) %>%
  summarise("Mean height"=mean(height))
```
A tibble: 37 × 2
 species Mean height
 <chr> <dbl>
1 Aleena 79.0000
2 Besalisk 198.0000
3 Cerean 198.0000
4 Chagrian 196.0000
5 Clawdite 168.0000
6 Droid NA
7 Dug 112.0000
8 Ewok 88.0000
```

# Group\_by & summarise & arrange

**Task:** find the mean height by species

```
```{r}
starwars %>%
  drop_na(species) %>%
  group_by(species) %>%
  drop_na(height) %>%
  summarise(Mean_height=mean(height)) %>%
  arrange(desc(Mean_height))
````
```

A tibble: 37 × 2

| species  | Mean_height |
|----------|-------------|
| <chr>    | <dbl>       |
| Quermian | 264.0000    |
| Wookiee  | 231.0000    |
| Kaminoan | 221.0000    |
| Kaleesh  | 216.0000    |
| Gungan   | 208.6667    |
| Pau'an   | 206.0000    |

Alternative

```
```{r}
starwars %>%
  drop_na(species) %>%
  group_by(species) %>%
  summarise(Mean_height=mean(height,na.rm=TRUE)) %>%
  arrange(desc(Mean_height))
````
```

# Mutate & case\_when

**Task:** create a new variable in starwars data set. The name is "size". Get a summary of mass, first.

If mass < 40, classify it as small, if <100 it should be regular, more than 100: big.

Notice the  
customized order

```
```{r}
starwars %>%
  drop_na(mass) %>%
  mutate(Size=case_when(mass<40~"Small",
                        mass<100~"Regular",
                        mass>100~"Big")) %>%
  select(name,species,Size) %>%
  arrange(match(Size,c("Small","Regular","Big")))
```
A tibble: 59 × 3
 name species Size
 <chr> <chr> <chr>
1 R2-D2 Droid Small
2 R5-D4 Droid Small
3 Yoda Yoda's species Small
4 Wicket Systri Warrick Ewok Small
5 Ratts Tyerel Aleena Small
6 Luke Skywalker Human Regular
```

# Task



1. Create a new df of starwars with only human and droids.
2. Classify them as small or tall, depending on their height (<180 small,  $\geq$  180 tall).
3. Create a contingency table that shows how many humans and droids are short and tall. In absolute frequency and relative frequency.
4. Create a mosaic plot to present the information.

# Task



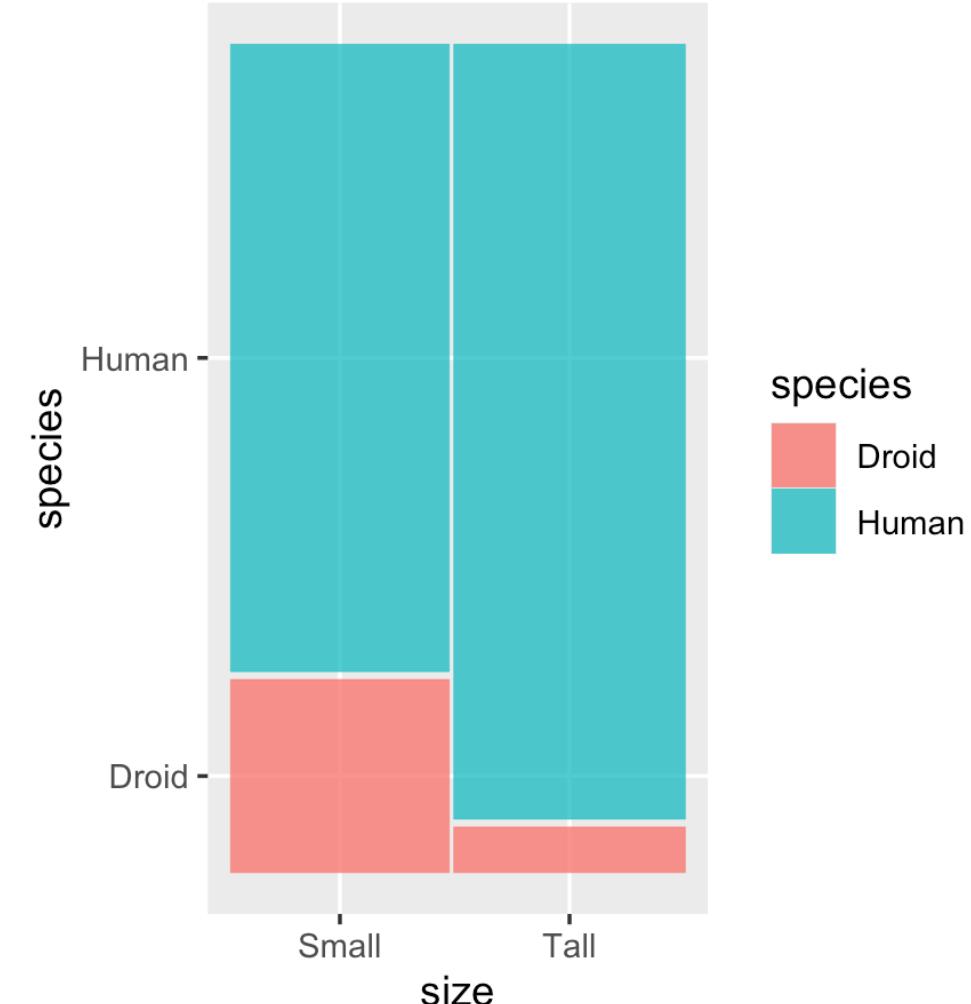
```
starwars3 <- starwars %>%
 filter(species=="Human" | species=="Droid") %>%
 mutate(size=case_when(height<180~"Small",
 height>=180~"Tall"))

table(starwars3$species,starwars3$size) %>%
 prop.table() %>%
 addmargins()

starwars3 <- starwars3 %>% drop_na(size)

ggplot(starwars3) + geom_mosaic(aes(x = product(species, size), fill = species))
```

**Of the total population:**  
14% are droids  
86% are humans  
    11% are small droids  
    3% are tall droids  
    37% are small humans  
    49% are tall humans.



# Statistics

collection data experiments theory

forecasting particularly increase mean disciplines total particular sampling particular sampling probability tools quantity prediction related applicable business access median a.k.a. without advance leads a.k.a. predictive describe tested surveys population science statisticians subject explanation basis predictions used holds referring roughly whose distinct singular pertaining discipline true results government direction deals number mathematical branch together may design uncertainty communicating calculated plural rather quality fields

# organization

experience application deduce part considered problems social working useful ways improve process soundness collecting since samples larger roots given grouped inductive observations provide one confused survey wrong empirical summarize successful logically companies usually wide vital comprising word applications art modeled expertise sciences aspects people gained models direct inferences variety studied interpreting randomness guiding natural based thinking planning starts consider versed someone help

# descriptive

scientific analysis

# interpretation

necessary using empirical summarize successful logically companies usually wide vital comprising word applications art modeled expertise sciences aspects people gained models direct inferences variety studied interpreting randomness guiding natural based thinking planning starts consider versed someone help

# statistical Inference

also others organizations applied difference variety studied interpreting randomness guiding natural based thinking planning starts consider versed someone help

# mathematics

moves concerned presentation parameters moves focus probabilities new misleading academic way opposite