# Tidy data :-

A way of standardising the organization of data within R.

Data Frame → Collection of columns.

> Columns should be named
> Data stored can be many different types, like numeric, factor, or character.
> Each column should contain same number of data items.

* Tibbles : Streamlined data frames

> Never change data types of the inputs
> Never change names of your variables
> Never create row names
> Make printing easier

# pulls only first 10 rows

* Tidy data standars →
> Variables are organized into columns
> Observations are organized into rows
> Each value must have its own cell

————— ✗ ————— ✗ ————— ✗ —————

install.packeages ("tidyverse")          glimpse() & str()
library (ggplot2)
data ("diamonds)                         Summary
View (diamonds)

library(tidyverse)  head (diamonds)      # first 6 rows.
                    str (diamonds)    # High level info cd_name & type
                    colnames (diamonds)   # only names
from ]
dplyr ]   →  mutate (diamonds, carat_2 = carat_2 * 100)
                          new column created.

Data Frame= Data Analysis default way of
interacting with data.

Cleaning up with the basics :→

3 Packages    Here , Skimr , Janitor

Here :→  makes  Referencing  easier.
Skimr :→ Summaring data is easy & skim quickly.
Janitor &→  functions for cleaning data.

pacman :: p_load (pacman, here, skimr, janitor, dplyr)
↓
Saved   10 lines

Palmer Penguin Package :  3 species , size, clutch size,
blood isotope ratio.

Install packages ("palmerpenguins")
library (palmerpenguins)

+ skim_without 'charts ()        glimpse ()
* head ()                        select ()

> skim_without_charts (penguins)     # Gives summary
> glimpse (penguins)
> head (penguins)
>   penguins    % > %
       select (species)        # only species column

>> penguins % > %.              # Everything but species.
       select (- species)

**\* [ Rename Column ]**

penguins %>%
     rename (island_new = island)

**\*** rename_with (penguins, tolower)
     ↓
     for consistency.

**\*** Clean Names : Consistents, no duplicats, etc.

clean_names (penguins)

---

Operators ⟶

| | | | | |
|---|---|---|---|---|
| + : | Add | % % | Modeulus |
| - : | Sub | % / % | Integer Div |
| \* : | Mul | ^ | Exponent |
| / : | Div | | |

**Arithmetic**

---

**Relational**

| | | | | |
|---|---|---|---|---|
| > | → | greater than | == | equal to |
| < | → | less then | != | not equal to |
| >= | → | greater than equal | | |
| <= | → | less than equal | | |

---

**Logical**

| | | | |
|---|---|---|---|
| & | → | Element-wise logical AND | $x \leftarrow c(3,5,7)$ |
| && | → | logical AND | $y \leftarrow c(2,4,6)$ |
| \| | → | Element-wise logical OR | $x < 5$ & $y < 5$ |
| \|\| | → | logical OR | ④ T F F → |
| ! | → | logical Not | $x < 5$ && $y < 5$ |
| | | | TRUE |

Assignment
Operator

| ← | leftward |
| <-- | leftward } Assignment |
| = | leftward |

— ] Rightward
--->> } Assignment

Page No
Date

\* Organizing data → arrange()
group_by()
filter()

sort by column

{ penguins %>% arrange (bill_length_mm)  # Ascending

penguins %>% arrange (- bill_length_mm)  # Descending

→ These results are only in the console to save create a
new dataframe

penguins2 ← penguins %>% arrange (bill_length_mm)
View (penguins 2)

— ✗ —

Group By

penguins %>% group by (island) %>% drop_na()
%>% summarize (mean_bill_length_mm =
mean (bill_length_mm))

penguins %>% group by (species, island) %>%
dropna() %>% summarize (max_bl = max (bill_length_mm),
mean_b = mm (bill_length_mm))

— ✗ —

Filter

penguins %>% filter (species == "Adelie")

4/3

(P) R let's you more back and forth betn analysis & visualization quickly.

> Plotly → wide Range; General Purpose

> RGL :→ Specific, 3D

> Lattice

> Diagraphs

> leaflet                    > gganimate

> Highchartes              > gg ridges

> Patchwork               > ggplot2

• ggplot2 is most popular visualization package in R.
• Use it singulars or with others
• less code more viz

Creater: statistician, developer : Hadley Wickham 2005

Inspiration: The Grammar of Graphics
           Schorarly study of Data Viz by leland Wilkinson

In the same way English grammar gives us rules to build any kind of sentence

It gives us rules to build any kind of visual.

Rule of Grplot 2:

> Start with the ggplot function & choose a dataset to work with
> Add a geom-function to display your data.
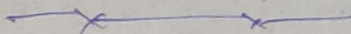> Map the variables you want to plot in the arguments of the aes() function.

Template →

Ex
```
ggplot (data = penguins) +
    geom-point (mapping = aes (x = flipper_length_mm,
                                y = body_mass_g))
```
↓

```
ggplot (data = <Data>) +
    <GEOM_FUNC> (mapping = aes (<AES mappings>))
```

Ex
```
ggplot (data = penguins) +
    geom-point (mapping = aes (x = bill_length_mm,
                                y = bill_depth_mm))
```

✗————✗

(✱) aes (x = "—||—", y = "—||—", color = species)

changes color of the different data points.

(✱) shape = species ; assigns different shape to each data point (clustering)

legend is generated automatically.

(✱) color = species & shape = species → can be used
simultaneous

write inside aes for chart in regard to variable
<u>outside for all!</u>

---

● alpha: controls transparency

mapping = aes(

ggplot (data = df) + geom-point ( x= " ", y = " ",
color = species, shape = species, size = species))

mapping = aes(
geom-point (x = " ", y = " ", ● ), color = "purple")

———×——————×————

---

* ggplot (data = df) +
  geom-smooth (mapping = aes (x, y,)) +
  geom-point (mapping = aes (x, y, ))

∫ | linetype = species | → different line for different
  species

———————×—————

geom-jitter ( mapping = aes (x, y, ))
scatter plot with random noise

Jittering helps us deal with overplotting

———×————×.

> In simpler words, when you learn the basic steps for creating a plot in ggplot2.

> You can reuse these steps to create lots of different kinds of plots.

+> You can add/remove layers of detail to your plot without changing its basic structure or underlying data.

---

Upcoming :- Aesthetics, Geoms, Facets, labels & annotations

⊛ Aesthetic :→ A visual property of an object in your plot.

ex. in scatter plots, aesthetics are :→ size, shape, color

⊛ Geom :→ Geometric object used to represent data
ex. points → scatter bars → barchart
line → line chart

⊛ Facet → let you display smaller group, or subsets, of your data
ex. separate plot for all variable,

⊛ label & Annotations → lets your customize your plot.

\* Hands ion →

   install·packages ("___")
   library ("H")

   pacman :: p_load (pacman, ggplot2, palmer penguins)

   penDF = data.frame (penguins)

\* \* \*   ggplot (data = penDF) + geom_point (mapping =
        aes(x=flipper_length_mm, y = body_mass_g))

   ggplot (data = penDF)    → create a plane

   [ Add layers ]

   geom_point ( ~~aes~~ mapping =aes( x= flipper_length_mm,
                  y = body_mass_g ))

   geom functions → different function for
   different plots.
   geom_point → Scatterplot | geom_bar → Bar chart

   Each geom takes a mapping argument and
   is always paired with aes function
   which takes x & y variable to map

Same code → ggplot (data = penguins, mapping = aes(x=flipper
Different                               y= body _H_ ))
syntax     + geom_point ()

library ("diamonds")    # already loaded
ggplot (data = diamonds)
geom_bar( mapping = aes(x = cut ))
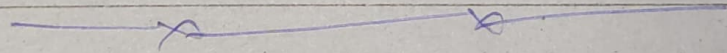
↓

by default calculates    row & numbers

↓

how many time    row element is present


geom_bar (mapping =    aes (x = cut, color = cut ))
                                    fill = cut.

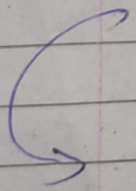↓

color: outlines the bare
fill & fills the bare



Damn



facet

Shows different views    of your data
> used for comparison!


> facet_wrap() : single variable
> facet_grid()


ggplot (data = penguins) +

geom_point (mapping = aes(x, y, color = Species)) +
facet_wrap (~ species)


# used when viz is too dense. 3 diff plots for 3 vars

## facet-grid () :→

```
ggplot (data = penguins)
   geom-point (mapping = aes (x = flmm, y = bmg,
                                 color = species)) +
     facet-grid ( sex N species)
```

Splits / divides viz into vertically by first variable
                           horizontally by second variable

< ✳ >  Can also use
       just one variable

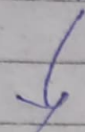```
facet-grid (~ sex)      or    facet-grid (~ species)
```

————————✗————————

Annotate :→

To add notes to a document or diagram to explain or comment upon it

+ labs ( title = "palmer Penguins : Body mass vs. Flipper len "))

↓

arguments → title, subtitle.
Caption : Data source

p ← ggplot () + geom-fin () + labs ()

p+ annotate ( "text", x= 220, y= 3500, label = "The hentoos are the Largest", color = "purple", fontface = "bold", size = 4.5, angle = 25)

—×—⟶

\* Saving your Visualization

➢ gg save ()
➢ Export option (plots pane)
          ➢ save as Img or pdf   (Total → 6)
                                  Image format

ggsare () → saves last plot & current graphi
                                  device size

ggs ave ("Name.png")

—×—