

# SESSION 3: DATA VISUALIZATION

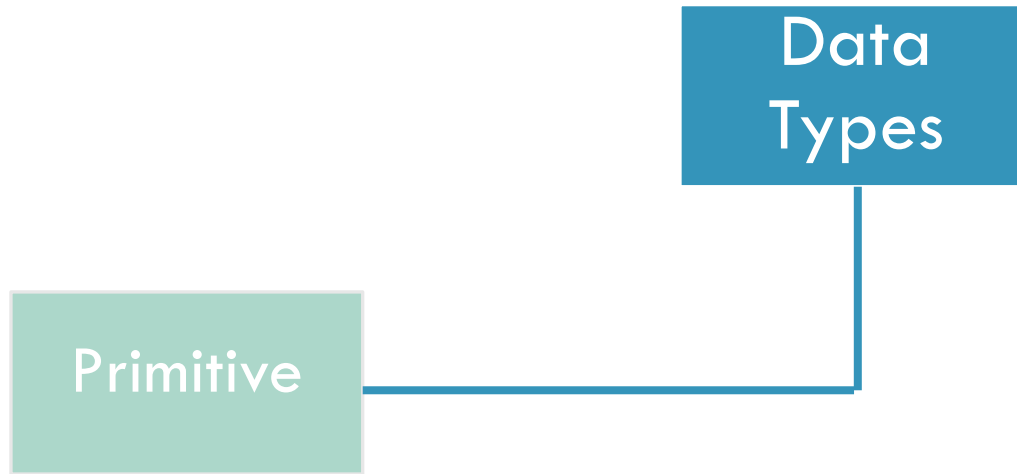
DR. SOFIA GIL-CLAVEL

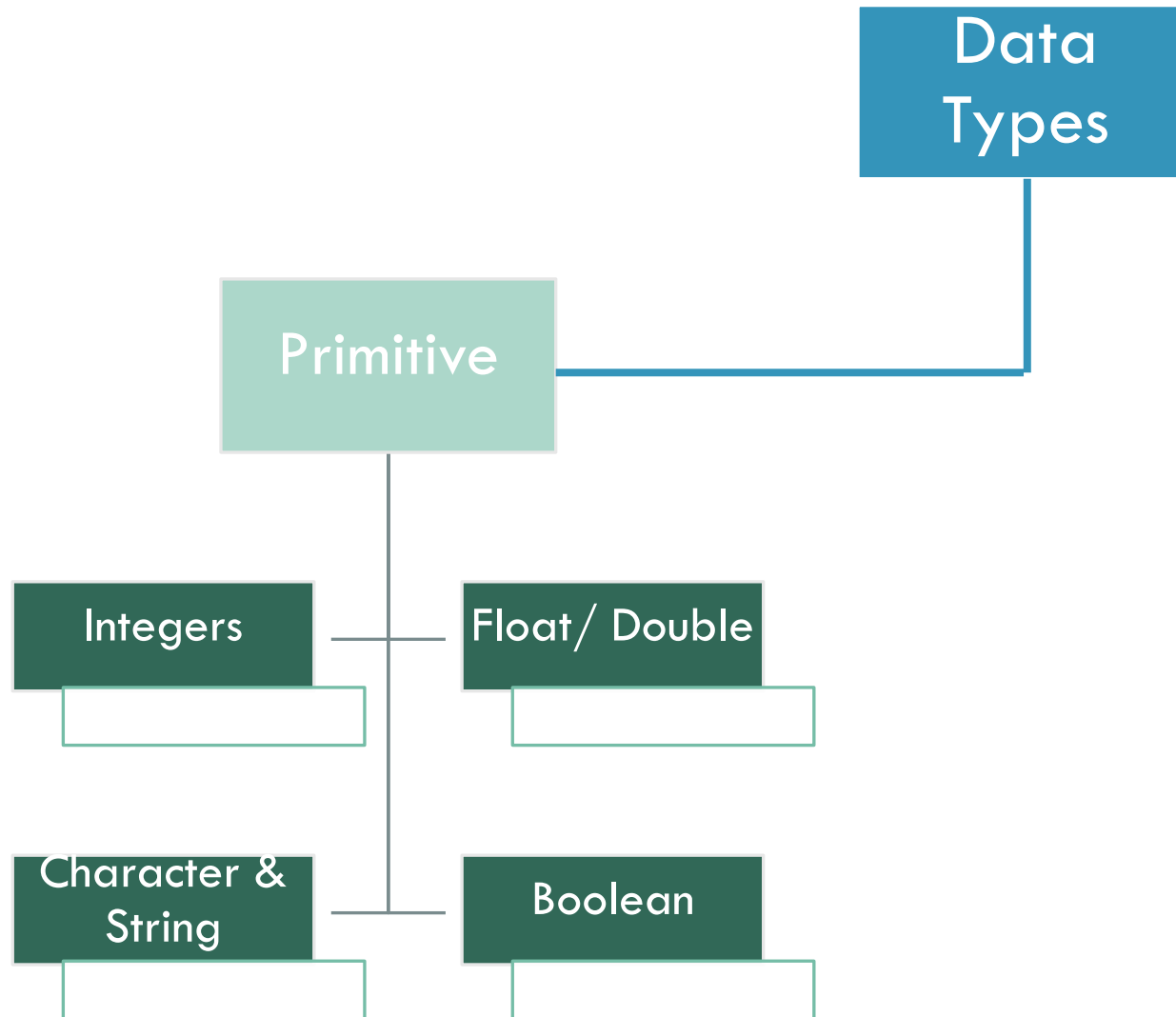
- ❖ Session 1-2: Recap
- ❖ Ggplot2 & the Grammar of Graphs

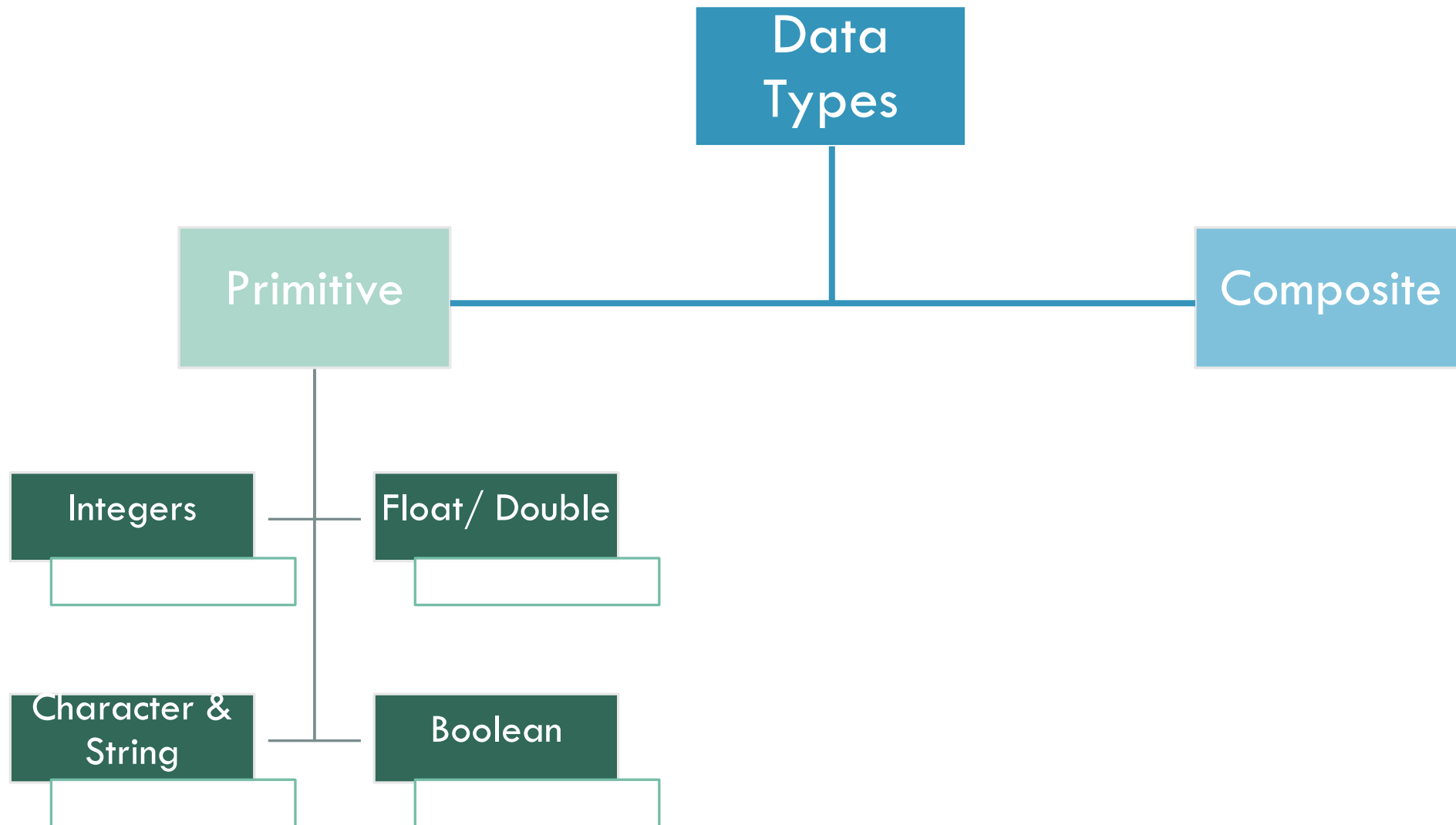
# 1. RECAP

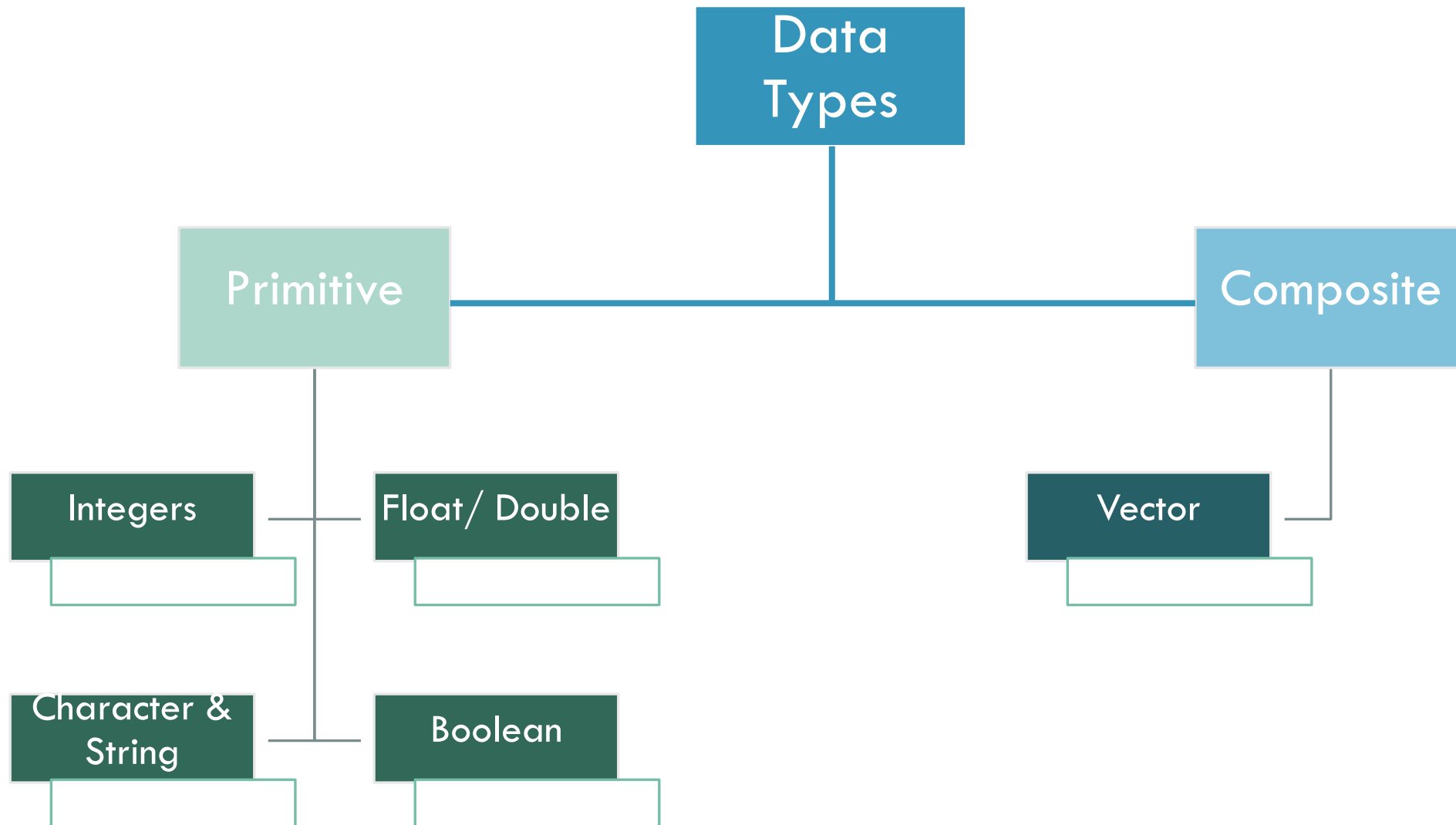
---

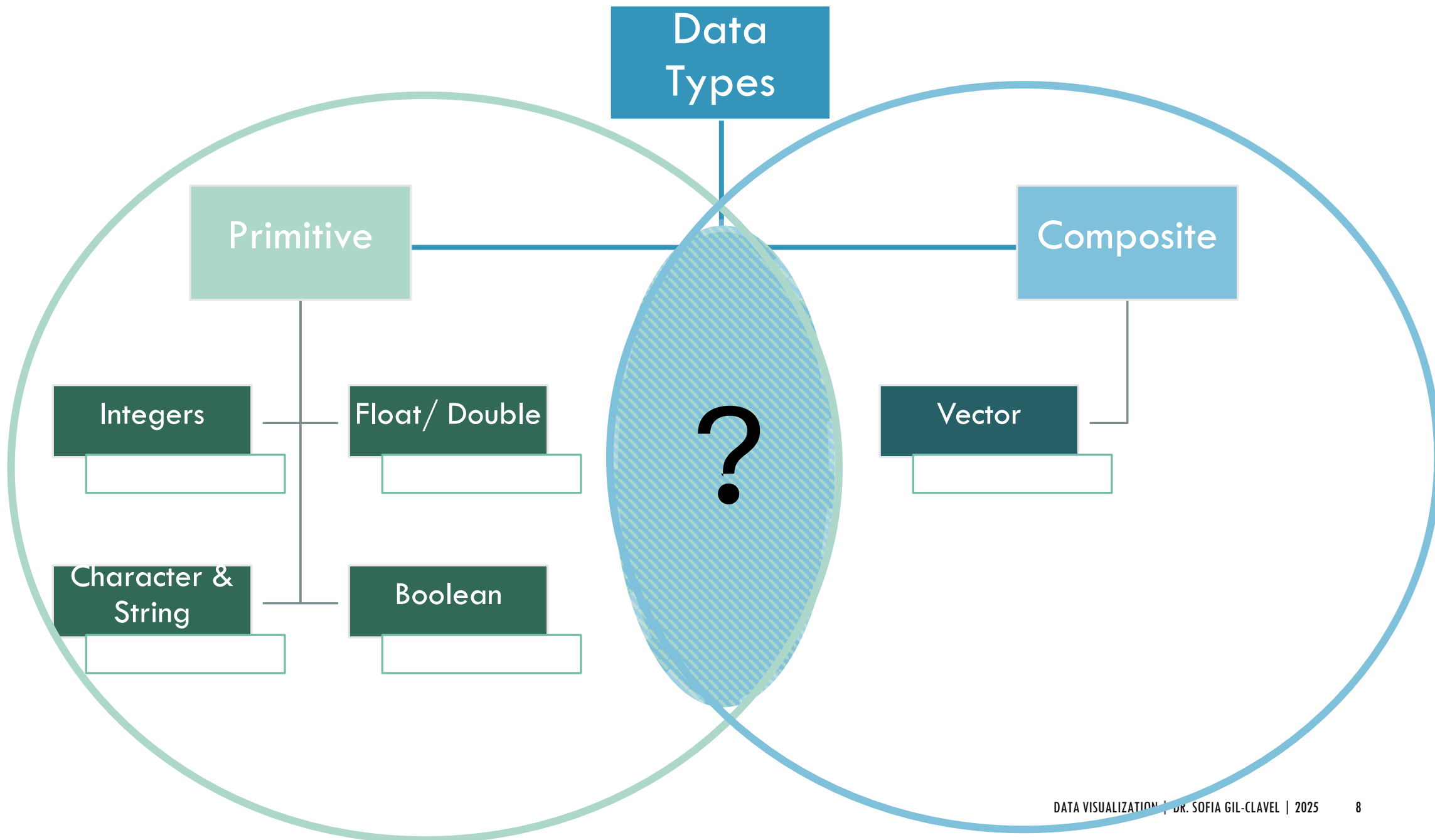
# Data Types



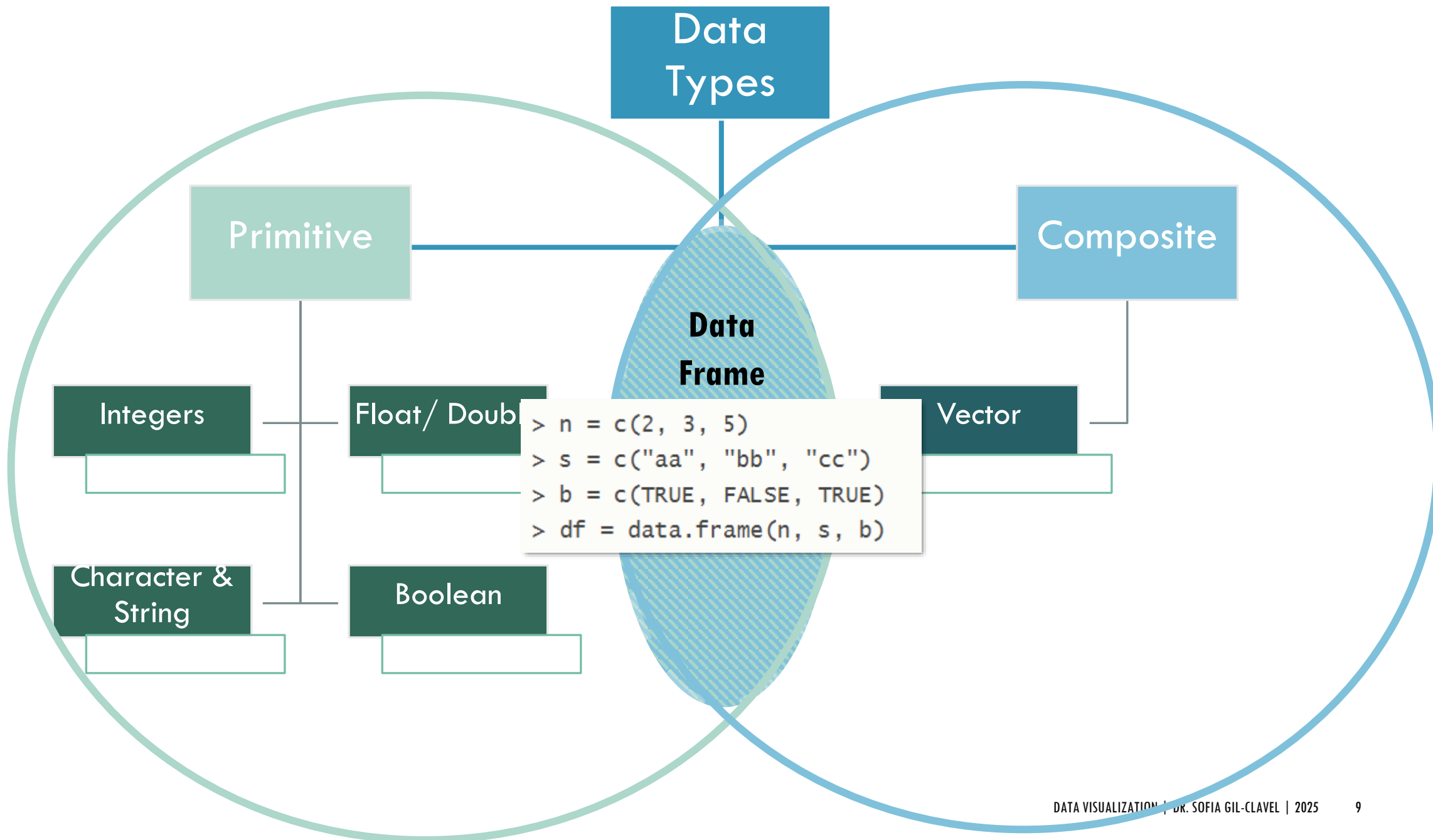












# PRIMITIVE DATA AND ITS OPERATORS

Operators					
Arithmetic		Comparative		Logic	
+	addition	<	less than	! x	logic NO
-	subtraction	>	more than	x & y	element-wise AND
*	multiplication	<=	less or equal than	x && y	single comparison AND
/	division	>=	more or equal than	x   y	element-wise OR
^	power	==	equal than	x    y	single comparison OR
%%	integer division	!=	different than	xor(x,y)	exclusive OR

# DATA PIPELINES

1.



```
> iris
# A tibble: 150 x 5
  Sepal.Length Sepal.Width Petal.Length Petal.Width Species
    <dbl>         <dbl>         <dbl>         <dbl>   <fct>
1         5.1           3.5           1.4           0.2 setosa
2         4.9           3           1.4           0.2 setosa
3         4.7           3.2           1.3           0.2 setosa
4         4.6           3.1           1.5           0.2 setosa
5          5           3.6           1.4           0.2 setosa
6         5.4           3.9           1.7           0.4 setosa
7         4.6           3.4           1.4           0.3 setosa
8          5           3.4           1.5           0.2 setosa
9         4.4           2.9           1.4           0.2 setosa
10        4.9           3.1           1.5           0.1 setosa
#> 140 more rows
#> print(n = ...) to see more rows
```

2.



## Basic piping

- `x %>% f` is equivalent to `f(x)`
- `x %>% f(y)` is equivalent to `f(x, y)`
- `x %>% f %>% g %>% h` is equivalent to `h(g(f(x)))`

3.



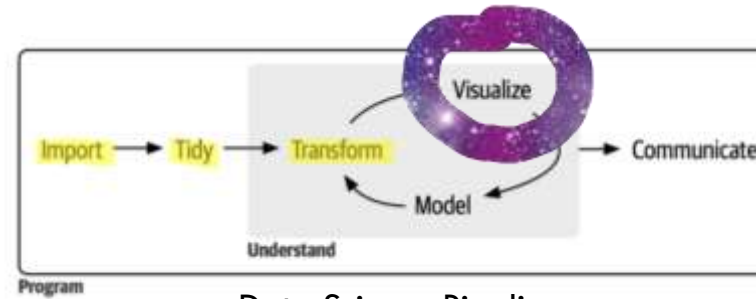
dplyr is a grammar of data manipulation, providing a consistent set of verbs that help you solve the most common data manipulation challenges:

- `mutate()` adds new variables that are functions of existing variables
- `select()` picks variables based on their names.
- `filter()` picks cases based on their values.
- `summarise()` reduces multiple values down to a single summary.
- `arrange()` changes the ordering of the rows.

## 2. GGPLOT2 & THE GRAMMAR OF GRAPHS

---

# ❖ GGPLOT: [HTTPS://WWW.TIDYVERSE.ORG/](https://www.tidyverse.org/)



ggplot2 is a system for declaratively creating graphics, based on [The Grammar of Graphics](#). You provide the data, tell ggplot2 how to map variables to aesthetics, what graphical primitives to use, and it takes care of the details.

Source: <https://ggplot2.tidyverse.org/>



# GGPLOT2

**ggplot2** is a package for creating graphs, which follows the grammar of graphs.

## What is the grammar of graphs?

The grammar of graphs tells us that "... Every statistical graph is a mapping of data to a set of geometric objects (points, lines, bars) that contain aesthetic attributes (color, shape, size). The graph can even have statistical transformations of the data, and these are drawn on a specific Cartesian plane."

# THE COMPONENTS OF A GRAPH

- 1) The **data** we want to visualize and a set of aesthetic transformations that describe how the variables in the data behave.
- 2) **Layers** made of geometric objects (which we'll call **geoms**) and statistical transformations (which we'll call **stats**). Geoms represent what you see on the graph: points, lines, polygons, etc. Stats are a summary of the data being observed.
- 3) **Scales** maps the values of the data to values in an aesthetic space, such as color, size, or shape. Scales draws the legends or axes, which provides reverse mapping that makes it possible to read the original data from the graph.
- 4) The facet specification (**facet**) describes how to divide data into subsets and how to display it in subgraphs within the same graph. This is also called a conditioner.
- 5) A **theme** controls the points that are displayed, such as font size and color.

A quick guide is here: <https://github.com/rstudio/cheatsheets/blob/main/data-visualization.pdf>



# 1) THE DATA

```
library(ggplot2)  
mpg
```

**Data**

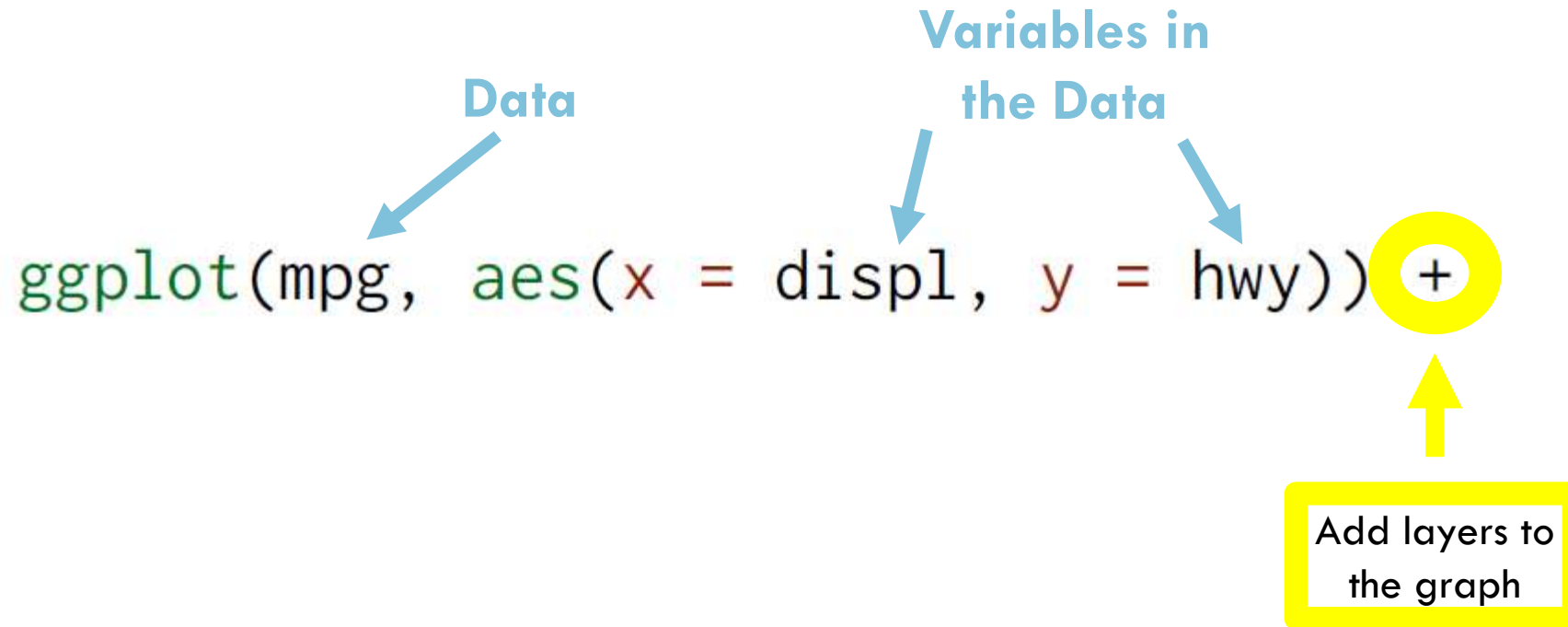
**Variables in the Data**

```
ggplot(mpg, aes(x = displ, y = hwy))
```



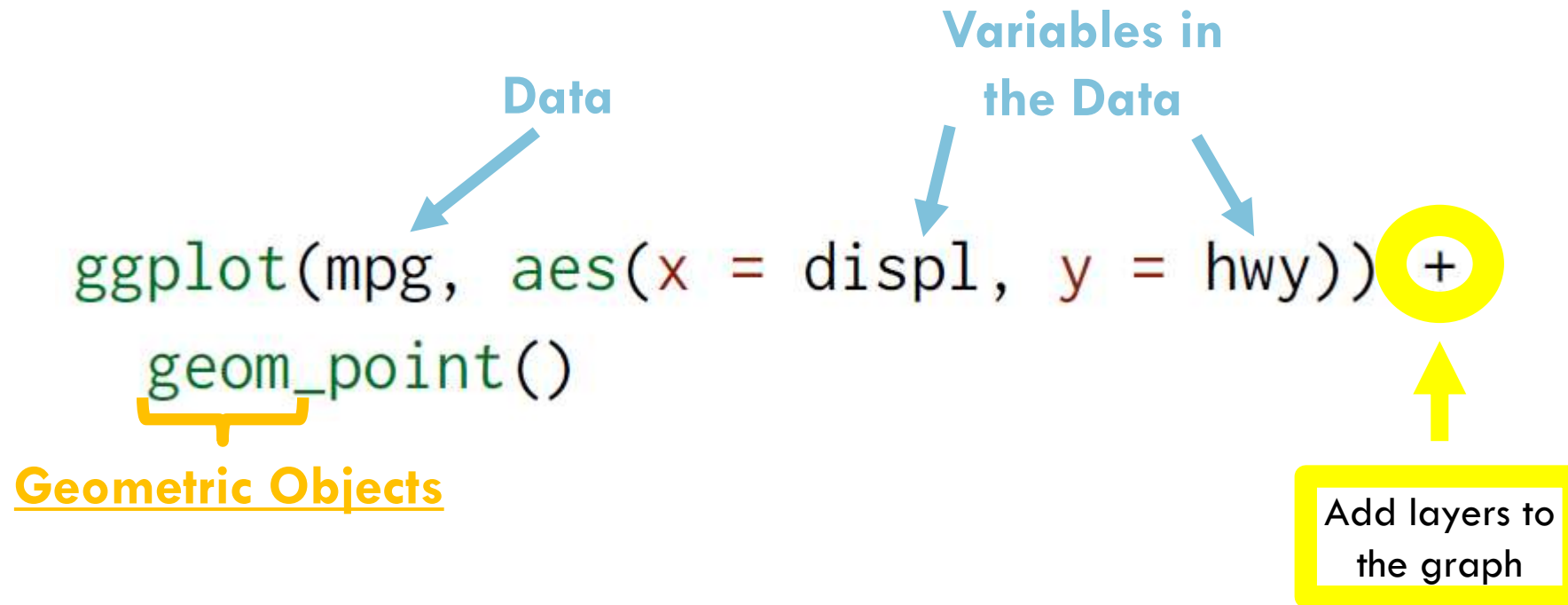
## 2) LAYERS

```
library(ggplot2)  
mpg
```



## 2) LAYERS

```
library(ggplot2)  
mpg
```



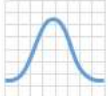
## 2) LAYERS: BASIC GEOMETRIES

### ONE VARIABLE continuous

```
c <- ggplot(mpg, aes(hwy)); c2 <- ggplot(mpg)
```



**c + geom\_area**(stat = "bin")  
x, y, alpha, color, fill, linetype, size



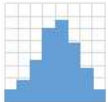
**c + geom\_density**(kernel = "gaussian")  
x, y, alpha, color, fill, group, linetype, size, weight



**c + geom\_dotplot**()  
x, y, alpha, color, fill



**c + geom\_freqpoly**()  
x, y, alpha, color, group, linetype, size



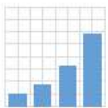
**c + geom\_histogram**(binwidth = 5)  
x, y, alpha, color, fill, linetype, size, weight



**c2 + geom\_qq**(aes(sample = hwy))  
x, y, alpha, color, fill, linetype, size, weight

### discrete

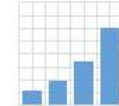
```
d <- ggplot(mpg, aes(fl))
```



**d + geom\_bar**()  
x, alpha, color, fill, linetype, size, weight

### one discrete, one continuous

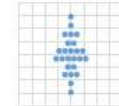
```
f <- ggplot(mpg, aes(class, hwy))
```



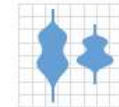
**f + geom\_col**()  
x, y, alpha, color, fill, group, linetype, size



**f + geom\_boxplot**()  
x, y, lower, middle, upper, ymax, ymin, alpha, color, fill, group, linetype, shape, size, weight



**f + geom\_dotplot**(binaxis = "y", stackdir = "center")  
x, y, alpha, color, fill, group



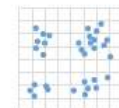
**f + geom\_violin**(scale = "area")  
x, y, alpha, color, fill, group, linetype, size, weight

### both discrete

```
g <- ggplot(diamonds, aes(cut, color))
```



**g + geom\_count**()  
x, y, alpha, color, fill, shape, size, stroke



**e + geom\_jitter**(height = 2, width = 2)  
x, y, alpha, color, fill, shape, size

`library(ggplot2)`  
`mpg`

# EXERCISE 2.1

□ From slide 19 choose one graph from the following groups and save them as:

*A=From the “one continues variable”*

*B=From the “one discrete one continuous variables”*

## 2) LAYERS: A NON-INTUITIVE CASE

What do the following functions do?

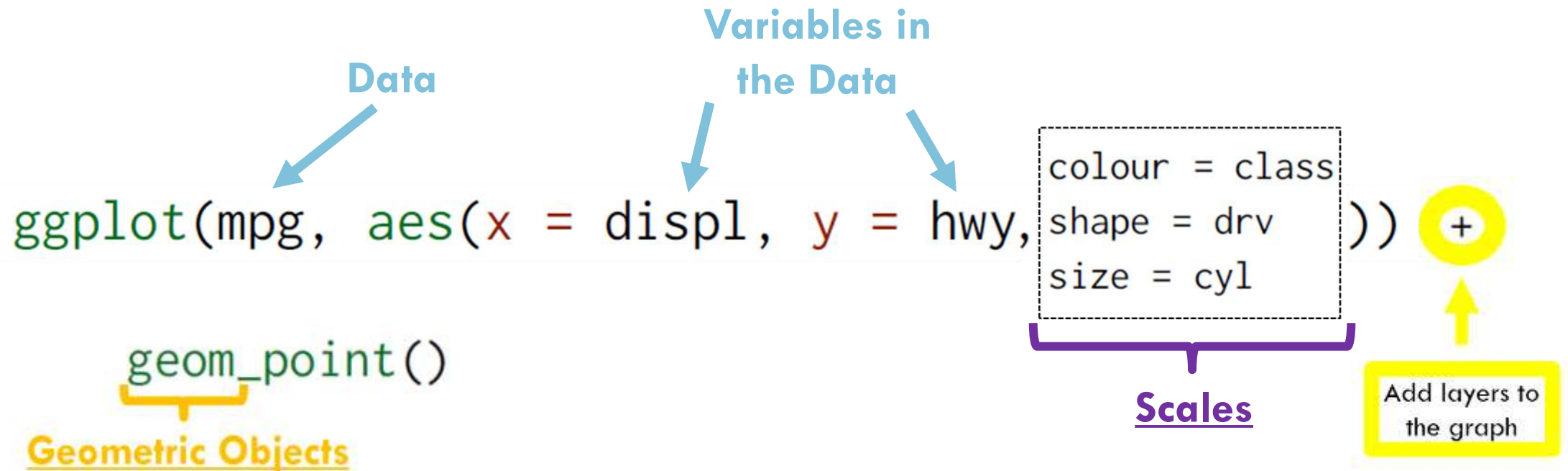
```
ggplot(mpg, aes(displ, hwy)) + geom_point(aes(colour = "blue"))
```

```
ggplot(mpg, aes(displ, hwy)) + geom_point(aes(colour="blue"))
```

```
ggplot(mpg, aes(displ, hwy)) + geom_point(aes(colour=class))
```

```
ggplot(mpg, aes(displ, hwy)) + geom_point(colour="blue")
```

### 3) SCALES: COLOR, SIZE, SHAPE, AND OTHER AESTHETIC ATTRIBUTES

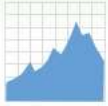




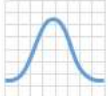
### 3) SCALES: COLOR, SIZE, SHAPE, ETC.

#### ONE VARIABLE continuous

```
c <- ggplot(mpg, aes(hwy)); c2 <- ggplot(mpg)
```



**c + geom\_area**(stat = "bin")  
x, y, alpha, color, fill, linetype, size



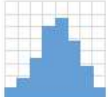
**c + geom\_density**(kernel = "gaussian")  
x, y, alpha, color, fill, group, linetype, size, weight



**c + geom\_dotplot**()  
x, y, alpha, color, fill



**c + geom\_freqpoly**()  
x, y, alpha, color, group, linetype, size



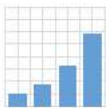
**c + geom\_histogram**(binwidth = 5)  
x, y, alpha, color, fill, linetype, size, weight



**c2 + geom\_qq**(aes(sample = hwy))  
x, y, alpha, color, fill, linetype, size, weight

#### discrete

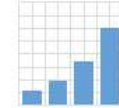
```
d <- ggplot(mpg, aes(fl))
```



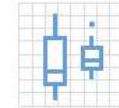
**d + geom\_bar**()  
x, alpha, color, fill, linetype, size, weight

#### one discrete, one continuous

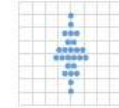
```
f <- ggplot(mpg, aes(class, hwy))
```



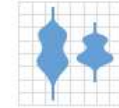
**f + geom\_col**()  
x, y, alpha, color, fill, group, linetype, size



**f + geom\_boxplot**()  
x, y, lower, middle, upper, ymax, ymin, alpha, color, fill, group, linetype, shape, size, weight



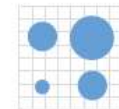
**f + geom\_dotplot**(binaxis = "y", stackdir = "center")  
x, y, alpha, color, fill, group



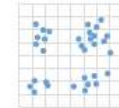
**f + geom\_violin**(scale = "area")  
x, y, alpha, color, fill, group, linetype, size, weight

#### both discrete

```
g <- ggplot(diamonds, aes(cut, color))
```



**g + geom\_count**()  
x, y, alpha, color, fill, shape, size, stroke



**e + geom\_jitter**(height = 2, width = 2)  
x, y, alpha, color, fill, shape, size

library(ggplot2)  
mpg

## EXERCISE 2.2

If possible, add some of the following scales to the graphs you chose in exercise 3.1.  
How does the graph change when...?

- You add “color”
- You add “shape”
- You add “size”

What happens when you choose more than one scale?

If possible, choose two of the previous scales for each graph created in exercise 3.1  
and update the variables A and B.



### 3) SCALES: SUBGROUPS POSITIONS

- `position_stack()`: stack overlapping bars (or areas) on top of each other.
- `position_fill()`: stack overlapping bars, scaling so the top is always at 1.
- `position_dodge()`: place overlapping bars (or boxplots) side-by-side.

```
dplot <- ggplot(diamonds, aes(color, fill = cut)) +  
  xlab(NULL) + ylab(NULL) + theme(legend.position = "none")
```

```
dplot + geom_bar()
```

```
dplot + geom_bar(position = "fill")
```

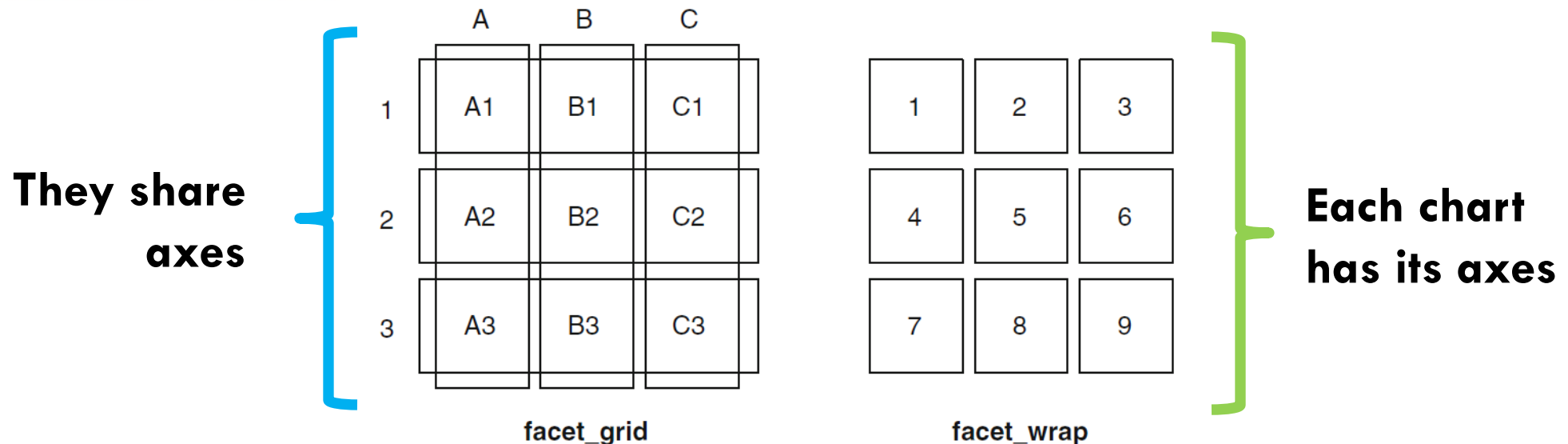
```
dplot + geom_bar(position = "dodge")
```

## EXERCISE 2.3

- ❑ For the “*one discrete*” option (`geom_bar`) of slide 20, use the different discrete positions. Remember to add the “fill” scale. What do the Y do to the graph?

## 4) FACETS

- `facet_null()`: a single plot, the default.
- `facet_wrap()`: “wraps” a 1d ribbon of panels into 2d.
- `facet_grid()`: produces a 2d grid of panels defined by variables which form the rows and columns.



## EXERCISE 2.4

From the **mpg** database, what variables would you choose to group the graph by facet?

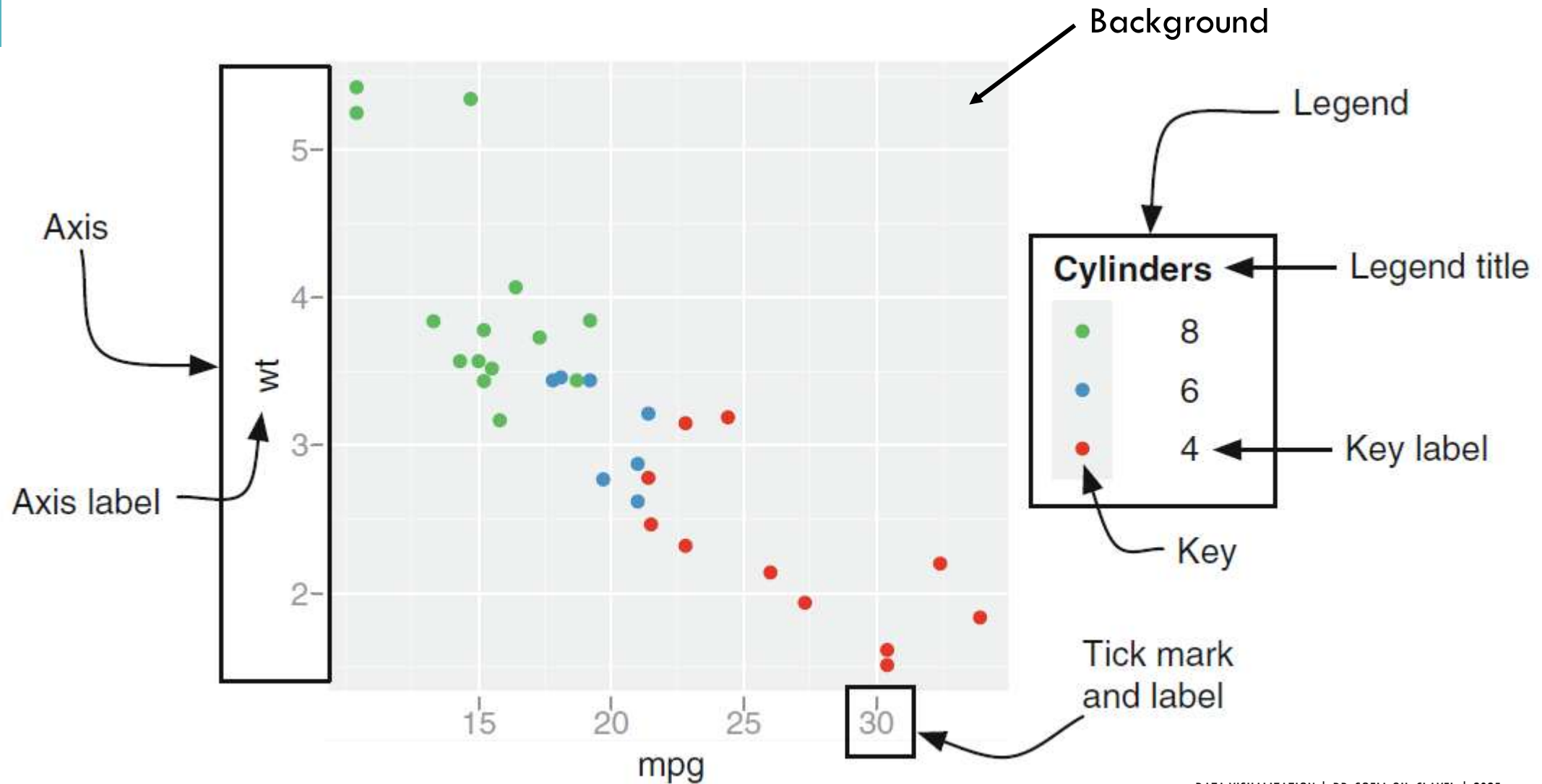
Apply the grouping to each of A and B graphs from exercise 3.2.

## 5) THEMES

**Themes** allows you to exercise fine **control** over **the non-data elements of your plot**.

The theme system does not affect how the data is rendered by geoms, or how it is transformed by scales. Themes don't change the perceptual properties of the plot, but **they do help you make the plot aesthetically pleasing** or match an existing style guide. **Themes give you control over** things like **fonts, ticks, panel strips, and backgrounds**.

# 5) THEMES: MODIFYING THE APPEARANCE

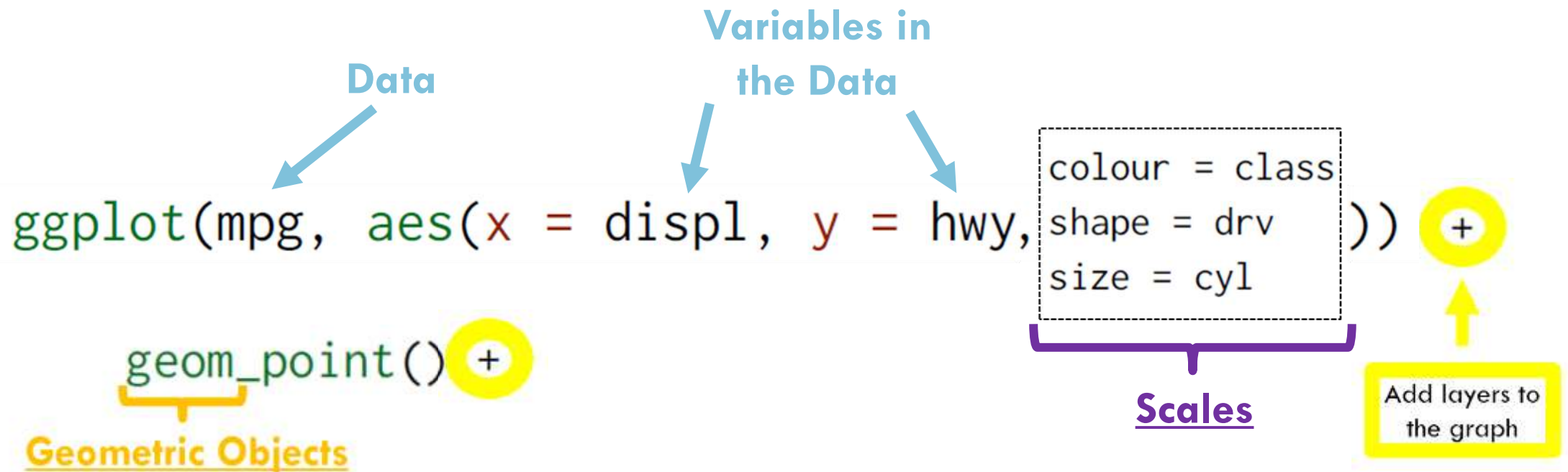


## 5) THEMES: PREDEFINED

- `theme_bw()`: a variation on `theme_grey()` that uses a white background and thin grey grid lines.
- `theme_linedraw()`: A theme with only black lines of various widths on white backgrounds, reminiscent of a line drawing.
- `theme_light()`: similar to `theme_linedraw()` but with light grey lines and axes, to direct more attention towards the data.
- `theme_dark()`: the dark cousin of `theme_light()`, with similar line sizes but a dark background. Useful to make thin coloured lines pop out.
- `theme_minimal()`: A minimalistic theme with no background annotations.
- `theme_classic()`: A classic-looking theme, with x and y axis lines and no gridlines.
- `theme_void()`: A completely empty theme.



## 5) THEMES FROM SCRATCH



```
theme(  
  plot.title = element_text(face = "bold", size = 12),  
  legend.background = element_rect(fill = "white", size = 4, colour = "white"),  
  legend.justification = c(0, 1),  
  legend.position = c(0, 1),  
  axis.ticks = element_line(colour = "grey70", size = 0.2),  
  panel.grid.major = element_line(colour = "grey70", size = 0.2),  
  panel.grid.minor = element_blank()  
)
```



## EXERCISE 2.5

Using the different **themes**' elements on your A and B graphs, create the most horrible graph anyone has ever seen and update the variables A and B.

# SAVING THE GRAPHS

ggplot2 provides a convenient shorthand to save graphs created with ggplot: **ggsave()**

```
ggplot(mpg, aes(displ, cty)) + geom_point()  
ggsave("output.pdf")
```

## EXERCISE 2.6

- ❑ In the “help” section, check the parameters of **ggsave**.
- ❑ Use **ggsave** to save your graphs A and B. Play around with the parameters until you manage to save your graph in a good quality format.

# REFERENCES

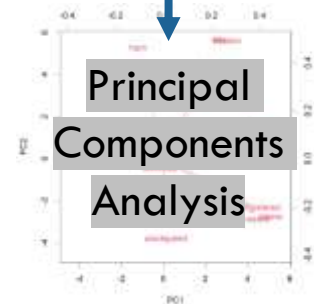
- ❑ Wickham, Hadley, Mine Çetinkaya-Rundel, and Garrett Grolemund. *R for data science*. " O'Reilly Media, Inc.", 2023. Accessed May 7, 2024.  
<https://r4ds.hadley.nz/>.
- ❑ Wickham, Hadley. *Ggplot2. Use R!* Cham: Springer International Publishing, 2016.  
<https://doi.org/10.1007/978-3-319-24277-4>.
- ❑ Wilke, C. *Fundamentals of Data Visualization: A Primer on Making Informative and Compelling Figures*. First edition. Sebastopol, CA: O'Reilly Media, 2019.  
<https://clauswilke.com/dataviz/>

# SESSION 4: UNSUPERVISED LEARNING

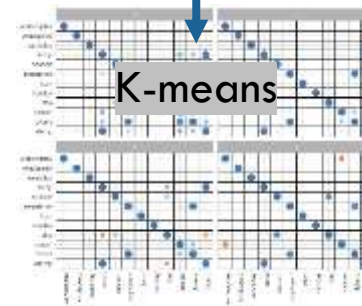
```
> iris
# A tibble: 150 × 5
  Sepal.Length Sepal.Width Petal.Length Petal.Width Species
    <dbl>         <dbl>         <dbl>         <dbl>    <fct>
1         5.1           3.5           1.4           0.2  setosa
2         4.9           3           1.4           0.2  setosa
3         4.7           3.2           1.3           0.2  setosa
4         4.6           3.1           1.5           0.2  setosa
5          5           3.6           1.4           0.2  setosa
6         5.4           3.9           1.7           0.4  setosa
7         4.6           3.4           1.4           0.3  setosa
8          5           3.4           1.5           0.2  setosa
9         4.4           2.9           1.4           0.2  setosa
10        4.9           3.1           1.5           0.1  setosa
# 140 more rows
# Use `print(n = ...)` to see more rows
```

Part of  
exploratory  
data analysis!

Dimensionality  
reduction



Clustering





Join us!



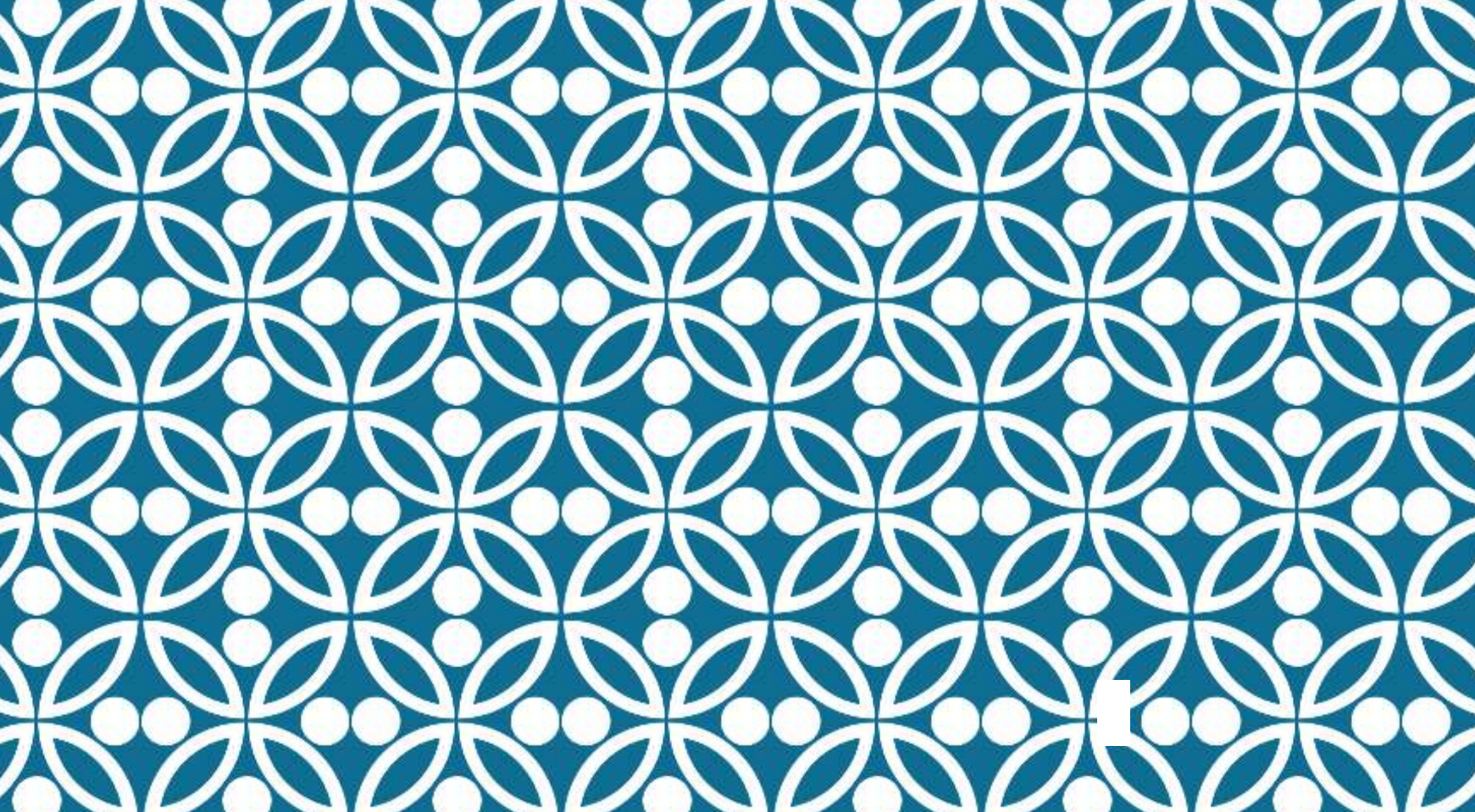
<https://forms.office.com/e/8Bgd2YsasJ>

All about the lab:

<https://societal-analytics.nl/>

Contact us at:

[analytics-lab.fsw@vu.nl](mailto:analytics-lab.fsw@vu.nl)



<https://sofiag1l.github.io/>

# THANKS!

Dr. Sofia Gil-Clavel