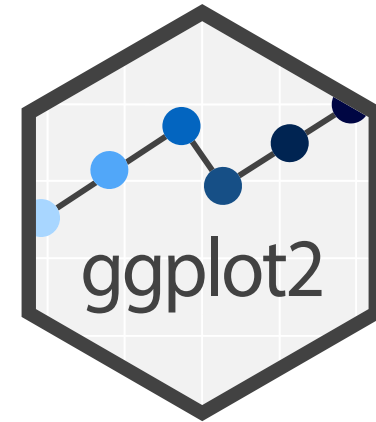# Data handling and visualisation in R

Marine Ecosystem Dynamics
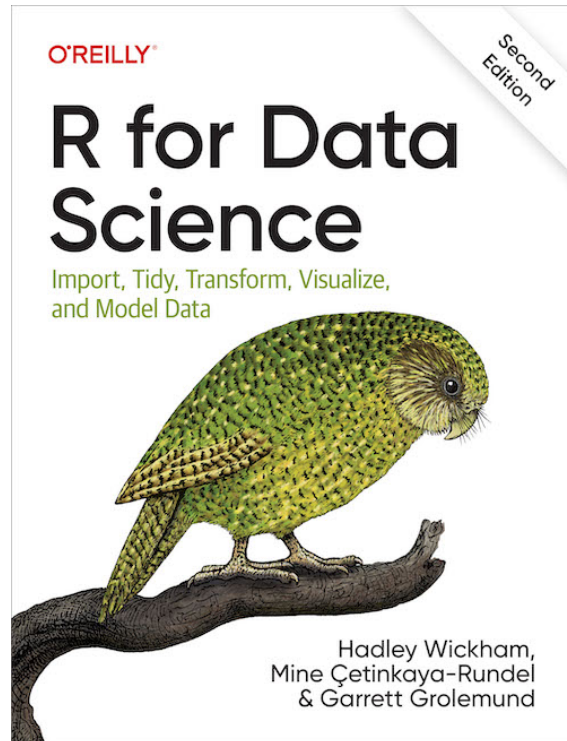
Kinlan M.G. Jan

Stockholm University

# Plan for today's lecture

- Introduction to `tidyverse`

- Pipe the data using `magrittr`

- Clean the data using `tidyr`

- Arrange the data using `dplyr`

- Plot using `ggplot2`

# Tidyverse



R for Data Science
Import, Tidy, Transform, Visualize, and Model Data

Hadley Wickham, Mine Çetinkaya-Rundel & Garrett Grolemund

- `tidyverse` is a collection of packages

- It is now a standard in data analysis

- It is easier to read and keep track of what is happening with the pipe operator `%>%`

The book is available online

# Pipe the data using `magrittr`

`%>%` takes the data from the left and place it to the right

- `x %>% function()` = `function(x)`

Without the pipe operator:

```
1  iris_subset <- iris[, c("Species", "Sepal.
2  iris_subset$Sepal_Ratio <- iris_subset$Sep
3
4  aggregate(Sepal_Ratio ~ Species, data = i
5                         FUN = function(x) c(A
6                                            st
```

With the pipe operator:

```
1  library(magrittr)
2  iris %>%
3    dplyr::select(Species, Sepal.Length, Sep
4    dplyr::mutate(Sepal_Ratio = Sepal.Length
5    dplyr::group_by(Species)%>%
6    dplyr::summarise(Average_ratio = mean(Se
7                     standard_deviation = so
```

R 4.1.0 introduced a native pipe operator, `|>`

Stockholm University

# Tidy the data with `tidyr`

A table is tidy if:

- Each variable is in its own column

- Each observation is in its own row

Key functions:

- `pivot_longer`
- `pivot_wider`
- `unite`
- `separate`



variables



observations



values

source: https://r4ds.had.co.nz/tidy-data.html

# Tidy the data with `tidyr` - `iris` example

| Sepal.Length | Sepal.Width | Petal.Length | Petal.Width | Species | id |
|---:|---:|---:|---:|:---|---:|
| 5.1 | 3.5 | 1.4 | 0.2 | setosa | 1 |
| 4.9 | 3.0 | 1.4 | 0.2 | setosa | 2 |
| 4.7 | 3.2 | 1.3 | 0.2 | setosa | 3 |
| 4.6 | 3.1 | 1.5 | 0.2 | setosa | 4 |

```
1  iris %<>% dplyr::mutate(id = 1:150)
2  iris |> head(4)
```

Stockholm University

# Tidy the data with `tidyr` - `pivot_longer`

| Species | id | Parameter | Size |
|---------|----|-----------|------|
| setosa | 1 | Sepal.Length | 5.1 |
| setosa | 1 | Sepal.Width | 3.5 |
| setosa | 1 | Petal.Length | 1.4 |
| setosa | 1 | Petal.Width | 0.2 |

```r
1  long_iris <- iris |>
2    tidyr::pivot_longer(1:4,
3                        names_to = "Parameter",
4                        values_to = "Size")
5  long_iris |> head(4)
```

Stockholm University

# Tidy the data with `tidyr` - `separate`

| Species | id | Organ | Measure | Size |
|---------|-----|-------|---------|------|
| setosa  | 1   | Sepal | Length  | 5.1  |
| setosa  | 1   | Sepal | Width   | 3.5  |
| setosa  | 1   | Petal | Length  | 1.4  |
| setosa  | 1   | Petal | Width   | 0.2  |

```r
1  sep_iris <- long_iris |>
2    tidyr::separate(Parameter, into = c("Organ", "Measure"))
3  sep_iris |> head(4)
```

Stockholm University

# Tidy the data with tidyr - pivot_wider

| Species | id | Measure | Sepal | Petal |
|---------|----|---------|-------|-------|
| setosa | 1 | Length | 5.1 | 1.4 |
| setosa | 1 | Width | 3.5 | 0.2 |
| setosa | 2 | Length | 4.9 | 1.4 |
| setosa | 2 | Width | 3.0 | 0.2 |

```r
1  wide_iris <- sep_iris |>
2    tidyr::pivot_wider(names_from = "Organ",
3                       values_from = "Size")
4  wide_iris |> head(4)
```

Stockholm
University

# Tidy the data with `tidyr` - `unite`

| Species | id | Measure | Sepal/ Petal |
|---------|-----|---------|--------------|
| setosa | 1 | Length | 5.1/1.4 |
| setosa | 1 | Width | 3.5/0.2 |
| setosa | 2 | Length | 4.9/1.4 |
| setosa | 2 | Width | 3/0.2 |

```
1  unite_iris <- wide_iris |>
2    tidyr::unite(col = "Sepal/Petal", c(Sepal, Petal), sep = "/")
3  unite_iris |>  head(4)
```

Stockholm
University

# Arrange the data with `dyplr`

`dplyr` simplifies the data manipulation with self-explanatory functions:

- `filter` observations based on their values

- `mutate` a new column as a function of others

- `select` variables based on their names

- `group_by` variable

- `summarise` the data

```
1  iris %>%
2    dplyr::filter(Petal.Length >= 1.4) %>%
3    dplyr::mutate(Sepal_Ratio = Sepal.Length/Sepal.Width) %>%
4    dplyr::select(Species, Sepal_Ratio) %>% # equaivalent to select(-c(Sepal.Length, Sepal.Wi
5    dplyr::group_by(Species) %>%
6    dplyr::summarise(Average_ratio = mean(Sepal_Ratio),
7                     standard_deviation = sd(Sepal_Ratio))
```
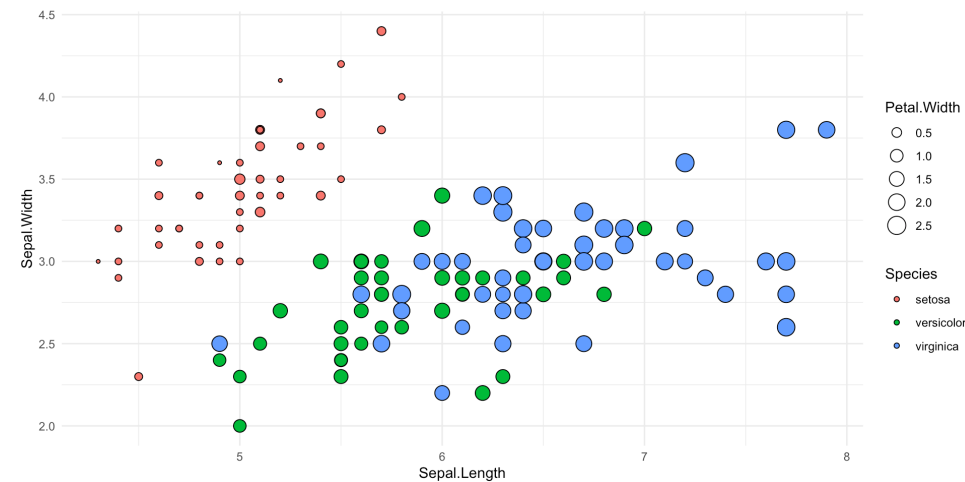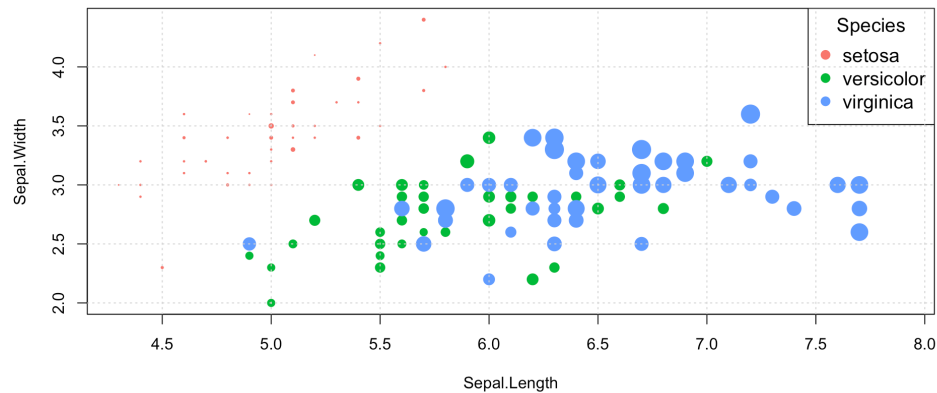
Stockholm
University

# Data visulalisation

It is very important to look at the data.

Totally different data might have similar statistics…



| statistics | value |
| --- | --- |
| mean_x | 54.27 |
| mean_y | 47.83 |
| sd_x | 16.77 |
| sd_y | 26.94 |

source: datasauRus

# Visualise the data with **ggplot2**



```
1  species_palette <- c("#F8766D", "#00BA38"
2  size_vector <- iris$Petal.Width
3  plot(x = iris$Sepal.Length,
4       y = iris$Sepal.Width,
5       col = species_palette[iris$Species],
6       bg = species_palette[iris$Species],
7       pch = 21,
8       cex = size_vector,
9       xlim = c(min(iris$Sepal.Length), max(
10      ylim = c(min(iris$Sepal.Width), max(
11      xlab = "Sepal.Length",
12      ylab = "Sepal.Width")
13 legend("topright", legend = levels(iris$Sp
14 grid(lwd = 1, lty = "dotted")
```

```
1  library(ggplot2)
2  ggplot(iris,
3        mapping = aes(x = Sepal.Length,
4                      y = Sepal.Width,
5                      fill = Species,
6                      size = Petal.Width))
7    geom_point(shape = 21,
8               col = 1) +
9    theme_minimal()
```
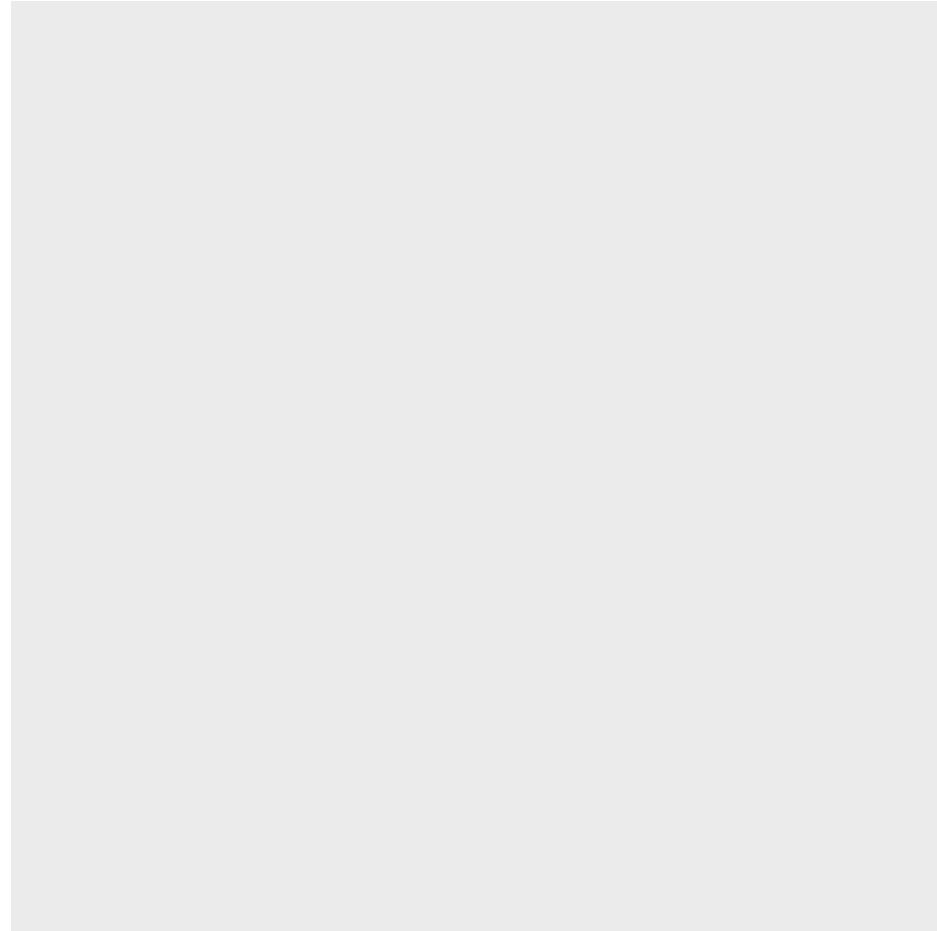
# Grammar of graphic (**gg**)

- Data

- Aesthetics - Visual characteristics (e.g., x, y, size)

- Geometry - How to represent the data (e.g., lines, point, boxplot)

- Statistics - What statistics to show

- Facets - Split the data

- Coordinates - Position of the geometry

- Themes - Visual changes

```
1  ggplot(data = <DATA>,
2         mapping = aes(<AESTHETICS>)) +
3     geom_<GEOMETRY>()
```
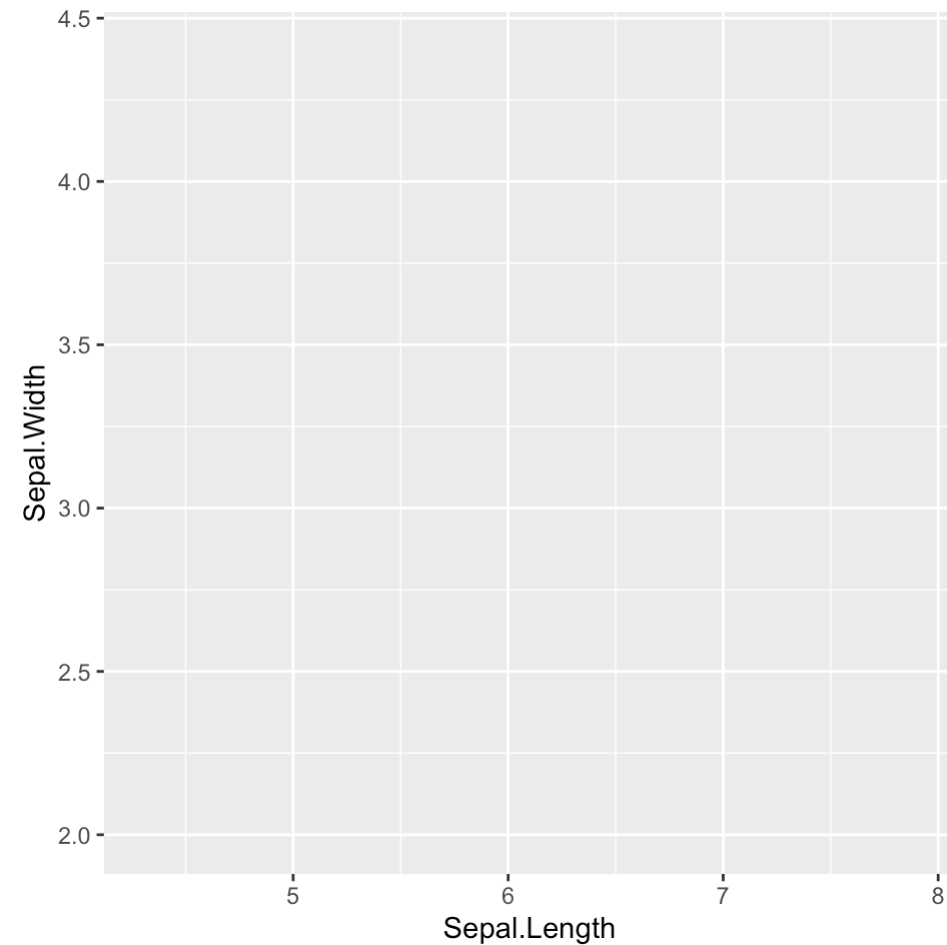
Stockholm University

# Let's plot using **ggplot2** - *Data*

```
1  ggplot(data = iris)
```

# Let's plot using **ggplot2** - *Aesthetics*

```
1  ggplot(data = iris,
2          mapping = aes(x = Sepal.Length,
3                           y = Sepal.Width))
```
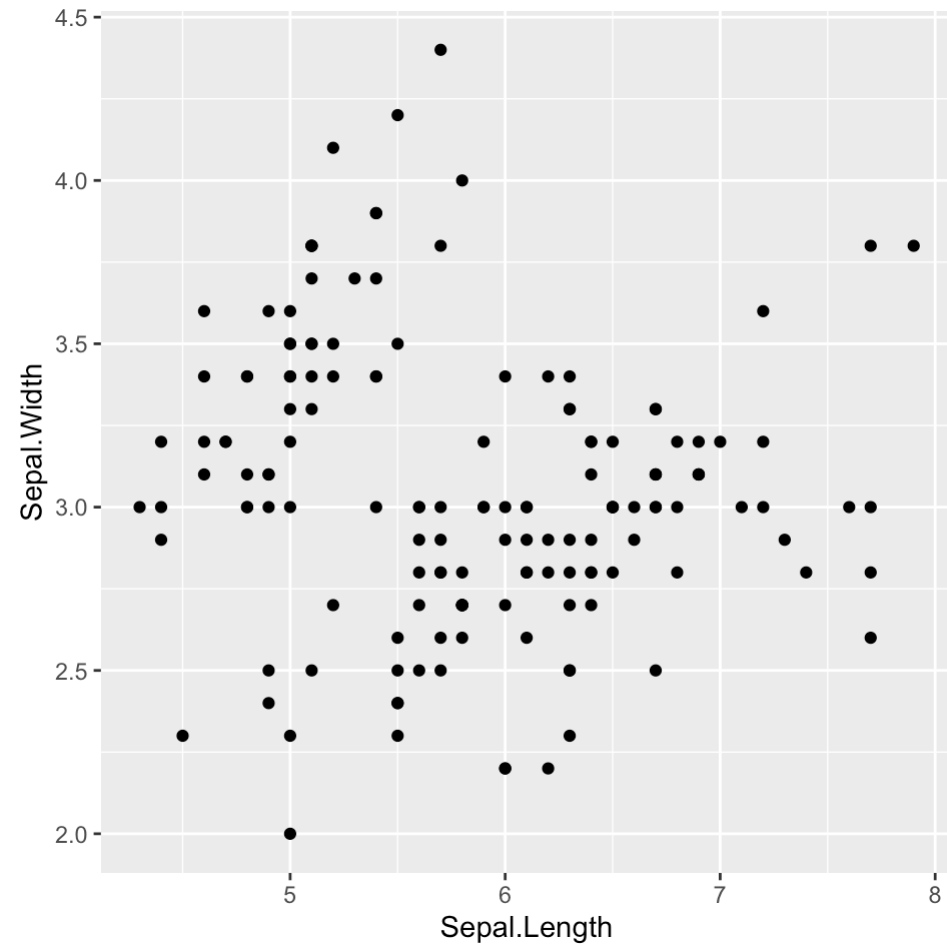


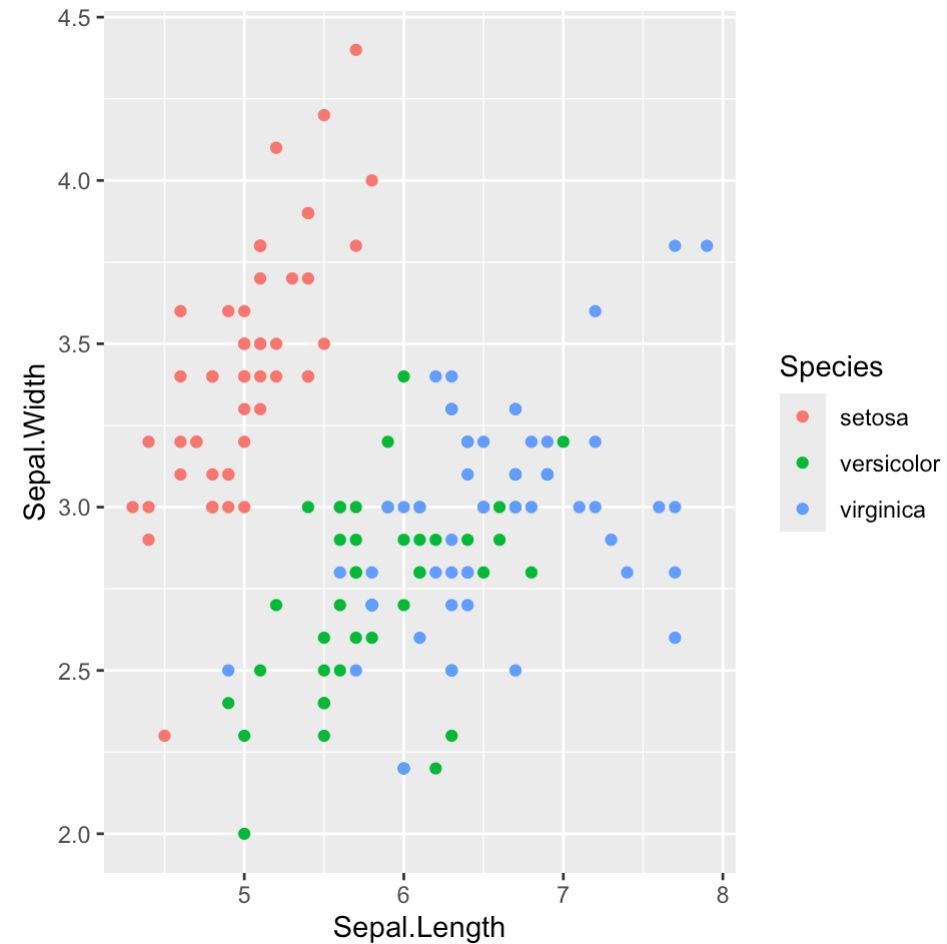Stockholm University

# Let's plot using **ggplot2** - *Geometry*

```
1  ggplot(data = iris,
2          mapping = aes(x = Sepal.Length,
3                            y = Sepal.Width)) +
4      geom_point()
```

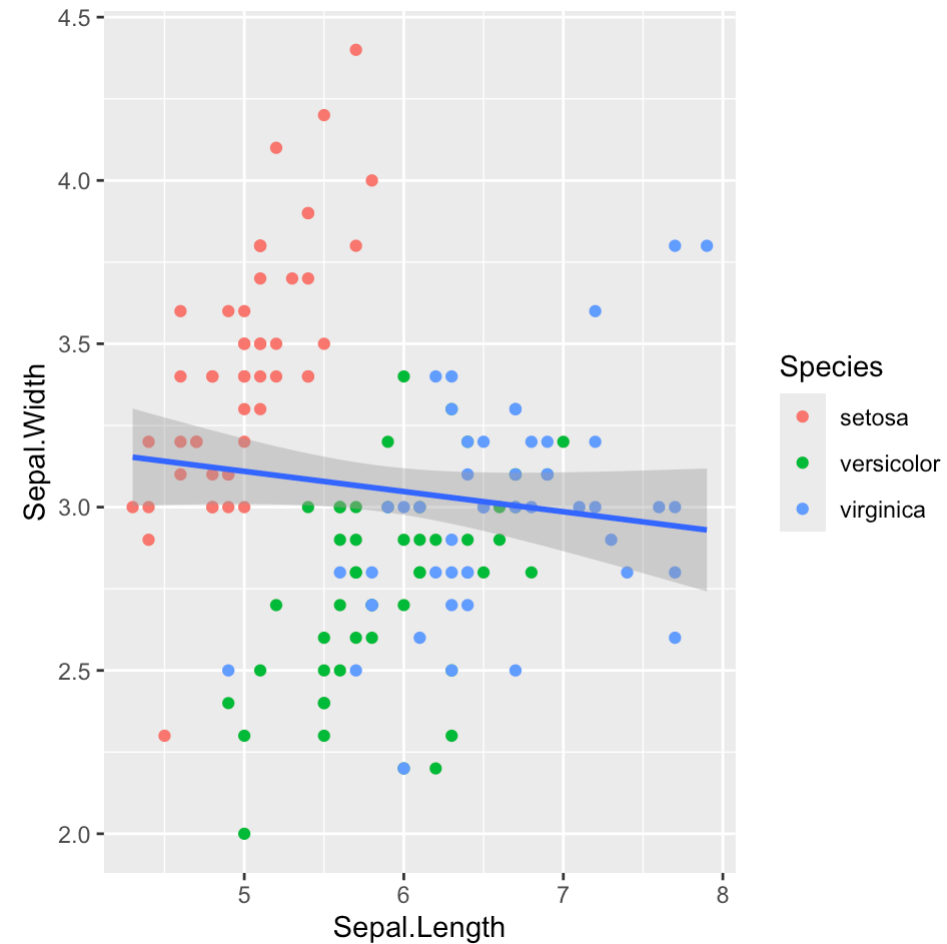# Let's plot using **ggplot2** - *Geometry*

```
1  ggplot(data = iris,
2         mapping = aes(x = Sepal.Length,
3                       y = Sepal.Width)) +
4     geom_point(mapping = aes(col = Species)
```
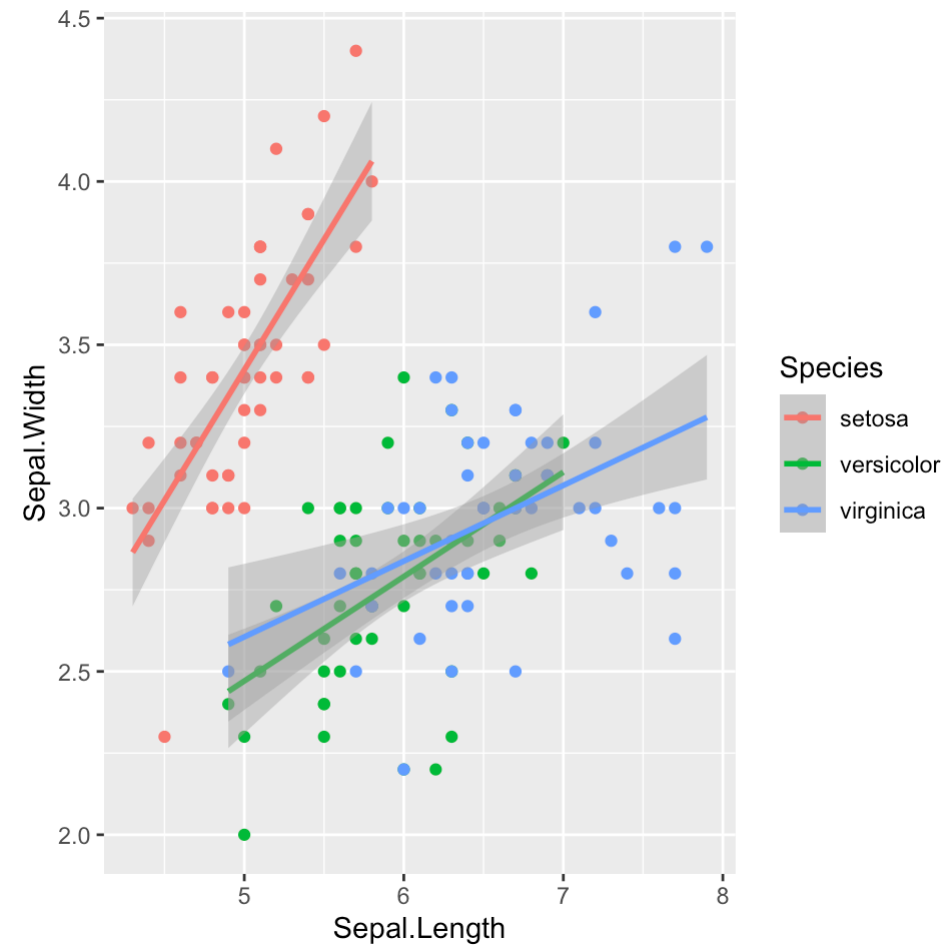


Stockholm
University

# Let's plot using **`ggplot2`** - *Statistics*

```
1  ggplot(data = iris,
2        mapping = aes(x = Sepal.Length,
3                      y = Sepal.Width)) +
4    geom_point(mapping = aes(col = Species)
5    stat_smooth(method = "lm")
```

# Let's plot using `ggplot2` - *Statistics*

```
1  ggplot(data = iris,
2         mapping = aes(x = Sepal.Length,
3                       y = Sepal.Width,
4                       col = Species)) +
5    geom_point() +
6    stat_smooth(method = "lm")
```
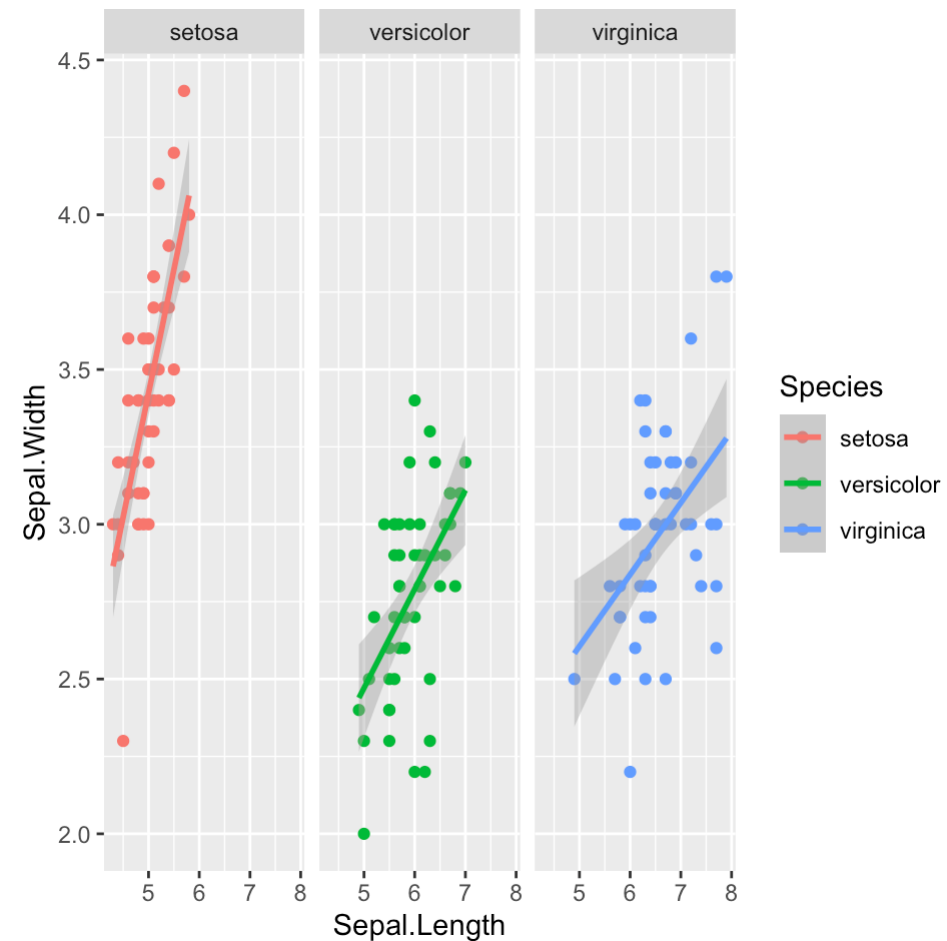
# Let's plot using **ggplot2** - *Facets*

```
1  ggplot(data = iris,
2         mapping = aes(x = Sepal.Length,
3                       y = Sepal.Width,
4                       col = Species)) +
5    geom_point() +
6    stat_smooth(method = "lm") +
7    facet_wrap(~Species)
```
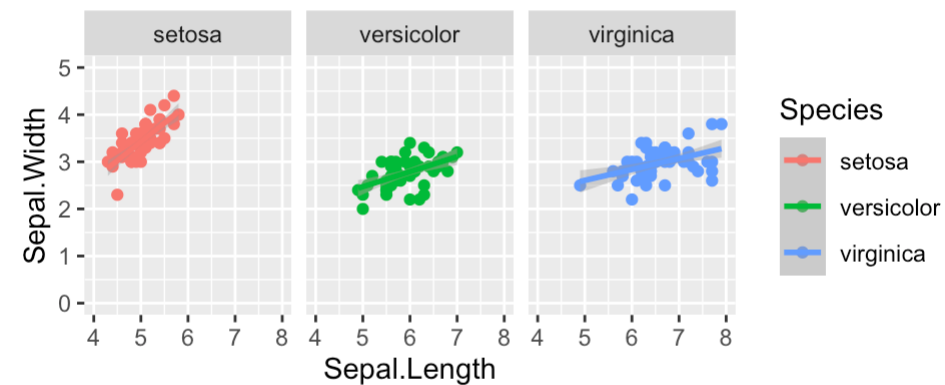


Stockholm
University

# Let's plot using ggplot2 - *Coordinates*

```
1  ggplot(data = iris,
2         mapping = aes(x = Sepal.Length,
3                       y = Sepal.Width,
4                       col = Species)) +
5    geom_point() +
6    stat_smooth(method = "lm") +
7    facet_wrap(~Species)+
8    coord_fixed(xlim = c(4,8), ylim = c(0,5)
```
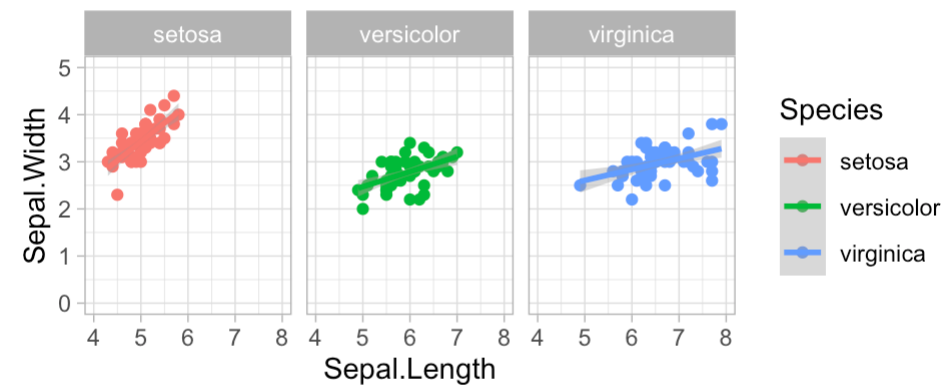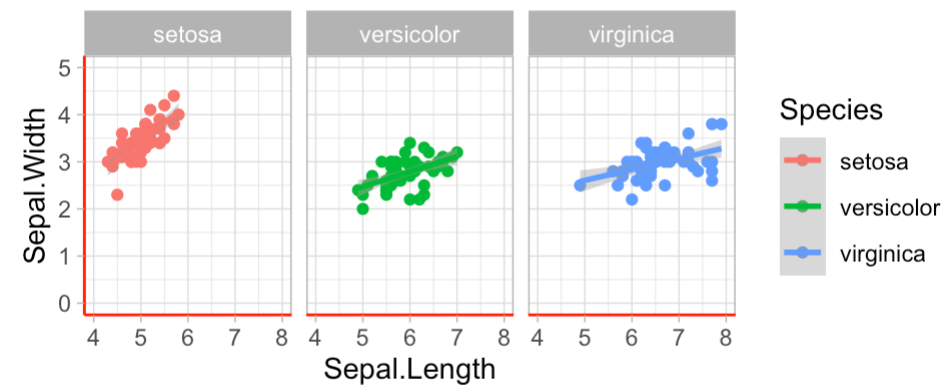
# Let's plot using **ggplot2** - *Themes*

```
1  ggplot(data = iris,
2        mapping = aes(x = Sepal.Length,
3                      y = Sepal.Width,
4                      col = Species)) +
5  geom_point() +
6  stat_smooth(method = "lm") +
7  facet_wrap(~Species) +
8  coord_fixed(xlim = c(4,8), ylim = c(0,5
9  theme_light()
```

# Let's plot using `ggplot2` - *Themes*

```
1  ggplot(data = iris,
2          mapping = aes(x = Sepal.Length,
3                        y = Sepal.Width,
4                        col = Species)) +
5     geom_point() +
6     stat_smooth(method = "lm") +
7     facet_wrap(~Species) +
8     coord_fixed(xlim = c(4,8), ylim = c(0,5
9     theme_light() +
10    theme(axis.line = element_line(color =
```
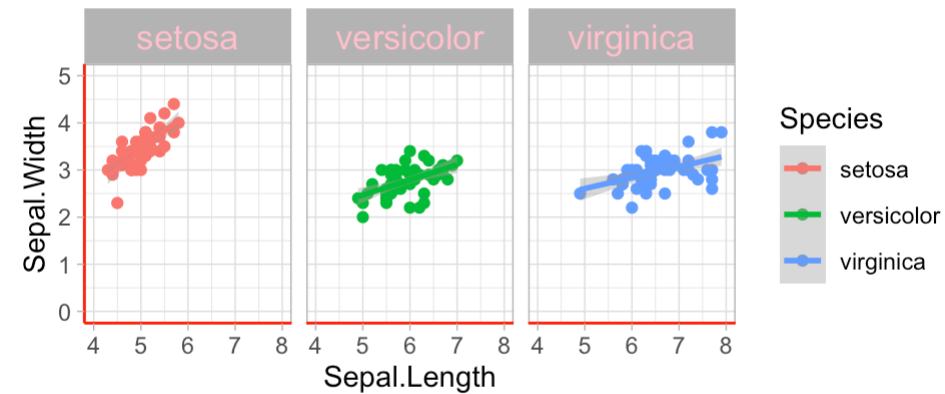
# Let's plot using **ggplot2** - *Themes*

```
1  ggplot(data = iris,
2         mapping = aes(x = Sepal.Length,
3                       y = Sepal.Width,
4                       col = Species)) +
5    geom_point() +
6    stat_smooth(method = "lm") +
7    facet_wrap(~Species) +
8    coord_fixed(xlim = c(4,8), ylim = c(0,5
9    theme_light() +
10   theme(axis.line = element_line(color =
11         strip.text = element_text(size =
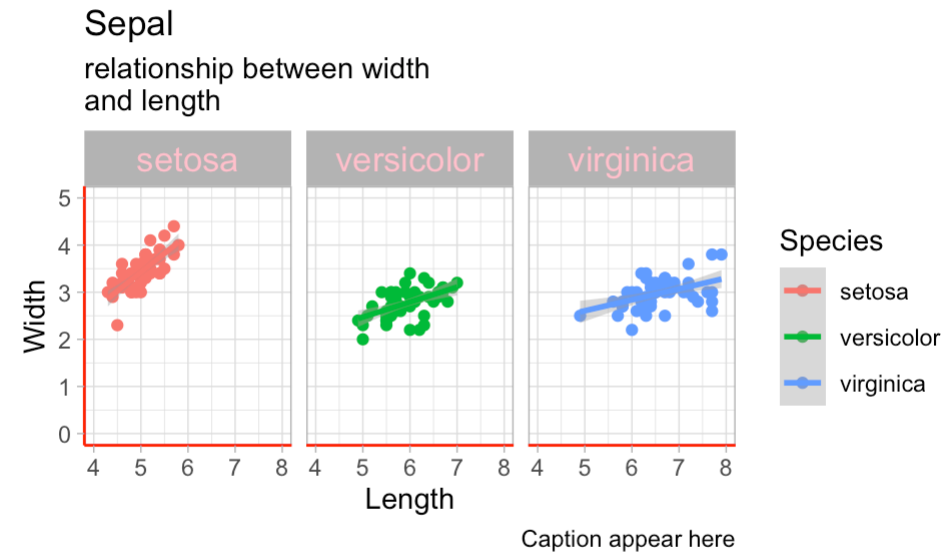```

# Let's plot using `ggplot2` - *Extra tips*

```
1  ggplot(data = iris,
2         mapping = aes(x = Sepal.Length,
3                       y = Sepal.Width,
4                       col = Species)) +
5    geom_point() +
6    stat_smooth(method = "lm") +
7    facet_wrap(~Species)+
8    coord_fixed(xlim = c(4,8), ylim = c(0,5)
9    theme_light() +
10   theme(axis.line = element_line(color =
11         strip.text = element_text(size =
12   labs(title = "Sepal", x = "Length" , y =
```

# Let's plot using `ggplot2` - *Extra tips*

```
1   ggplot(data = iris,
2          mapping = aes(x = Sepal.Length,
3                        y = Sepal.Width,
4                        col = Species)) +
5     geom_point() +
6     stat_smooth(method = "lm") +
7     facet_wrap(~Species)+
8     coord_fixed(xlim = c(4,8), ylim = c(0,5)
9     theme_light() +
10    theme(axis.line = element_line(color =
11          strip.text = element_text(size =
12    labs(title = "Sepal", x = "Length" , y =
13    scale_color_manual(values = c("forestgre
```



Sepal
relationship between width and length

Caption appear here