

Universidade Federal do Rio Grande do Norte
Instituto Metr pole Digital
IMD0601 - Bioestat stica

Visualiza  o dos dados em R

Prof. Dr. Tetsu Sakamoto
Instituto Metr pole Digital - UFRN
Sala A224, ramal 182
Email: tetsu@imd.ufrn.br



Baixe a aula (e os arquivos)

- Para aqueles que não clonaram o repositório:

```
> git clone https://github.com/tetsufmbio/IMD0601.git
```

- Para aqueles que já tem o repositório local:

```
> cd /path/to/IMD0601
```

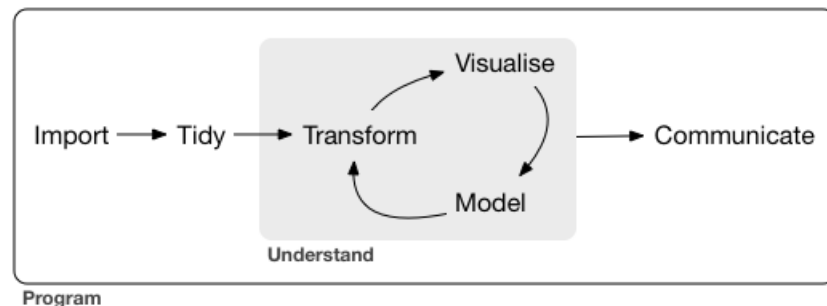
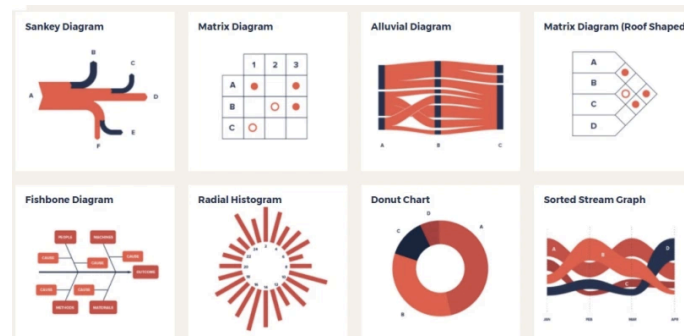
```
> git pull
```

Visualização dos dados

Cientistas de dados;

Combina as áreas de:

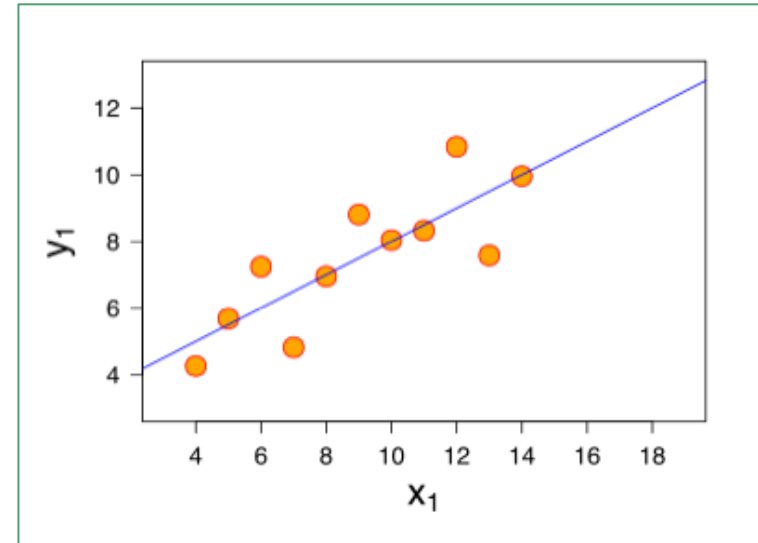
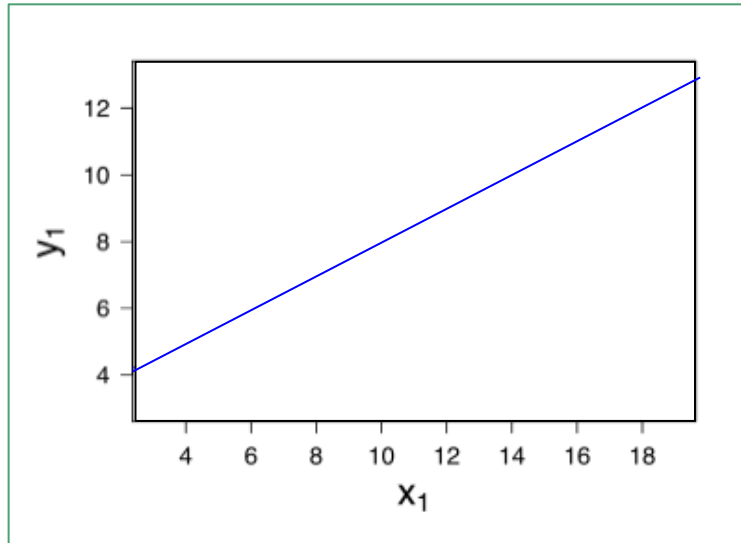
- **Estatística:** análise gráfica dos dados (Representação e interpretação dos dados);
- **Design:** Princípios de design (gráficos atrativos e que promove o melhor entendimento e comunicação);



<https://r4ds.had.co.nz/introduction.html>

Visualização dos dados

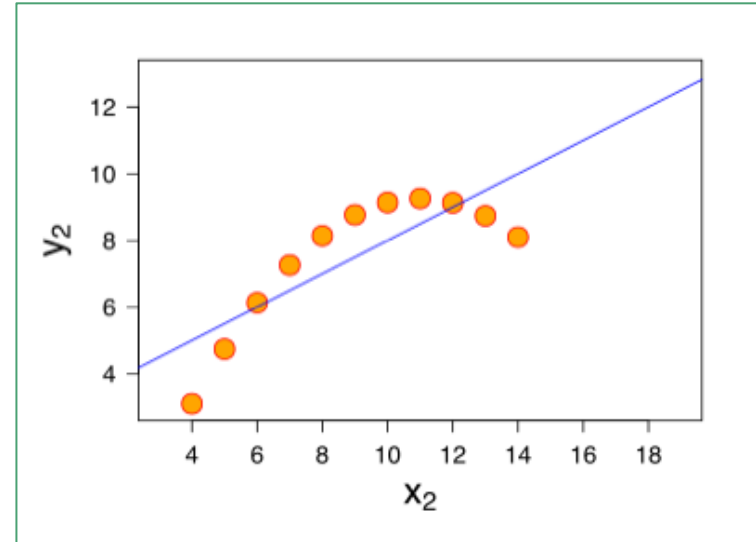
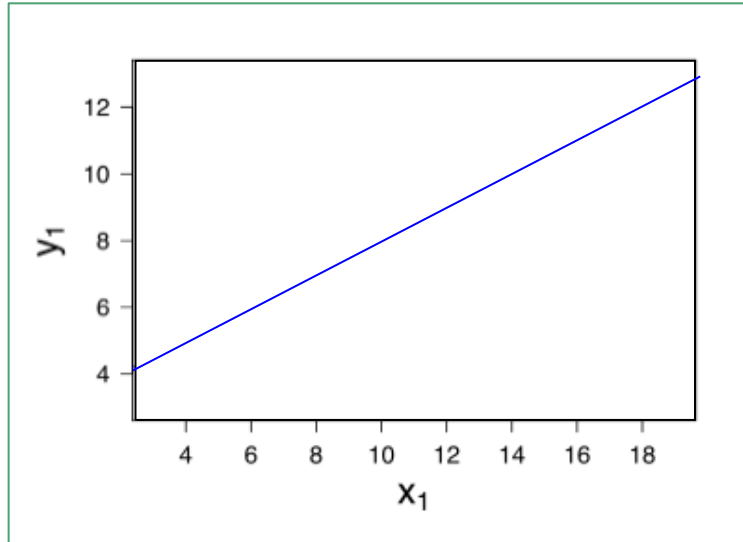
Quarteto de Anscombe



Anscombe, FJ. Graphs in Statistical Analysis. American Statistician, 1973.

Visualização dos dados

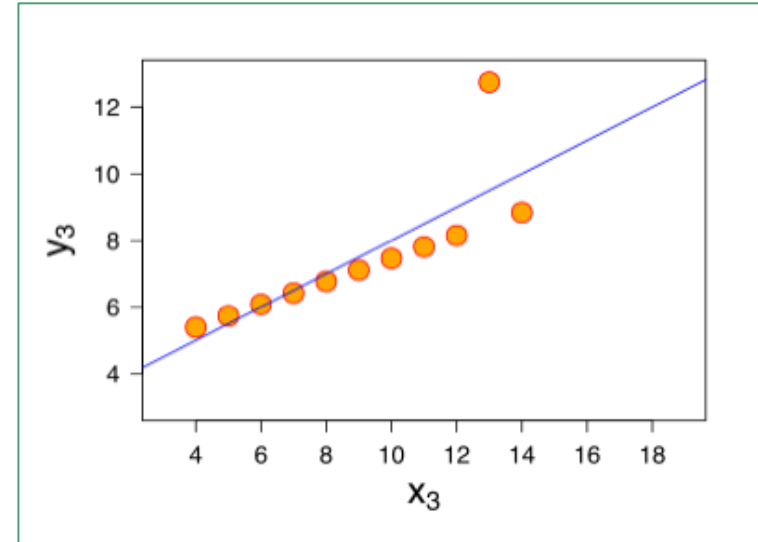
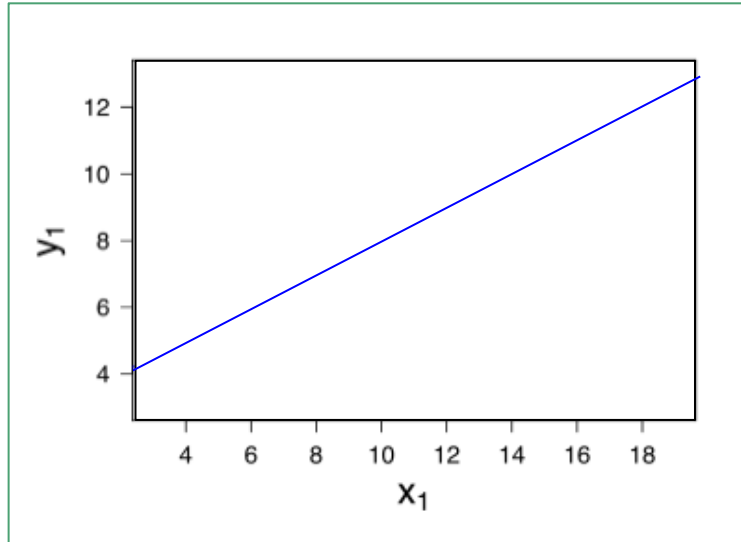
Quarteto de Anscombe



Anscombe, FJ. Graphs in Statistical Analysis. American Statistician, 1973.

Visualização dos dados

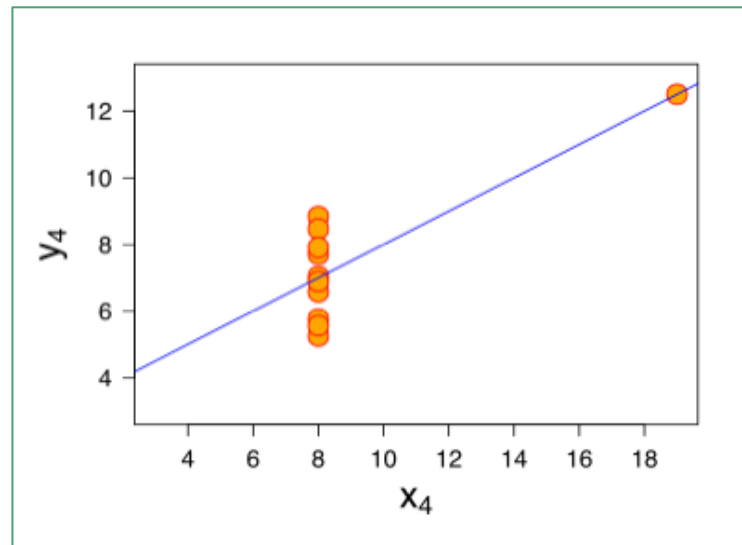
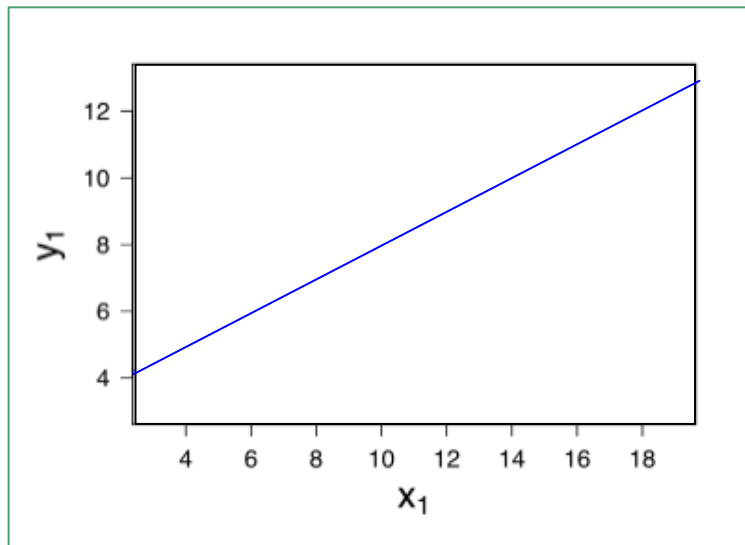
Quarteto de Anscombe



Anscombe, FJ. Graphs in Statistical Analysis. American Statistician, 1973.

Visualização dos dados

Quarteto de Anscombe



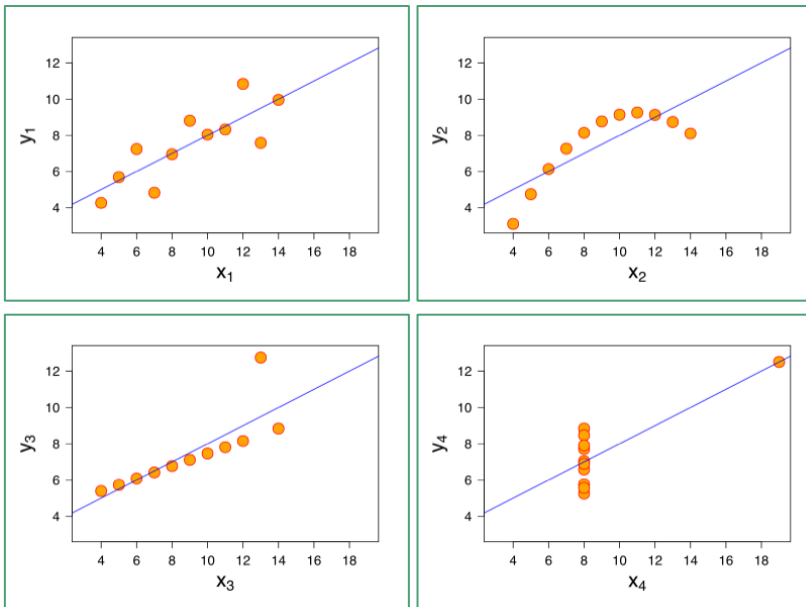
Anscombe, FJ. Graphs in Statistical Analysis. American Statistician, 1973.

Visualização dos dados

Quarteto de Anscombe



Propriedade	Valor	Precisão
Média de x	9	exato
Variância de x	11	exato
Média de y	7,50	até 2 casas decimais
Variância de y	4,125	$\pm 0,003$
Correlação entre x e y	0,816	até 3 casas decimais
Reta de regressão linear	$y = 3,00 + 0,500x$	até 2 e 3 casas decimais, respectivamente
Coefficiente de determinação da regressão linear: R^2	0,67	até 2 casas decimais



Anscombe, F.J. Graphs in Statistical Analysis. American Statistician, 1973.

Visualização dos dados

Gráficos exploratórios

- Gerado facilmente;
- Dados pesados;
- Para uma audiência específica (você e seus colegas);
- Análise gráfica dos dados;
- Servem para analisar e confirmar.

R Base Graphics

Gráficos explanatórios

- Laboriosos;
- Específicos para determinados dados;
- Para uma audiência ampla (publicação ou apresentação);
- Parte comunicativa do processo.
- Servem para informar e persuadir.

Ggplot2, ggvis
lattice

R Base Graphics

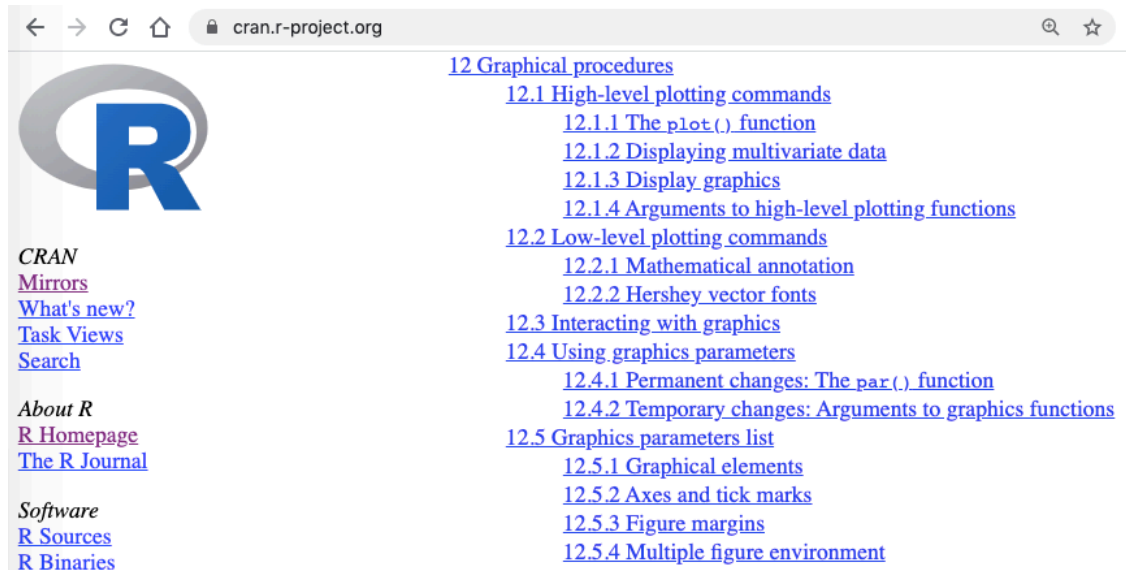
R Base Graphics: An Idiot's Guide

PROGRAMMING, R, UDACITY

SWIRL – R PROGRAMMING – LESSON 15 – BASE GRAPHICS

HOW-TO & USEFUL STUFF

Base Graphics in R: A Detailed Idiot's Guide

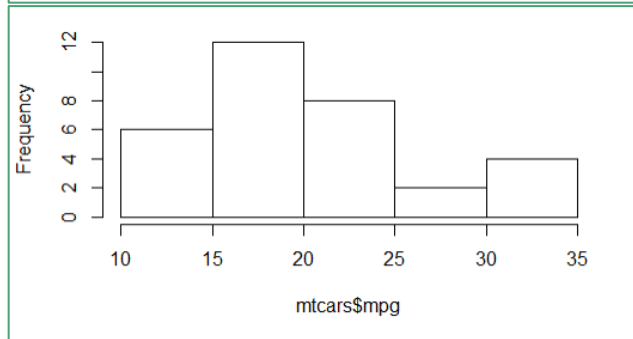
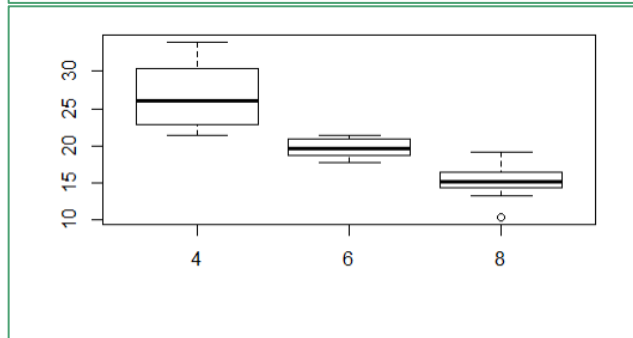
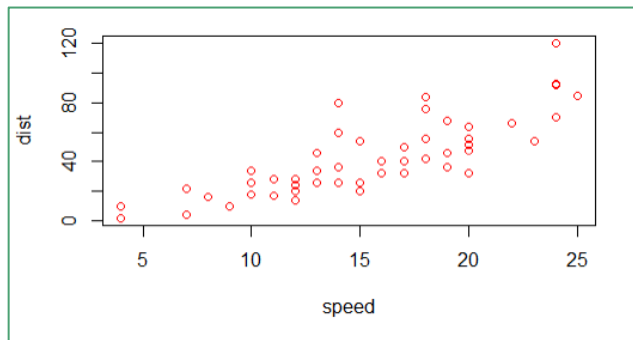


R Base Graphics

```
plot(x = cars$speed, y = cars$dist,  
     xlab = "Speed", ylab = "Stopping  
     Distance", col = 2)
```

```
boxplot(formula = mpg~cyl, data =  
mtcars)
```

```
hist(mtcars$mpg)
```





Overview

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.

Installation

```
# The easiest way to
install.packages("tidyverse")
```

```
# Alternatively, install
install.packages("ggplot2")
```

```
# Or the development
install.packages("devtools::install_github('tidyverse/ggplot2')")
```

```
# install.packages('devtools::install_github('tidyverse/ggplot2')')
```

Data Visualization with ggplot2 : : CHEAT SHEET

Basics

ggplot2 is based on the **grammar of graphics**, the idea that you can build every graph from the same components: a **data set**, a **coordinate system**, and **geoms**—visual marks that represent data points.

To display values, map variables in the data to visual properties of the geom (**aesthetics**) like **size**, **color**, and **x** and **y** locations.

Complete the template below to build a graph.

```
ggplot(data = <DATA>) +
  <COORDINATE_FUNCTION> +
  <GEOM_FUNCTION> +
  <SCALE_FUNCTION> +
  <THEME_FUNCTION>
```

ggplot(data = mpg, aes(x = cty, y = hwy)) Begins a plot that you finish by adding layers. Add one geom function per layer.

aesthetic mappings (aes()) Create a complete plot with given data, geom, and mappings. Supplies many useful defaults.

last plot() Returns the last plot

Geoms

Use a geom function to represent data points, use the geom's aesthetic properties to represent variables. Each function returns a layer.

GRAPHICAL PRIMITIVES
a <- ggplot(economics, aes(date, unemploy))
b <- ggplot(seals, aes(x = long, y = lat))
a + geom_blank() (Useful for expanding limits)
b + geom_curve(aes(yend = lat + 1, xend = long - 1, curvature = 1), x = xend, y = yend, alpha, angle, color, curvature, linetype, size)
a + geom_path(linetype = "dotted", linejoin = "round", linetype = 1)
x, y, alpha, color, group, linetype, size
a + geom_polygon(aes(group = group))
x, y, alpha, color, fill, group, linetype, size
b + geom_rect(aes(xmin = long, ymin = lat, xmax = long + 1, ymax = lat + 1), x = xmin, y = ymin, y = ymax, alpha, color, fill, group, linetype, size)
a + geom_ribbon(aes(ymin = unemploy - 900, ymax = unemploy + 900), x = x, y = y, alpha, color, fill, group, linetype, size)

LINE SEGMENTS
common aesthetics: x, y, alpha, color, linetype, size
b + geom_abline(aes(intercept = 0, slope = 1))
b + geom_hline(aes(intercept = lat))
b + geom_vline(aes(intercept = long))
b + geom_segment(aes(yend = lat + 1, xend = long + 1))
b + geom_spoke(aes(angle = 1.1155, radius = 1))

ONE VARIABLE continuous
c <- ggplot(mpg, aes(hwy))
c + geom_area(aes(fill = "hwy"))
c + geom_density(kernel = "gaussian")
c + geom_dotplot()
c + geom_freqpoly(x, y, alpha, color, group, linetype, size, weight)

TWO VARIABLES
continuous x, continuous y
e <- ggplot(mpg, aes(cty, hwy))
e + geom_label(aes(label = cty, nudges_x = 1, nudges_y = 1, check_overlap = TRUE), x, y, label, alpha, angle, color, family, fontface, hjust, lineheight, size, vjust)
e + geom_jitter(height = 2, width = 2)
e + geom_point(x, y, alpha, color, fill, shape, size, stroke)
e + geom_quantile(x, y, alpha, color, group, linetype, size, weight)
e + geom_rug(sides = "bt", x, y, alpha, color, linetype, size)
e + geom_smooth(method = lm, x, y, alpha, color, fill, group, linetype, size, weight)
e + geom_text(aes(label = cty, nudges_x = 1, nudges_y = 1, check_overlap = TRUE), x, y, label, alpha, angle, color, family, fontface, hjust, lineheight, size, vjust)

discrete x, continuous y
f <- ggplot(mpg, aes(class, hwy))
f + geom_col(x, y, alpha, color, fill, group, linetype, size, weight)
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, linetype, size, weight)
f + geom_violin(scale = "area", x, y, alpha, color, fill, group, linetype, size, weight)

discrete x, discrete y
g <- ggplot(diamonds, aes(cut, color))

continuous bivariate distribution
h <- ggplot(diamonds, aes(carat, price))
h + geom_bin2d(binwidth = c(0.25, 500))
h + geom_density2d()
h + geom_hex()
continuous function
i <- ggplot(economics, aes(date, unemploy))
i + geom_area(x, y, alpha, color, fill, linetype, size)
i + geom_line(x, y, alpha, color, group, linetype, size)
i + geom_step(direction = "hv", x, y, alpha, color, group, linetype, size)

visualizing error
d1 <- data.frame(grp = c("A", "B"), fit = 4.5, se = 1.2)
j <- ggplot(d1, aes(grp, fit, ymin = fit - se, ymax = fit + se))
j + geom_crossbar(latten = 2)
j + geom_errorbar(x, ymax, ymin, alpha, color, fill, group, linetype, size)
j + geom_linerange(x, ymin, ymax, alpha, color, group, linetype, size)
j + geom_pointrange(x, y, ymin, ymax, alpha, color, fill, group, linetype, shape, size)

maps
data = data.frame(murder = USArrests\$Murder, state = tolower(rownames(USArrests)))
map <- map_data("state")
k <- ggplot(data, aes(fill = murder))

Links

Download from CRAN at

<https://cloud.r-project.org/package=ggplot2>

Browse source code at

<https://github.com/tidyverse/ggplot2/>

Report a bug at

<https://github.com/tidyverse/ggplot2/issues>

Learn more at

<https://r4ds.had.co.nz/data-visualisation.html>

Extensions at

<https://exts.ggplot2.tidyverse.org/gallery/>

ggplot2

Hadley Wickham

“The Grammar of Graphics”

Adiciona camadas nos gráficos para
melhor visualização dos dados;

```
library(ggplot2)
```



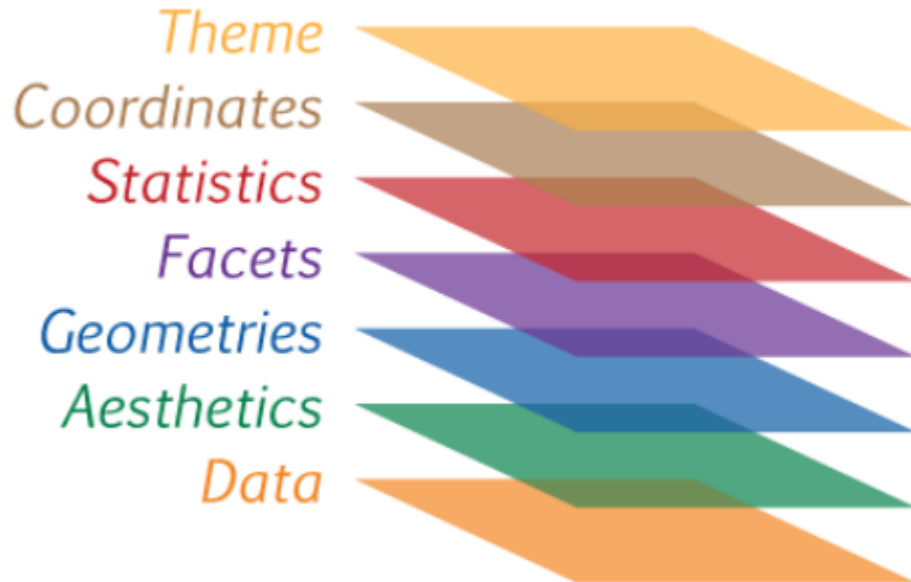
ggplot2

Data

Dados (tabela) onde se encontram as variáveis a serem representados graficamente.

```
ggplot(data=mtcars)
```

```
ggplot(mtcars)
```



ggplot2

Aesthetics

Permite especificar as variáveis que queremos utilizar na representação gráfica.

```
ggplot(mtcars, aes(x=mpg,  
y=wt))
```



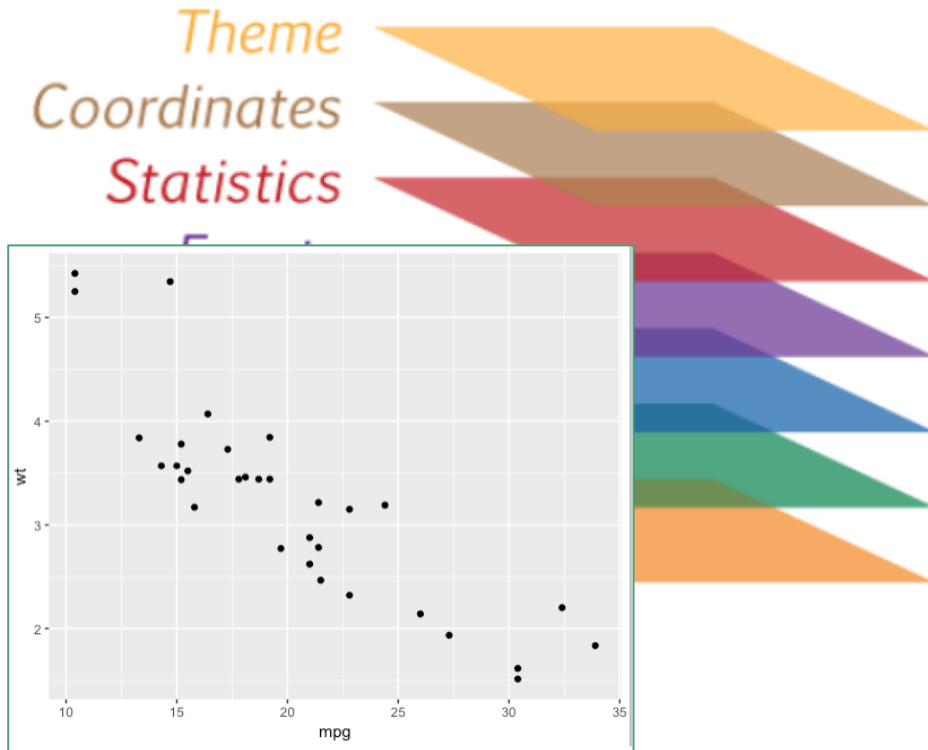
ggplot2

Geometries

Camada que indica a forma como os dados devem ser apresentados no gráfico.

```
ggplot(mtcars, aes(x=mpg,  
y=wt)) + geom_point()
```

```
g <- ggplot(mtcars, aes(x=mpg,  
y=wt)) + geom_point()
```



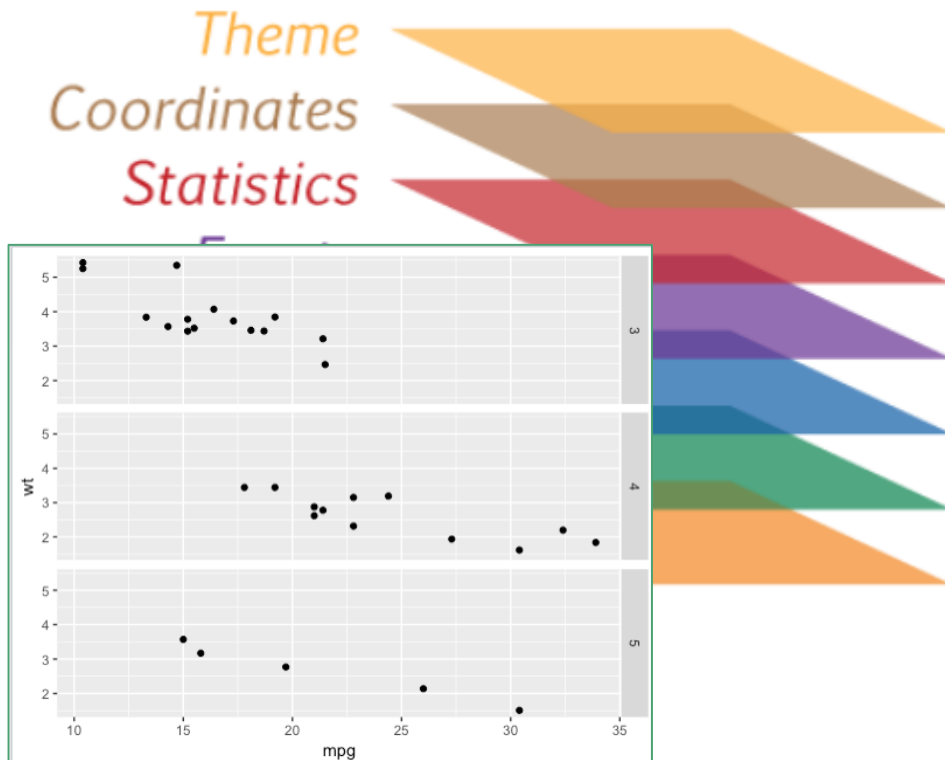
ggplot2

Facets

Permite colocar múltiplos gráficos em um canvas.

```
g + facet_grid(gear~.)
```

```
g <- g + facet_grid(gear~.)
```



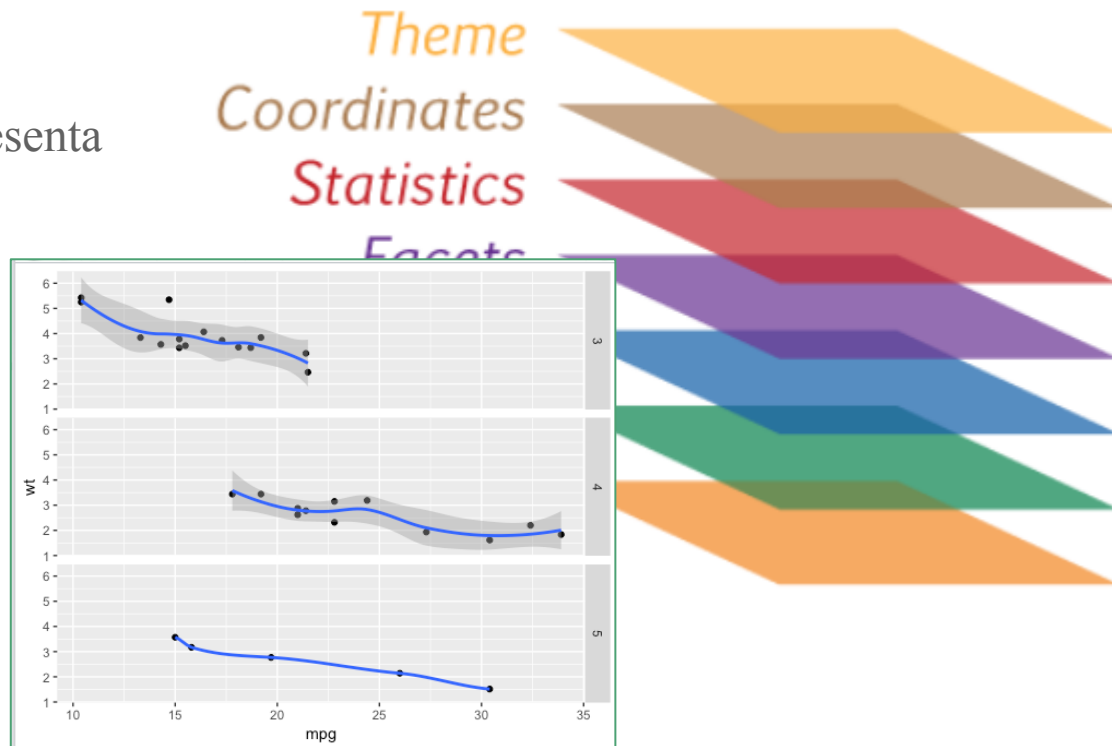
ggplot2

Statistics

Adiciona uma camada que representa uma análise estatística.

```
g + stat_smooth()
```

```
g <- g + stat_smooth()
```



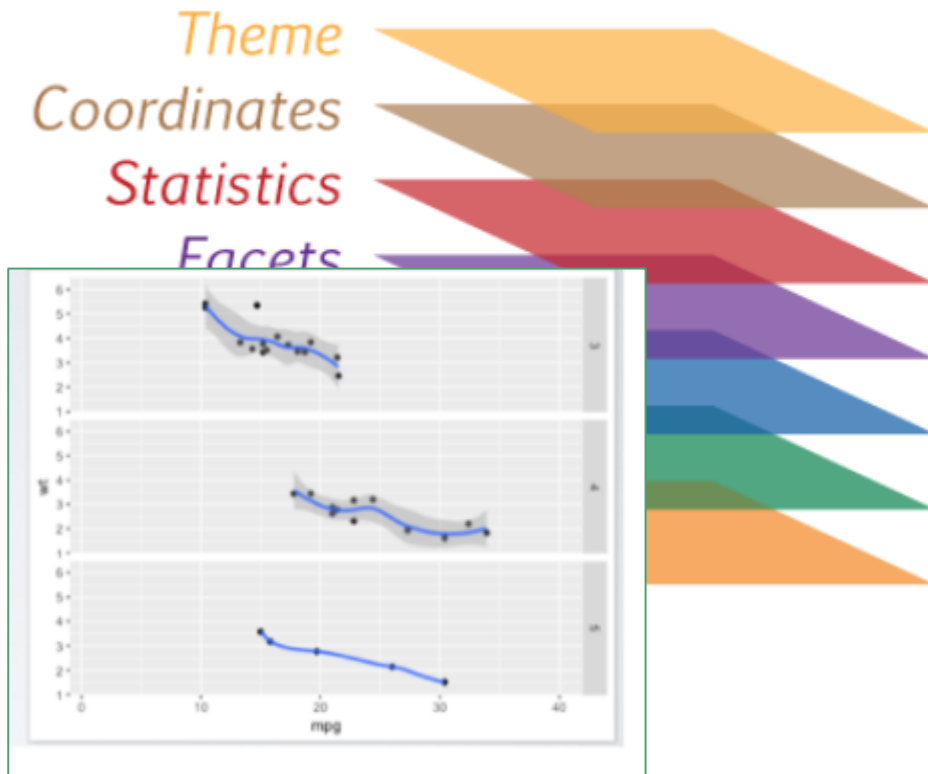
ggplot2

Coordinates

Camada que controla como as posições devem ser mapeadas no gráfico (limites do eixo x e y).

```
g + coord_cartesian(xlim =  
c(1, 40))
```

```
g <- g + coord_cartesian(xlim  
= c(1, 40))
```



ggplot2

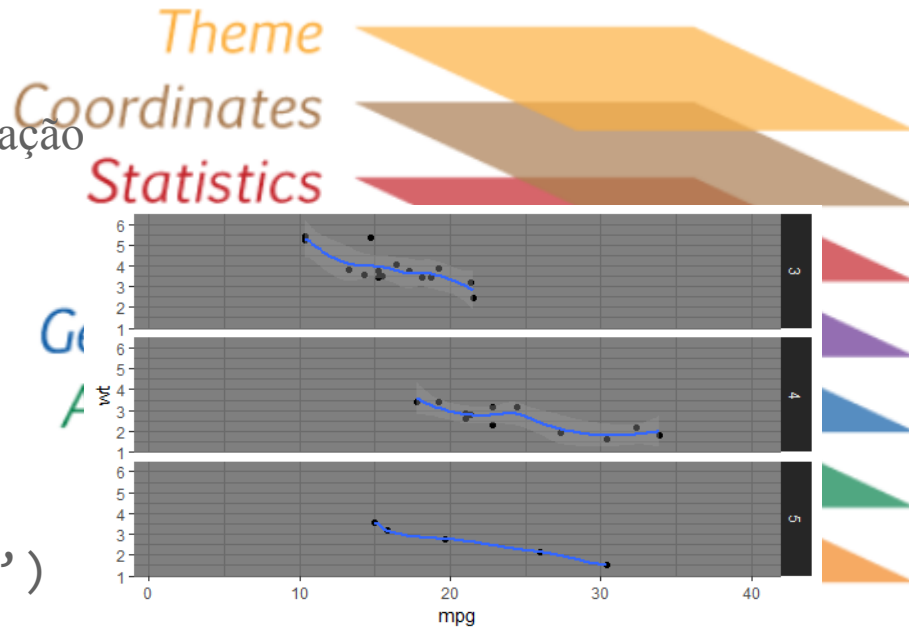
Themes

Camada que permite enriquecer a apresentação do gráfico (rótulos, fonte, cor, etc).

```
g + theme_dark()
```

```
g <- g + theme_dark()
```

```
install.packages('ggThemeAssist')
```



Iris dataset

```
data(iris)
```

```
str(iris)
```



Iris Versicolor

Iris Setosa

Iris Virginica

```
'data.frame':  150 obs. of  5 variables:
```

```
$ Sepal.Length: num  5.1 4.9 4.7 4.6 5 5.4 4.6 5 4.4 4.9 ...
```

```
$ Sepal.Width : num  3.5 3 3.2 3.1 3.6 3.9 3.4 3.4 2.9 3.1 ...
```

```
$ Petal.Length: num  1.4 1.4 1.3 1.5 1.4 1.7 1.4 1.5 1.4 1.5 ...
```

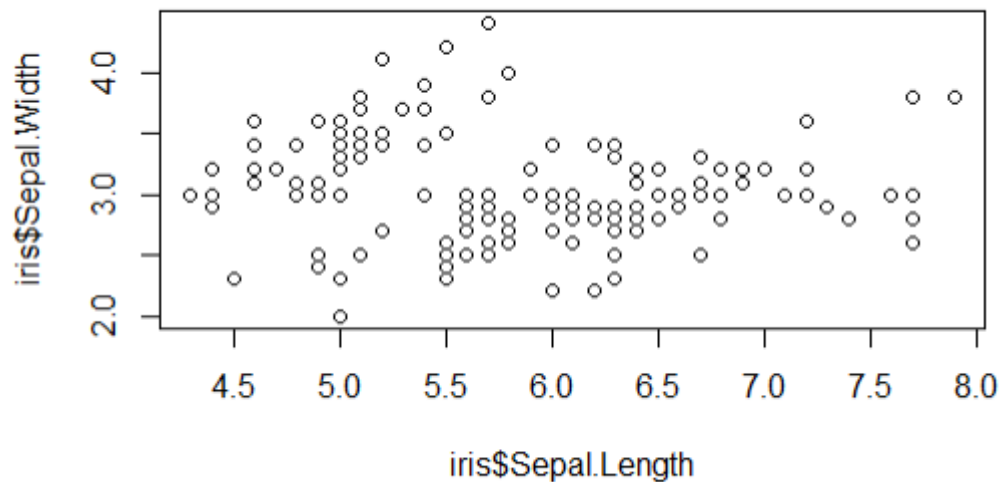
```
$ Petal.Width : num  0.2 0.2 0.2 0.2 0.2 0.4 0.3 0.2 0.2 0.1 ...
```

```
$ Species      : Factor w/ 3 levels "setosa","versicolor",...: 1 1
```

```
1 1 1 1 1 1 1 1 ...
```

Base plot

Sepal length X Sepal width



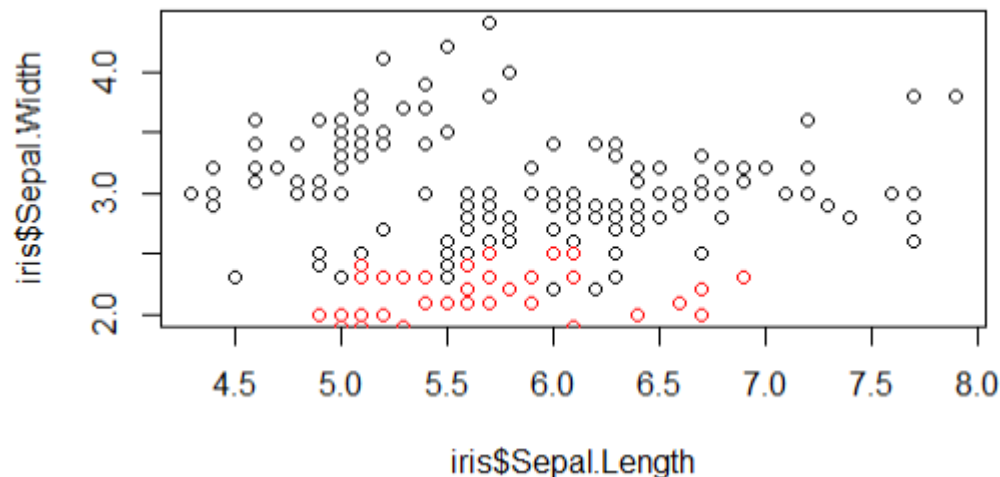
```
> plot(iris$Sepal.Length, iris$Sepal.Width)
```

Base plot

Como adicionar os dados de Petal length e Petal width?

```
> plot(iris$Sepal.Length, iris$Sepal.Width)
```

```
> points(iris$Petal.Length, iris$Petal.Width, col = "red")
```



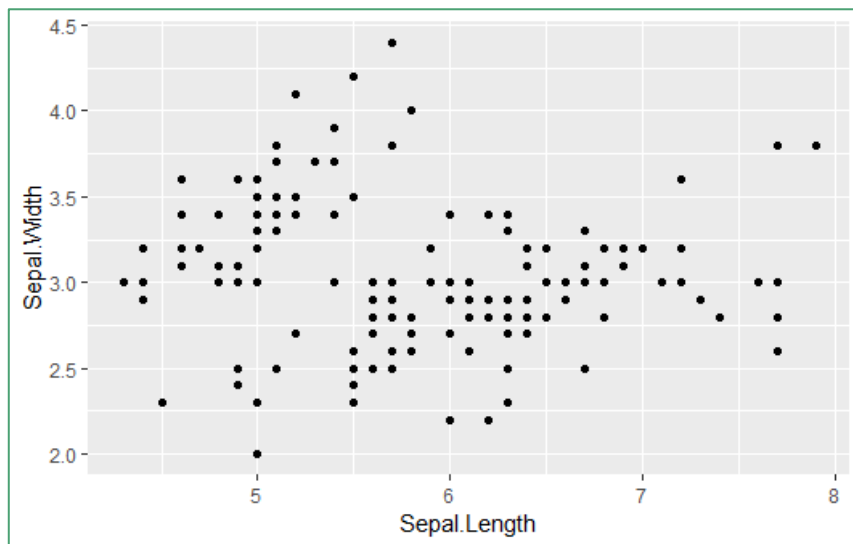
Limitações:

- O gráfico não é redesenhado;
- O gráfico gerado é uma imagem;
- A legenda é adicionada manualmente;
- Não há uma framework unificado para gerar diferentes tipos de gráficos.

ggplot2

Sepal length X Sepal width

```
> ggplot(iris, aes(x = Sepal.Length, y = Sepal.Width))  
  
+ geom_point()
```



ggplot2

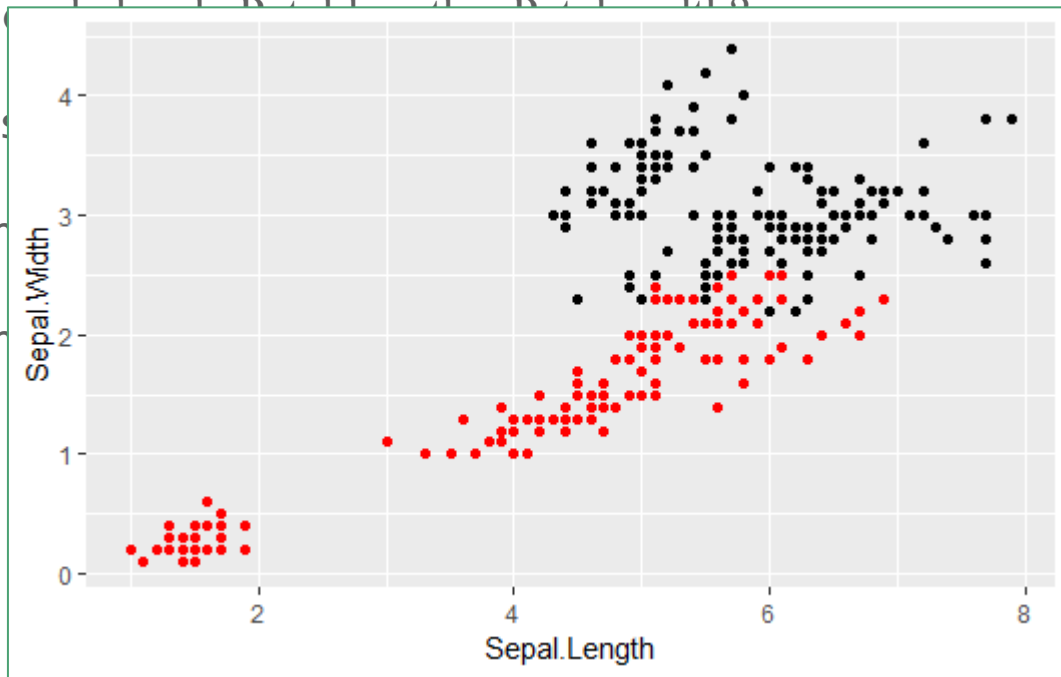
Como adicionar os dados de Petal length e Petal width?

```
> ggplot(iris, aes(x = Sepal.Length, y = Sepal.Width)) +  
  geom_point() +  
  geom_point(aes(x = Petal.Length, y = Petal.Width),  
col = "red")
```

ggplot2

Como adicionar

```
> ggplot(iris) +  
  geom_point() +  
  geom_point(col = "red")
```



ggplot2

Exercício

Adicione uma coluna na tabela **iris** que corresponda a um identificador único de cada observação.

```
> iris$Flower <- 1:nrow(iris)
```

Crie uma tabela onde as variáveis Length e Width estejam cada uma em uma coluna, como abaixo.

	Species	Flower	part	Length	Width
1	setosa	1	Petal	1.4	0.2
2	setosa	1	Sepal	5.1	3.5
...					

ggplot2

Exercício

```
> library(tidyr)

> iris$Flower <- 1:nrow(iris)

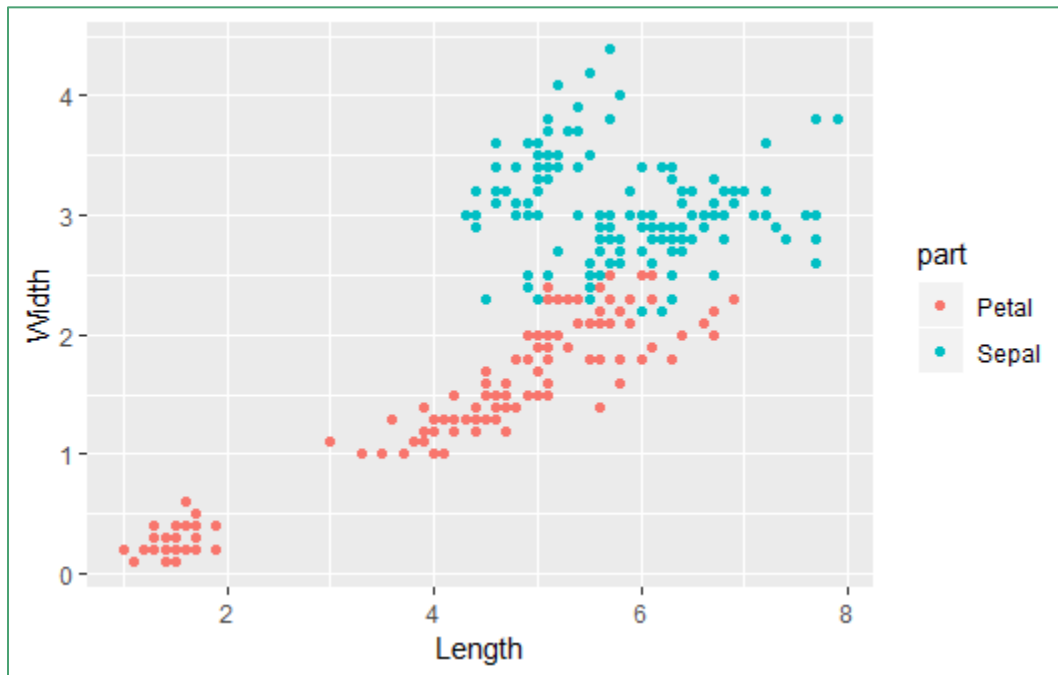
> iris.wide <- gather(iris, part_measure, val, -Species, -Flower )

> iris.wide <- separate(iris.wide, part_measure,
c("part", "measure"))

> iris.wide <- spread(iris.wide, measure, val)
```

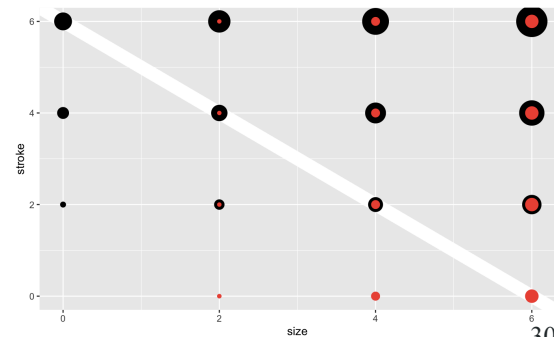
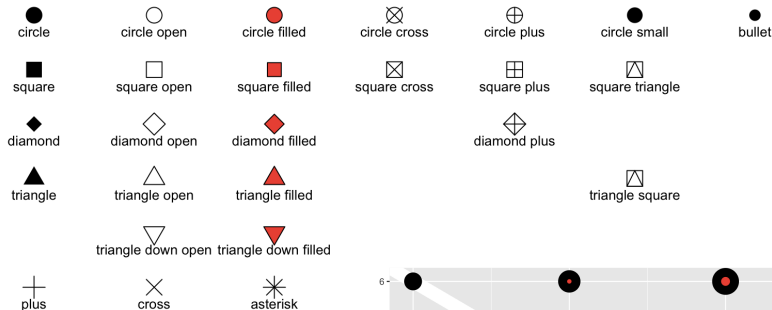
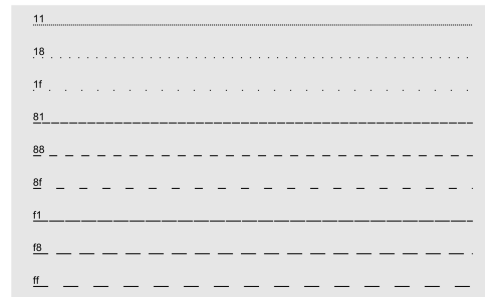
ggplot2

```
ggplot(iris.wide, aes(Length, Width, col = part)) + geom_point()
```



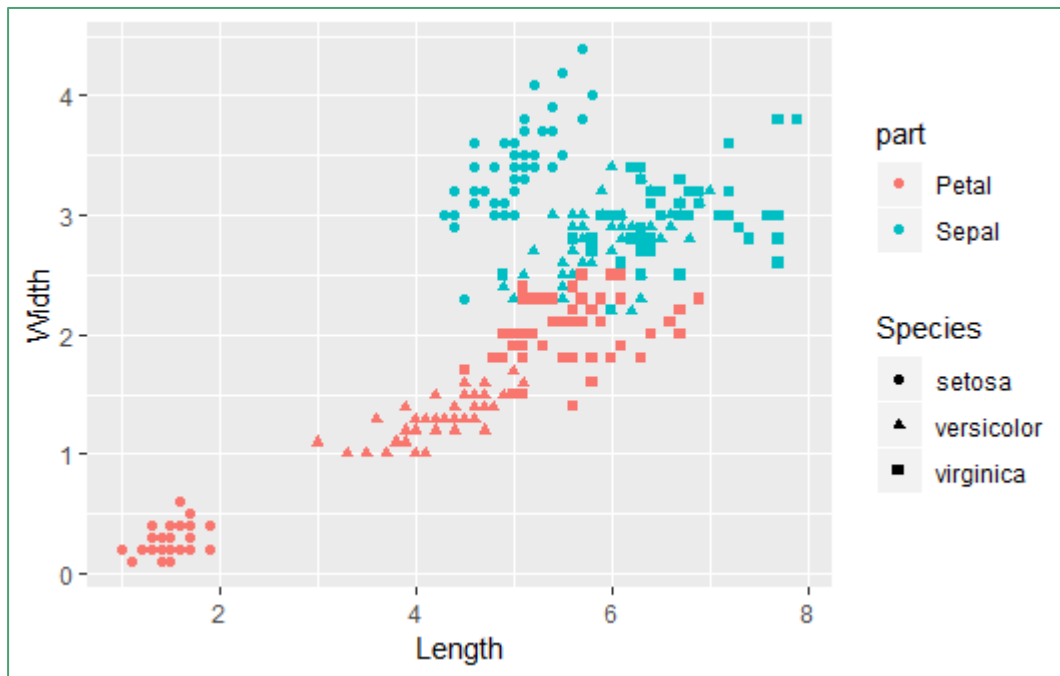
Parâmetros típicos da estética

- **x** → posição no eixo x;
- **y** → posição no eixo y;
- **col** → cor dos pontos ou formas;
- **fill** → cor a ser preenchida;
- **size** → diâmetro do ponto, largura da linha;
- **alpha** → transparência;
- **linetype** → padrão de tracejamento;
- **labels** → texto no gráfico;
- **shape** → formas;



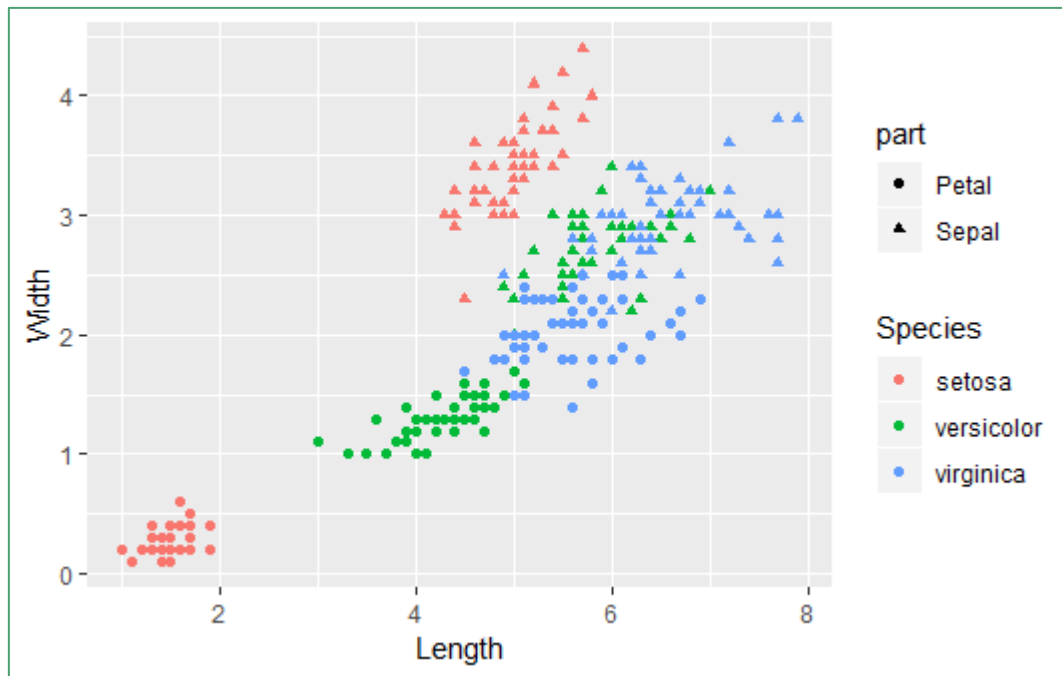
ggplot2

```
ggplot(iris.wide, aes(Length, Width, col = part, shape = Species))  
+ geom_point()
```



ggplot2

```
ggplot(iris.wide, aes(Length, Width, col = Species, shape = part))  
+ geom_point()
```



A sintaxe do ggplot

```
ggplot(data = <DATA>) +  
  <GEOM_FUNCTION>(  
    mapping = aes(<MAPPINGS>),  
    stat = <STAT>,  
    position = <POSITION>  
  ) +  
  <COORDINATE_FUNCTION> +  
  <FACET_FUNCTION>
```

Data	{variables of interest}				
Aesthetics	<i>x-axis</i> <i>y-axis</i>	<i>colour</i> <i>fill</i>	<i>size</i> <i>labels</i>	<i>alpha</i> <i>shape</i>	<i>line width</i> <i>line type</i>
Geometries	<i>point</i>	<i>line</i>	<i>histogram</i>	<i>bar</i>	<i>boxplot</i>
Themes	<i>non-data ink</i>				
Statistics	<i>binning</i>	<i>smoothing</i>	<i>descriptive</i>	<i>inferential</i>	
Coordinates	<i>cartesian</i>	<i>fixed</i>	<i>polar</i>	<i>limits</i>	
Facets	<i>columns</i>	<i>rows</i>			

Referências

- Aula baseada no curso “**Data Visualization with ggplot2 (Part 1)**” de Rick Scavetta: <https://www.datacamp.com/courses/data-visualization-with-ggplot2-1>
- Ciência de Dados com R – IBPAD: <https://www.ibpad.com.br/o-que-fazemos/publicacoes/introducao-ciencia-de-dados-com-r/>
- Tidyverse: <https://ggplot2.tidyverse.org/index.html> (link para cheat sheet!)
- R for data Science: <https://r4ds.had.co.nz/>
- <https://skillgaze.com/2017/10/31/understanding-different-visualization-layers-of-ggplot/>