

# Probabilidade e Estatística com R

Fernando Náufel

(versão de 20/11/2021)

---

## Sumário

---

<b>Apresentação</b>	<b>3</b>
Referências recomendadas	5
Exercício	5
<b>1 O Que É Estatística?</b>	<b>6</b>
1.1 Vídeo 1	6
1.2 Exercícios	6
1.3 Vídeo 2	9
1.4 Exercícios	9
<b>2 Introdução a R</b>	<b>10</b>
2.1 Vídeo 1	10
2.2 Vídeo 2	10
2.3 Exercícios	10
<b>3 Visualização com ggplot2</b>	<b>12</b>
3.1 Vídeo 1	12
3.2 Componentes de um gráfico ggplot2	12
3.3 Conjunto de dados	15
3.4 Gráficos de dispersão ( <i>scatter plots</i> )	21
3.5 Vídeo 2	38
3.6 Histogramas e cia.	38
3.7 Ogiva	43
3.8 Ramos e folhas	44
3.9 Exercícios	45
<b>4 Visualização com ggplot2 (continuação)</b>	<b>49</b>
4.1 Vídeo 1	49
4.2 <i>Boxplots</i>	49
4.3 Vídeo 2	59
4.4 Gráficos de barras e de colunas	59

4.5	Gráficos de linha e séries temporais . . . . .	71
4.6	Referências sobre visualização e R . . . . .	75
5	Medidas	76

---

## Apresentação

---

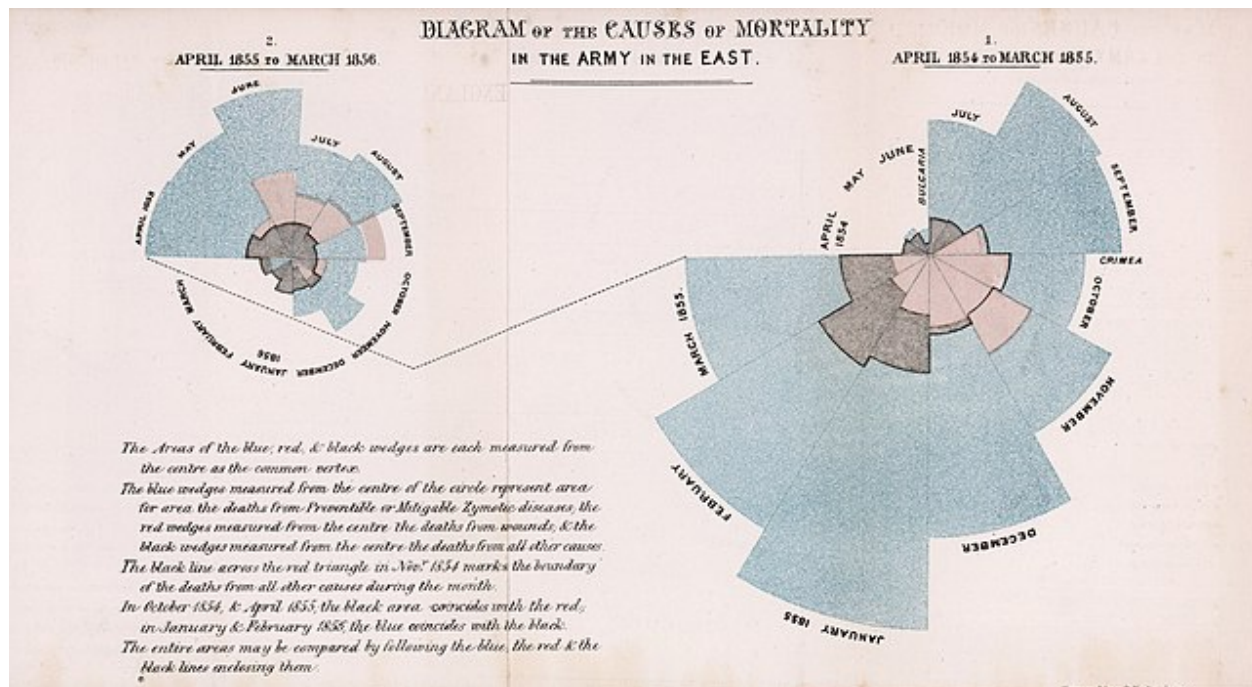
### Atenção



Este material ainda está em construção.

Pode haver mudanças a qualquer momento.

Verifique, no rodapé da página *web* ou na capa do arquivo pdf, a data desta versão.



Este livro/site foi iniciado em 2020, durante a pandemia de COVID-19, quando a Universidade Federal Fluminense (UFF) funcionou em regime de ensino remoto durante mais de um ano.

Para atender os alunos do curso de Probabilidade e Estatística do curso de graduação em Ciência da Computação da UFF, decidi gravar aulas em vídeo e disponibilizar os arquivos usados nelas. Foram esses arquivos que deram origem a este livro/site.

Este livro/site foi construído para pessoas que já saibam programar, embora não necessariamente em R.

Para tirar o máximo proveito deste material, você deve fazer o seguinte:

1. Assistir aos vídeos contidos em cada capítulo. A *playlist* completa está em <https://www.youtube.com/playlist?list=PL7SRLwLs7ocaV-Y1vrVU3W7mZnnS0qkWW>.
2. Instalar o R no seu computador ou abrir uma conta no RStudio Cloud, para poder usar o R *online*. Você encontra instruções para fazer isto no [capítulo de introdução a R](#).
3. Baixar, neste repositório do Github<sup>1</sup>, o código-fonte deste livro/site, para poder rodar e alterar os exemplos.
4. Seguir os *links* para outras fontes *online* que abordam assuntos que não são cobertos em detalhes neste curso.

<sup>1</sup><https://github.com/fnaufel/probestr>

5. Fazer os exercícios. Ao longo do tempo, acrescentarei *links* para vídeos explicando as soluções.



Se você estiver lendo este material na web, você pode clicar nos comandos e funções que aparecem nos blocos de código em R para abrir páginas da documentação sobre eles.

Se você preferir ler este livro em pdf, ou se quiser imprimi-lo, faça o *download* do arquivo aqui<sup>a</sup>.

<sup>a</sup><https://github.com/fnaufel/probesttr/blob/master/docs/probesttr.pdf>

---

## Referências recomendadas

---

### Em português

- Sillas Gonzaga, *Introdução a R para Visualização e Apresentação de Dados*, [http://sillasgonzaga.com/material/curso\\_visualizacao/index.html](http://sillasgonzaga.com/material/curso_visualizacao/index.html)
- Allan Vieira de Castro Quadros, *Introdução à Análise de Dados em R utilizando Tidyverse*, [https://allanvc.github.io/book\\_IADR-T/](https://allanvc.github.io/book_IADR-T/)
- Paulo Felipe de Oliveira, Saulo Guerra, Robert McDonnel, *Ciência de Dados com R – Introdução*, <https://cdr.ibpad.com.br/index.html>
- Curso R, *Ciência de Dados em R*, <https://livro.curso-r.com/>

---

### Em inglês

- Garrett Golemund, Hadley Wickham, *R for Data Science*, <https://r4ds.had.co.nz/>
- Chester Ismay, Albert Y. Kim, *A ModernDive into R and the Tidyverse*, <https://moderndive.com/>

---

## Exercício

1. Pesquise sobre a imagem do início deste capítulo. Ela foi criada em 1858 por Florence Nightingale.

## O Que É Estatística?

---

### 1.1

---

#### Vídeo 1

[https://youtu.be/6Q\\_XSoLCIpc](https://youtu.be/6Q_XSoLCIpc)

### 1.2

---

#### Exercícios

1. Você está interessado em estimar a altura de todos os homens da sua faculdade. Para isso, você decide medir as alturas de todos os homens da sua turma de Estatística.
  - Qual é a amostra?
  - Qual é a população?
2. Um instituto de pesquisa entrevista um grupo de 1000 pessoas, perguntando a cada uma se ela vai votar a favor do candidato  $A$  na próxima eleição. Dos entrevistados, 600 responderam que sim. A proporção 0,6 (ou 60%) é uma estatística ou um parâmetro?
3. Você vê alguma diferença entre as cinco situações abaixo? Quais das situações são equivalentes em termos da probabilidade de conseguir 10 cartas do mesmo naipe?
  - a. Usando um baralho normal, você retira 10 cartas e registra as cartas retiradas.

- b. Usando um baralho normal, você repete a seguinte sequência de ações 10 vezes: retirar uma carta do baralho, registrar a carta retirada e repor a carta no baralho.
  - c. Usando uma caixa contendo todas as cartas de 1 milhão de baralhos reunidos, você retira 10 cartas e registra as cartas retiradas.
  - d. Usando uma caixa contendo todas as cartas de 1 milhão de baralhos reunidos, você repete a seguinte sequência de ações 10 vezes: retirar uma carta da caixa, registrar a carta retirada e repor a carta na caixa.
  - e. Usando um baralho *infinito*, você retira 10 cartas e registra as cartas retiradas.
  - f. Usando um baralho *infinito*, você repete a seguinte sequência de ações 10 vezes: retirar uma carta do baralho, registrar a carta retirada e repor a carta no baralho.
4. Qual a graça dos quadrinhos na Figura 1.1, que também aparecem no vídeo<sup>1</sup>?



Figura 1.1: <http://xkcd.com/552/>

5. Qual a graça dos quadrinhos na Figura 1.2?
6. Veja este vídeo sobre o cavalo Hans:

<https://youtu.be/G3VkCmdUfZE>

Qual a relação entre esta história e a necessidade de duplo cegamento?

---

<sup>1</sup>[https://youtu.be/6Q\\_XSoLCIpc?t=1385](https://youtu.be/6Q_XSoLCIpc?t=1385)





LIMITAÇÕES DE ESTUDOS COM CEGAMENTO

Figura 1.2: <http://xkcd.com/1462/>

## 1.3

---

### Vídeo 2

<https://youtu.be/492VASxlDRo>

## 1.4

---

### Exercícios

1. Por que não faz sentido calcular a média dos CEPs de um grupo de pessoas?
2. Uma temperatura de  $-40$  graus Celsius é igual a uma temperatura de  $-40$  graus Fahrenheit?
3. Uma temperatura de zero graus Celsius é igual a uma temperatura de zero graus Fahrenheit?
4. Uma variação de temperatura de 1 grau Celsius é igual a uma variação de temperatura de 1 grau Fahrenheit?
5. Um saldo bancário de zero reais é igual a um saldo bancário de zero dólares?
6. Um produto de 1 milhão de reais custa o mesmo que um produto de 1 milhão de dólares?
7. Meses representados por números de 1 a 12 são dados de que nível?

### Introdução a R

---

#### 2.1

---

##### Vídeo 1

<https://youtu.be/1kXQDNqm41c>

#### 2.2

---

##### Vídeo 2

<https://youtu.be/3GEc1oiKDrU>

#### 2.3

---

##### Exercícios

1. Para criar sua conta no RStudio Cloud, acesse <https://rstudio.cloud/>.
2. Se você preferir instalar o R no seu computador, acesse
  - <https://cran.r-project.org/> para baixar e instalar o R, e
  - <https://rstudio.com/products/rstudio/download/> para baixar e instalar o RStudio, um IDE específico para R.

3. Abra o RStudio Cloud ou o seu RStudio instalado localmente.
4. Crie um novo projeto. Sempre trabalhe em projetos para ter seus arquivos organizados.

5. Para instalar o `swirl` (pacote do R para exercícios interativos)<sup>1</sup>, execute o seguinte comando no console do RStudio:

```
install.packages("swirl")
```

6. Para instalar os exercícios de introdução a R, execute os seguintes comandos no console do RStudio:

```
library(swirl)
install_course_github('fnaufel', 'introR')
```

7. Mude o idioma para português e execute o `swirl`.

```
select_language('portuguese', append_rprofile = TRUE)
swirl()
```

8. Na primeira execução, você vai precisar se identificar (qualquer nome serve). Com essa identificação, o `swirl` vai registrar o seu progresso nas lições.
9. No `swirl`, as perguntas são mostradas no console. Você também deve responder no console.
10. Às vezes, um *script* será aberto no editor de textos para que você complete um programa. Quando seu programa estiver pronto, salve o arquivo e digite `submit()` no console para o `swirl` processar o *script*.
11. O `swirl` dá instruções claras no console. Na dúvida, digite `info()` no *prompt* do R (`>`).
12. Se, em vez do *prompt* do R, o console mostrar reticências (`. . .`), tecle *Enter*.
13. Se nada funcionar, tecle *ESC*.
14. Para sair do `swirl()`, digite `bye()` no *prompt* do R.
15. Para voltar para os exercícios, digite

```
library(swirl)
swirl()
```

---

<sup>1</sup><https://swirlstats.com/>

---

### Visualização com ggplot2

---



Busque mais informações sobre os pacotes `tidyverse` e `ggplot2` nas referências recomendadas.

#### 3.1

---

##### Vídeo 1

<https://youtu.be/OBpNjqIIyhI>

#### 3.2

---

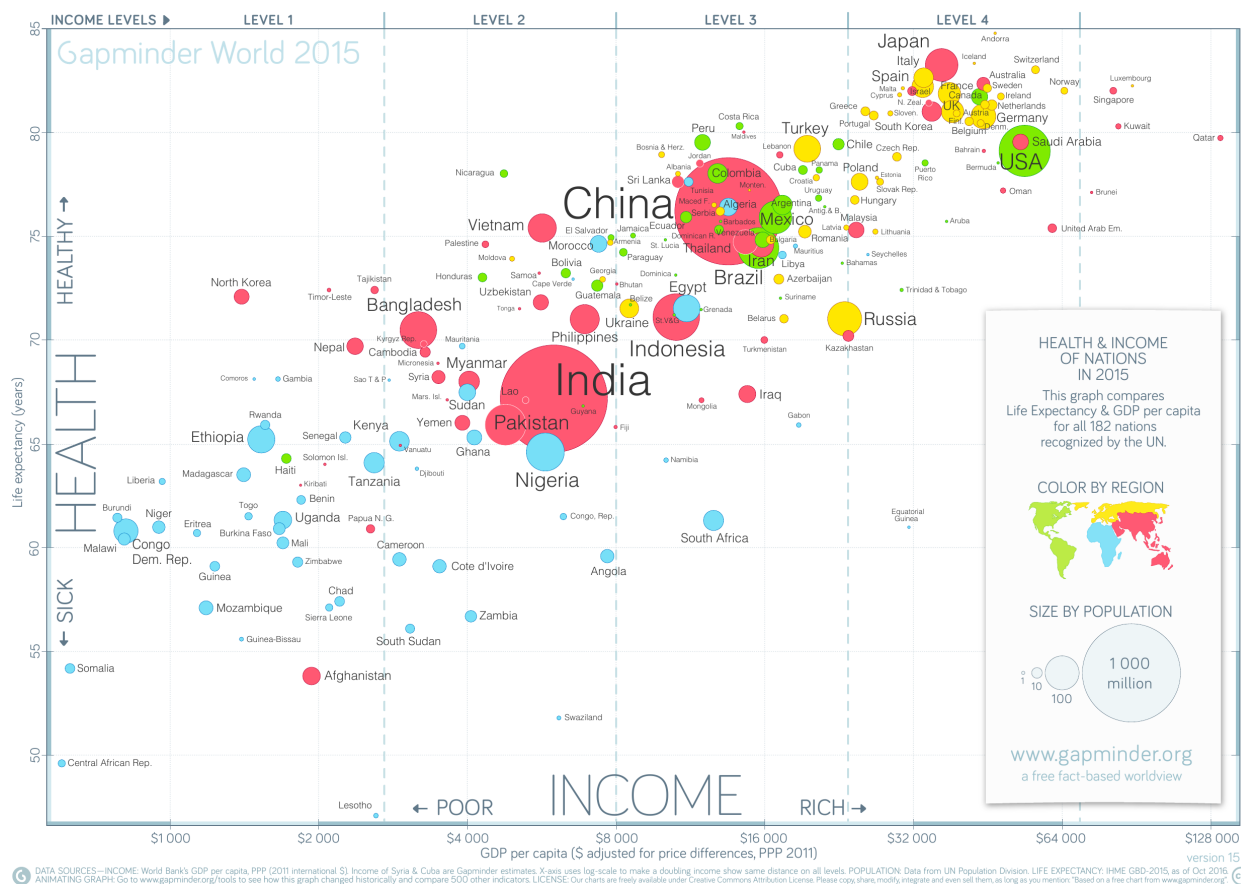
### Componentes de um gráfico ggplot2

#### 3.2.1

---

##### Geometrias e mapeamentos estéticos (*mappings*)

- Observe o gráfico abaixo, obtido de <https://www.gapminder.org/downloads/updated-gapminder-world-poster-2015/>.



- O gráfico mostra como, em cada país, a saúde (mais precisamente, a expectativa de vida) se relaciona com a riqueza (mais precisamente, o PIB *per capita*).
- Além da expectativa de vida e o do PIB *per capita*, o gráfico traz mais informações sobre cada país.
- Cada país é representado por um ponto (a **geometria**).
- Informações sobre cada país são representadas por características do ponto correspondente (as **estéticas**):

Variável	Geometria	Estética
PIB <i>per capita</i>	ponto	posição x
Expectativa de vida	ponto	posição y
População	ponto	tamanho
Continente	ponto	cor

- Você pode usar outras estéticas para representar informações:
  - Cor de preenchimento.
  - Cor do traço.
  - Tipo do traço (sólido, pontilhado, tracejado etc.).
  - Forma (círculo, quadrado, triângulo etc.).

- Opacidade.
- etc.
- Você pode usar outras geometrias:
  - Linhas.
  - Barras ou colunas.
  - Caixas.
  - etc.

### 3.2.2

---

#### Escalas (*scales*)

- As escalas controlam os detalhes da aparência da geometria e do mapeamento (eixos, cores etc.).
- Os eixos do gráfico acima são escalas **contínuas**, com valores reais.
- Observe o eixo horizontal. Os valores não aumentam linearmente, mas sim exponencialmente: cada passo à direita equivale a *dobrar* o valor do PIB. O eixo horizontal segue uma **escala logarítmica**.
- Os tamanhos dos pontos formam uma escala **discreta**, com 4 valores possíveis (veja a legenda no canto inferior direito do gráfico).
- As cores também formam uma escala discreta.

### 3.2.3

---

#### Rótulos (*labels*)

- O gráfico também representa informação na forma de texto.
- Além de rótulos (por exemplo, o texto que identifica cada eixo), **o texto também pode, ele mesmo, ser uma geometria, com suas próprias estéticas**: observe como o nome de cada país é escrito em um tamanho proporcional à sua população.

### 3.2.4

---

#### Outros componentes

- Coordenadas:
  - Este gráfico usa **coordenadas cartesianas**, com eixos  $x$  e  $y$ .
  - Existem gráficos que usam um sistema de **coordenadas polares**.
- Temas:
  - Incluem todos os elementos “decorativos”: cor de fundo, linhas de grade, etc. Ajudam a facilitar a leitura e a interpretação.

- No gráfico acima, um detalhe interessante do tema é a divisão de cada eixo em segmentos claros e segmentos escuros.
- Legendas (*guides*).
- Facetas:
  - Às vezes, um gráfico é composto por múltiplos subgráficos.
  - Cada subgráfico é uma **faceta**.
  - Facetas evitam que informações demais sejam apresentadas no mesmo lugar.

### 3.3

## Conjunto de dados

- Nossos exemplos de gráficos vão usar dados sobre o sono de diversos mamíferos.
- O conjunto de dados se chama `msleep` e está incluído no pacote `ggplot2`.
- Para ver a documentação, digite

```
library(ggplot2)
?msleep
```

- Vamos atribuir o conjunto de dados à variável `df`:

```
df <- msleep
df
## # A tibble: 83 x 11
##   name      genus  vore  order  conservation sleep_total sleep_rem
##   <chr>     <chr>  <chr> <chr>   <chr>          <dbl>    <dbl>
## 1 Cheetah   Acinon~ carni Carniv~ lc           12.1      NA
## 2 Owl monkey Aotus   omni  Primat~ <NA>        17        1.8
## 3 Mountain be~ Aplodo~ herbi Rodent~ nt         14.4      2.4
## 4 Greater sho~ Blarina omni  Sorico~ lc          14.9      2.3
## 5 Cow       Bos     herbi Artiod~ domesticated  4         0.7
## 6 Three-toed ~ Bradyp~ herbi Pilosa <NA>        14.4      2.2
## # ... with 77 more rows, and 4 more variables: sleep_cycle <dbl>,
## #   awake <dbl>, brainwt <dbl>, bodywt <dbl>
```

- Vamos examinar a estrutura — usando R base:

```
str(df)
## tibble [83 x 11] (S3: tbl_df/tbl/data.frame)
##  $ name      : chr [1:83] "Cheetah" "Owl monkey" "Mountain beaver" ...
##  $ genus     : chr [1:83] "Acinonyx" "Aotus" "Aplodontia" ...
##  $ vore      : chr [1:83] "carni" "omni" "herbi" ...
##  $ order     : chr [1:83] "Carnivora" "Primates" "Rodentia" ...
```



```
## $ conservation: chr [1:83] "lc" NA "nt" ...
## $ sleep_total : num [1:83] 12,1 17 14,4 14,9 4 14,4 8,7 7 ...
## $ sleep_rem : num [1:83] NA 1,8 2,4 2,3 0,7 2,2 1,4 NA ...
## $ sleep_cycle : num [1:83] NA NA NA 0,133 ...
## $ awake : num [1:83] 11,9 7 9,6 9,1 20 9,6 15,3 17 ...
## $ brainwt : num [1:83] NA 0,0155 NA 0,00029 0,423 NA NA NA ...
## $ bodywt : num [1:83] 50 0,48 1,35 0,019 ...
```

- Podemos usar `glimpse`, uma função do `tidyverse`:

```
glimpse(df)
## Rows: 83
## Columns: 11
## $ name <chr> "Cheetah", "Owl monkey", "Mountain beaver", "Gre~
## $ genus <chr> "Acinonyx", "Aotus", "Aplodontia", "Blarina", "B~
## $ vore <chr> "carni", "omni", "herbi", "omni", "herbi", "herb~
## $ order <chr> "Carnivora", "Primates", "Rodentia", "Soricomorp~
## $ conservation <chr> "lc", NA, "nt", "lc", "domesticated", NA, "vu", ~
## $ sleep_total <dbl> 12,1, 17,0, 14,4, 14,9, 4,0, 14,4, 8,7, 7,0, 10,~
## $ sleep_rem <dbl> NA, 1,8, 2,4, 2,3, 0,7, 2,2, 1,4, NA, 2,9, NA, 0~
## $ sleep_cycle <dbl> NA, NA, NA, 0,1333333, 0,6666667, 0,7666667, 0,3~
## $ awake <dbl> 11,9, 7,0, 9,6, 9,1, 20,0, 9,6, 15,3, 17,0, 13,9~
## $ brainwt <dbl> NA, 0,01550, NA, 0,00029, 0,42300, NA, NA, NA, 0~
## $ bodywt <dbl> 50,000, 0,480, 1,350, 0,019, 600,000, 3,850, 20,~
```

- Para examinar só as primeiras linhas do *data frame*:

```
head(df)
## # A tibble: 6 x 11
##   name      genus vore order conservation sleep_total sleep_rem
##   <chr>      <chr> <chr> <chr>      <chr>          <dbl>      <dbl>
## 1 Cheetah    Acinon~ carni Carniv~ lc              12.1        NA
## 2 Owl monkey Aotus   omni  Primat~ <NA>          17          1.8
## 3 Mountain be~ Aplodo~ herbi Rodent~ nt              14.4         2.4
## 4 Greater sho~ Blarina omni  Sorico~ lc              14.9         2.3
## 5 Cow        Bos     herbi Artiod~ domesticated    4           0.7
## 6 Three-toed ~ Bradyp~ herbi Pilosa <NA>          14.4         2.2
## # ... with 4 more variables: sleep_cycle <dbl>, awake <dbl>,
## #   brainwt <dbl>, bodywt <dbl>
```

- Para examinar o *data frame* interativamente:

```
view(df)
```

- Podemos produzir um sumário dos dados usando o pacote *summarytools* (que já foi carregado neste documento):

```
df %>% dfSummary() %>% print()
```

Variável	Estatísticas / Valores	Freqs (% de Válidos)	Faltante
name [character]	1. African elephant 2. African giant pouched rat 3. African striped mouse 4. Arctic fox 5. Arctic ground squirrel 6. Asian elephant 7. Baboon 8. Big brown bat 9. Bottle-nosed dolphin 10. Brazilian tapir [ 73 outros ]	1 ( 1,2%) 1 ( 1,2%) 1 ( 1,2%) 1 ( 1,2%) 1 ( 1,2%) 1 ( 1,2%) 1 ( 1,2%) 1 ( 1,2%) 1 ( 1,2%) 1 ( 1,2%) 73 (88,0%)	0 (0,0%)
genus [character]	1. Panthera 2. Spermophilus 3. Equus 4. Vulpes 5. Acinonyx 6. Aotus 7. Aplodontia 8. Blarina 9. Bos 10. Bradypus [ 67 outros ]	3 ( 3,6%) 3 ( 3,6%) 2 ( 2,4%) 2 ( 2,4%) 1 ( 1,2%) 1 ( 1,2%) 1 ( 1,2%) 1 ( 1,2%) 1 ( 1,2%) 1 ( 1,2%) 67 (80,7%)	0 (0,0%)
vore [character]	1. carni 2. herbi 3. insecti 4. omni	19 (25,0%) 32 (42,1%) 5 ( 6,6%) 20 (26,3%)	7 (8,4%)
order [character]	1. Rodentia 2. Carnivora 3. Primates 4. Artiodactyla 5. Soricomorpha 6. Cetacea 7. Hyracoidea 8. Perissodactyla 9. Chiroptera 10. Cingulata [ 9 outros ]	22 (26,5%) 12 (14,5%) 12 (14,5%) 6 ( 7,2%) 5 ( 6,0%) 3 ( 3,6%) 3 ( 3,6%) 3 ( 3,6%) 2 ( 2,4%) 2 ( 2,4%) 13 (15,7%)	0 (0,0%)
conservation [character]	1. cd 2. domesticated 3. en 4. lc 5. nt 6. vu	2 ( 3,7%) 10 (18,5%) 4 ( 7,4%) 27 (50,0%) 4 ( 7,4%) 7 (13,0%)	29 (34,9%)

Variável	Estatísticas / Valores	Freqs (% de Válidos)	Faltante
sleep_total [numeric]	Média (dp) : 10,4 (4,5) mín < mediana < máx: 1,9 < 10,1 < 19,9 IQE (CV) : 5,9 (0,4)	65 valores distintos	0 (0,0%)
sleep_rem [numeric]	Média (dp) : 1,9 (1,3) mín < mediana < máx: 0,1 < 1,5 < 6,6 IQE (CV) : 1,5 (0,7)	32 valores distintos	22 (26,5%)
sleep_cycle [numeric]	Média (dp) : 0,4 (0,4) mín < mediana < máx: 0,1 < 0,3 < 1,5 IQE (CV) : 0,4 (0,8)	22 valores distintos	51 (61,4%)
awake [numeric]	Média (dp) : 13,6 (4,5) mín < mediana < máx: 4,1 < 13,9 < 22,1 IQE (CV) : 5,9 (0,3)	65 valores distintos	0 (0,0%)
brainwt [numeric]	Média (dp) : 0,3 (1) mín < mediana < máx: 0 < 0 < 5,7 IQE (CV) : 0,1 (3,5)	53 valores distintos	27 (32,5%)
bodywt [numeric]	Média (dp) : 166,1 (786,8) mín < mediana < máx: 0 < 1,7 < 6654 IQE (CV) : 41,6 (4,7)	82 valores distintos	0 (0,0%)

- Vemos que há muitos NA em diversas variáveis. Para nossos exemplos simples de visualização, vamos usar as colunas

- name
- genus
- order
- sleep\_total
- awake
- bodywt
- brainwt

- Mas... a coluna que mostra a dieta (vore) tem só 7 NA. Quais são?

```
df %>%
  filter(is.na(vore)) %>%
  select(name)
## # A tibble: 7 x 1
##   name
##   <chr>
## 1 Vesper mouse
```

```
## 2 Desert hedgehog
## 3 Deer mouse
## 4 Phalanger
## 5 Rock hyrax
## 6 Mole rat
## # ... with 1 more row
```

- OK. Vamos manter a coluna `vore` também, apesar dos NA. Quando formos usar esta variável, tomaremos cuidado.
- Também... a coluna `bodywt` tem 0 como valor mínimo. Como assim?

```
df %>%
  filter(bodywt < 1) %>%
  select(name, bodywt) %>%
  arrange(bodywt)
## # A tibble: 35 x 2
##   name                bodywt
##   <chr>              <dbl>
## 1 Lesser short-tailed shrew 0.005
## 2 Little brown bat        0.01
## 3 Greater short-tailed shrew 0.019
## 4 Deer mouse             0.021
## 5 House mouse            0.022
## 6 Big brown bat          0.023
## # ... with 29 more rows
```

- Ah, sem problema. A função `dfSummary` arredondou estes pesos para 0. Os valores de verdade ainda estão na *tibble*.
- Vamos criar uma *tibble* nova, só com as colunas que nos interessam:

```
sono <- df %>%
  select(
    name, order, genus, vore, bodywt,
    brainwt, awake, sleep_total
  )
```

- Vamos ver o sumário:

```
sono %>% dfSummary() %>% print()
```

Variável	Estatísticas / Valores	Freqs (% de Válidos)	Faltante
name [character]	1. African elephant 2. African giant pouched rat 3. African striped mouse 4. Arctic fox 5. Arctic ground squirrel 6. Asian elephant 7. Baboon 8. Big brown bat 9. Bottle-nosed dolphin 10. Brazilian tapir [ 73 outros ]	1 ( 1,2%) 1 ( 1,2%) 1 ( 1,2%) 1 ( 1,2%) 1 ( 1,2%) 1 ( 1,2%) 1 ( 1,2%) 1 ( 1,2%) 1 ( 1,2%) 1 ( 1,2%) 73 (88,0%)	0 (0,0%)
order [character]	1. Rodentia 2. Carnivora 3. Primates 4. Artiodactyla 5. Soricomorpha 6. Cetacea 7. Hyracoidea 8. Perissodactyla 9. Chiroptera 10. Cingulata [ 9 outros ]	22 (26,5%) 12 (14,5%) 12 (14,5%) 6 ( 7,2%) 5 ( 6,0%) 3 ( 3,6%) 3 ( 3,6%) 3 ( 3,6%) 2 ( 2,4%) 2 ( 2,4%) 13 (15,7%)	0 (0,0%)
genus [character]	1. Panthera 2. Spermophilus 3. Equus 4. Vulpes 5. Acinonyx 6. Aotus 7. Aplodontia 8. Blarina 9. Bos 10. Bradypus [ 67 outros ]	3 ( 3,6%) 3 ( 3,6%) 2 ( 2,4%) 2 ( 2,4%) 1 ( 1,2%) 1 ( 1,2%) 1 ( 1,2%) 1 ( 1,2%) 1 ( 1,2%) 1 ( 1,2%) 67 (80,7%)	0 (0,0%)
vore [character]	1. carni 2. herbi 3. insecti 4. omni	19 (25,0%) 32 (42,1%) 5 ( 6,6%) 20 (26,3%)	7 (8,4%)
bodywt [numeric]	Média (dp) : 166,1 (786,8) mín < mediana < máx: 0 < 1,7 < 6654 IQE (CV) : 41,6 (4,7)	82 valores distintos	0 (0,0%)
brainwt [numeric]	Média (dp) : 0,3 (1) mín < mediana < máx: 0 < 0 < 5,7 IQE (CV) : 0,1 (3,5)	53 valores distintos	27 (32,5%)

Variável	Estatísticas / Valores	Freqs (% de Válidos)	Faltante
awake [numeric]	Média (dp) : 13,6 (4,5) mín < mediana < máx: 4,1 < 13,9 < 22,1 IQE (CV) : 5,9 (0,3)	65 valores distintos	0 (0,0%)
sleep_total [numeric]	Média (dp) : 10,4 (4,5) mín < mediana < máx: 1,9 < 10,1 < 19,9 IQE (CV) : 5,9 (0,4)	65 valores distintos	0 (0,0%)

### 3.4

#### Gráficos de dispersão (*scatter plots*)

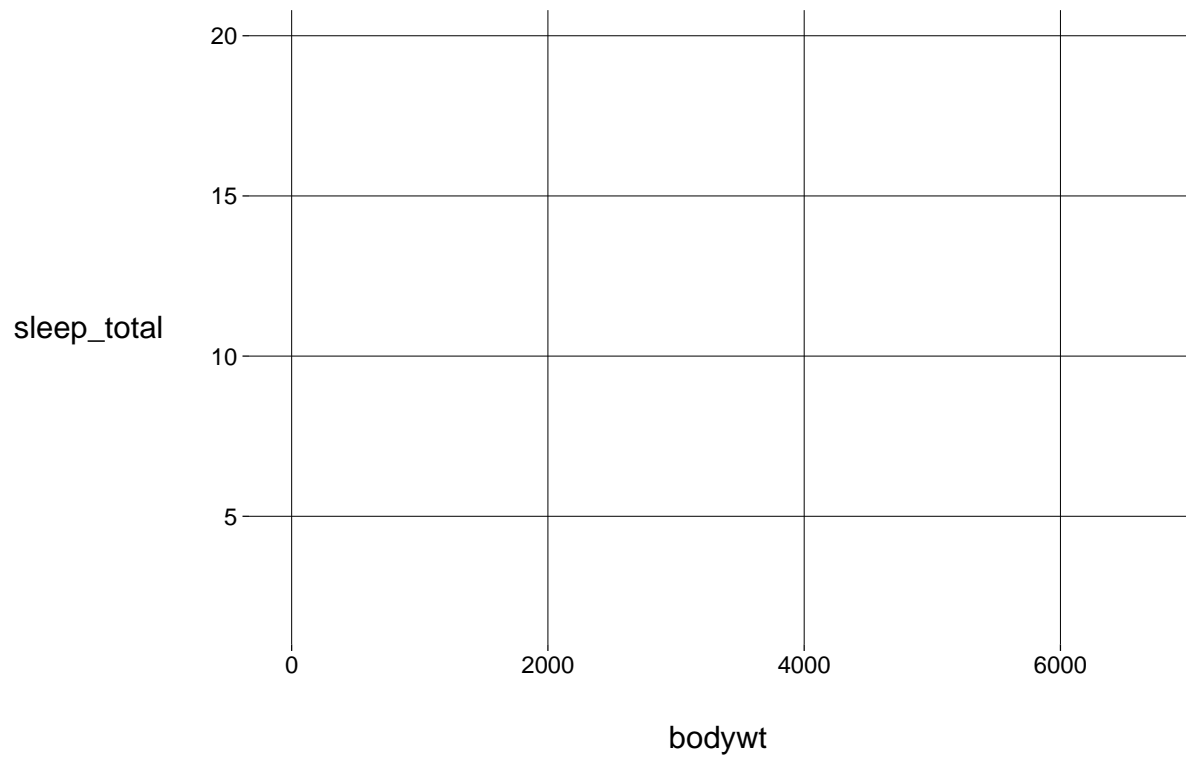
- Servem para visualizar a *relação* entre **duas variáveis quantitativas**.
- **Essa relação não é necessariamente de causa e efeito.**
- Isto é, a variável do eixo horizontal não determina, necessariamente, os valores da variável do eixo vertical.
- Pense em **associação**, **correlação**, não em causalidade.
- Troque as variáveis de eixo, se ajudar a deixar isto claro.

#### 3.4.1

##### Horas de sono e peso corporal

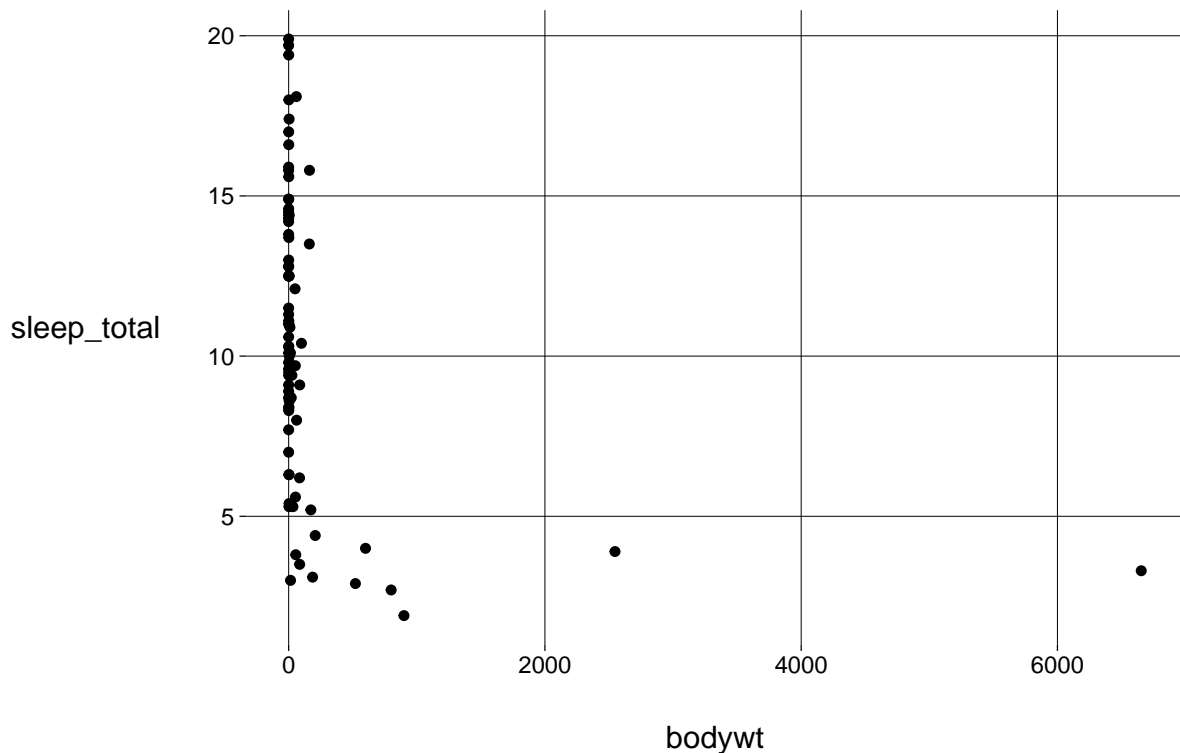
- Como as variáveis `sleep_total` e `bodywt` estão relacionadas?

```
sono %>%
  ggplot(aes(x = bodywt, y = sleep_total))
```



- O que houve? Cadê os pontos?
- O problema foi que só especificamos o mapeamento estético (com `aes`, que são as iniciais de *aesthetics*). **Faltou a geometria.**

```
sono %>%  
  ggplot(aes(x = bodywt, y = sleep_total)) +  
  geom_point()
```



- Que horror.
- A única coisa que percebemos aqui é que os mamíferos muito pesados dormem menos de 5 horas por noite.
- Estes animais muito pesados estão estragando a escala do eixo  $x$ .
- Que animais são estes?

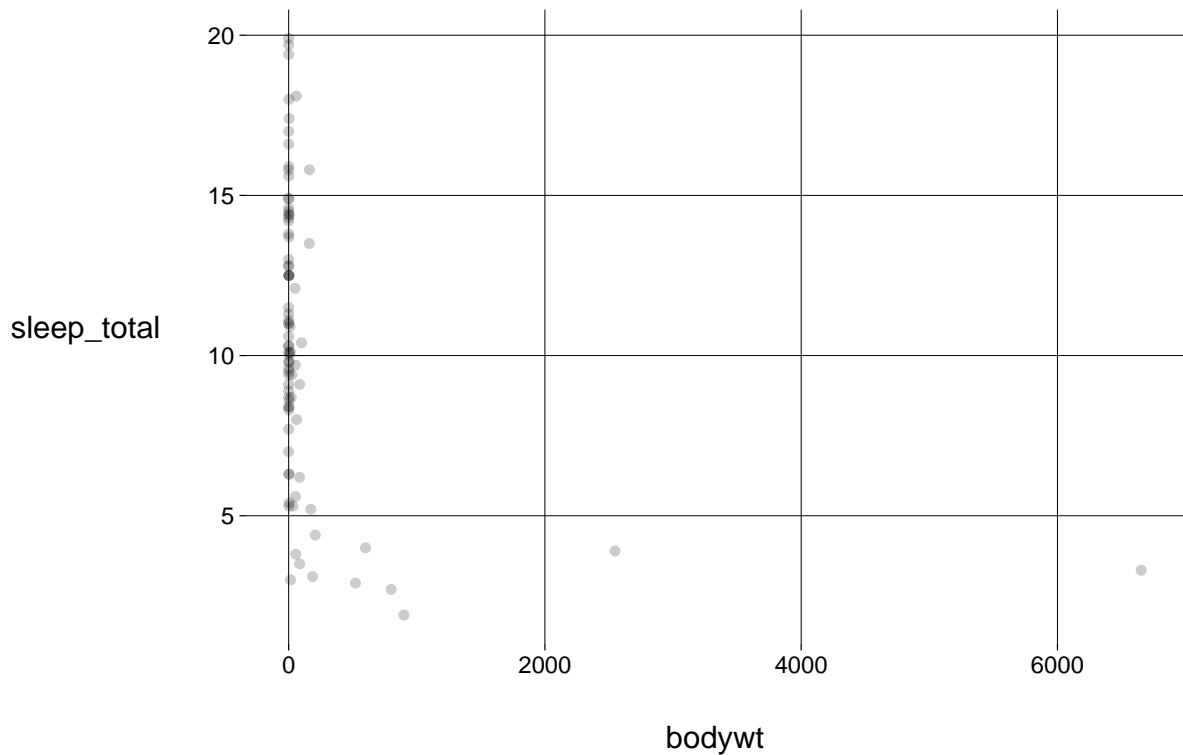
```
sono %>%
  filter(bodywt > 250) %>%
  select(name, bodywt) %>%
  arrange(bodywt)
## # A tibble: 6 x 2
##   name          bodywt
##   <chr>         <dbl>
## 1 Horse          521
## 2 Cow            600
## 3 Pilot whale   800
## 4 Giraffe       900.
## 5 Asian elephant 2547
## 6 African elephant 6654
```

- Além disso, há muitos pontos sobrepostos. Em bom português, temos um problema de *overplotting*.
- Existem diversas maneiras de lidar com isso.
- A primeira delas é alterando a opacidade dos pontos. Isto é um ajuste na geometria



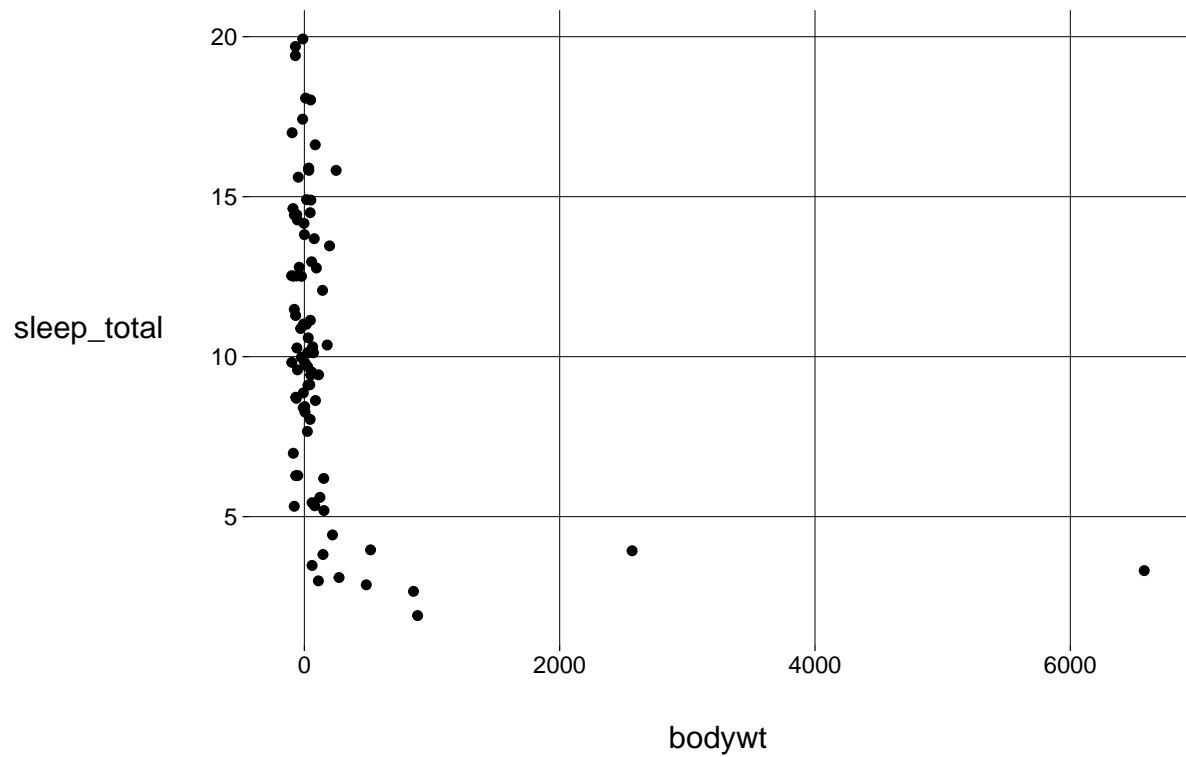
apenas, pois a opacidade, aqui, não representa informação nenhuma.

```
sono %>%  
  ggplot(aes(x = bodywt, y = sleep_total)) +  
    geom_point(alpha = 0.2)
```



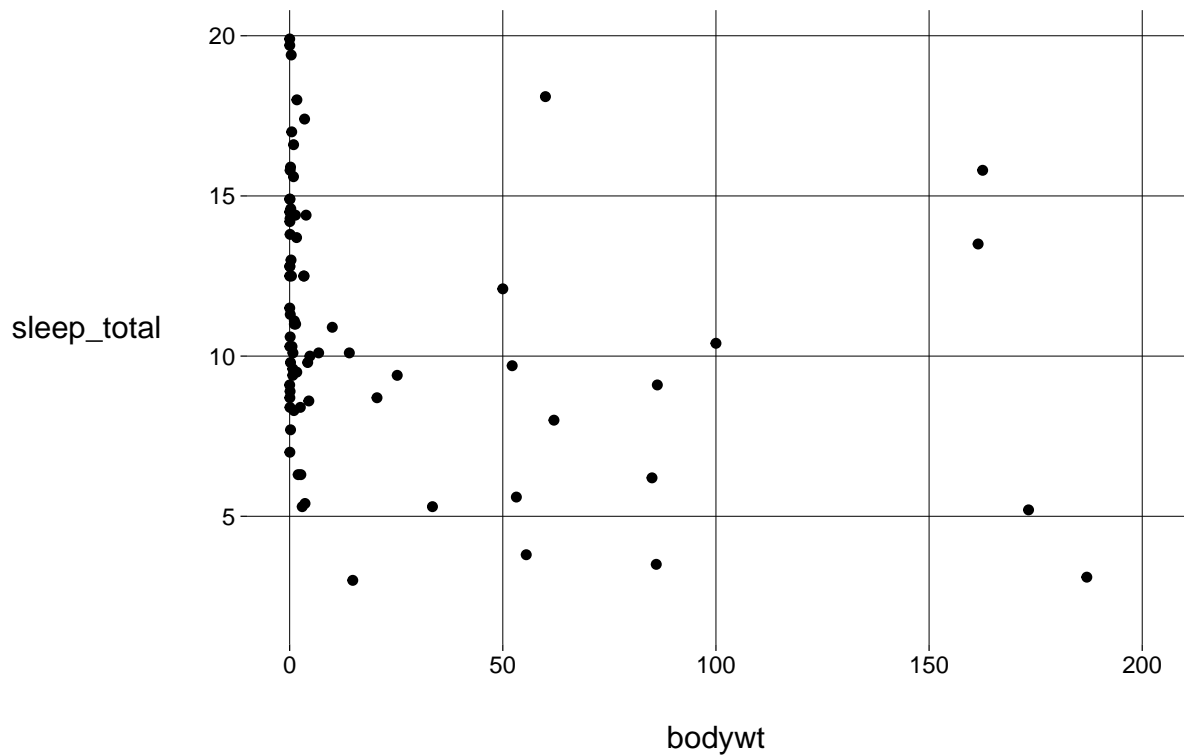
- Outra maneira é usar `geom_jitter` em vez de `geom_point`. “*Jitter*” significa “tremar”. As posições dos pontos são ligeiramente perturbadas, para evitar colisões. Perdemos precisão, mas a visualização fica melhor.

```
sono %>%  
  ggplot(aes(x = bodywt, y = sleep_total)) +  
    geom_jitter(width = 100)
```



- Vamos mudar os limites do gráfico para nos concentrarmos nos animais menos pesados. Observe que **isto é um ajuste na escala.**

```
sono %>%  
  ggplot(aes(x = bodywt, y = sleep_total)) +  
    geom_point() +  
    scale_x_continuous(limits = c(0, 200))  
## Warning: Removed 7 rows containing missing values (geom_point).
```



- Nestes limites, a relação entre horas de sono e peso não é mais tão pronunciada.

### 3.4.2

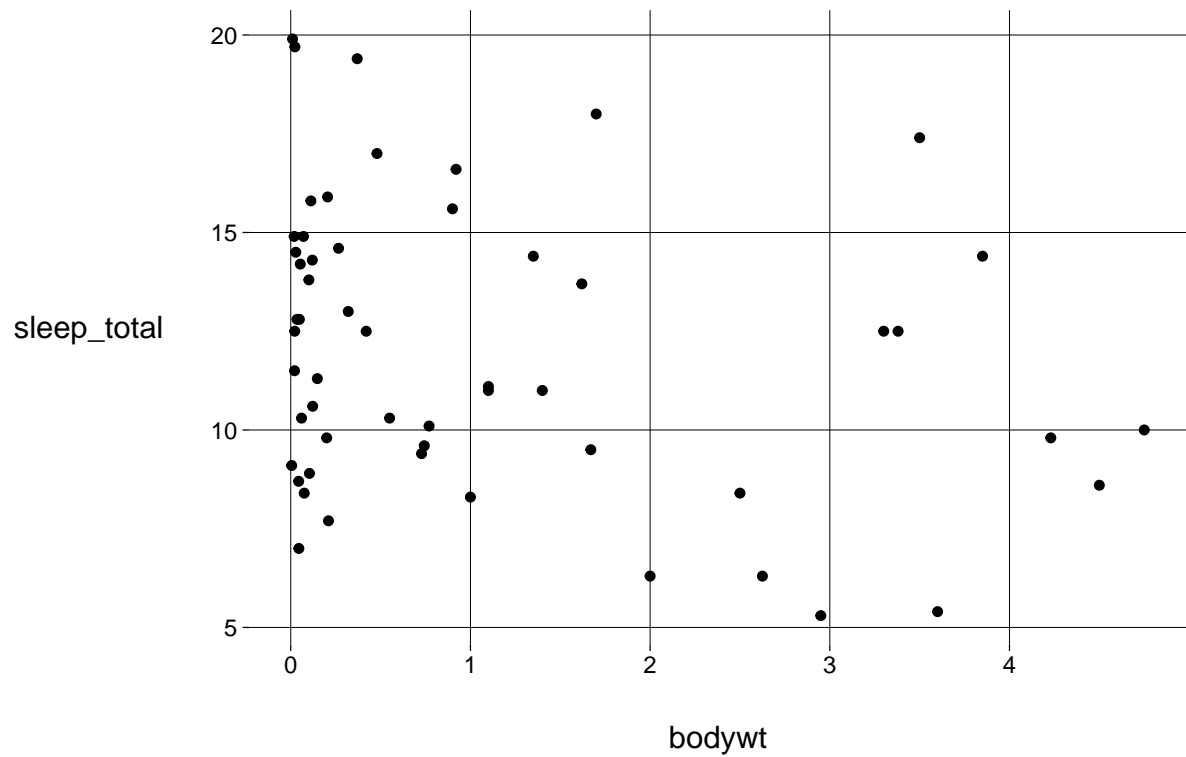
#### Horas de sono e peso corporal para animais pequenos

- Vamos restringir o gráfico a animais com no máximo 5kg.

```
limite <- 5
```

- Em vez de mudar a escala do gráfico, vamos filtrar as linhas do *data frame*:

```
sono %>%
  filter(bodywt < limite) %>%
  ggplot(aes(x = bodywt, y = sleep_total)) +
  geom_point()
```

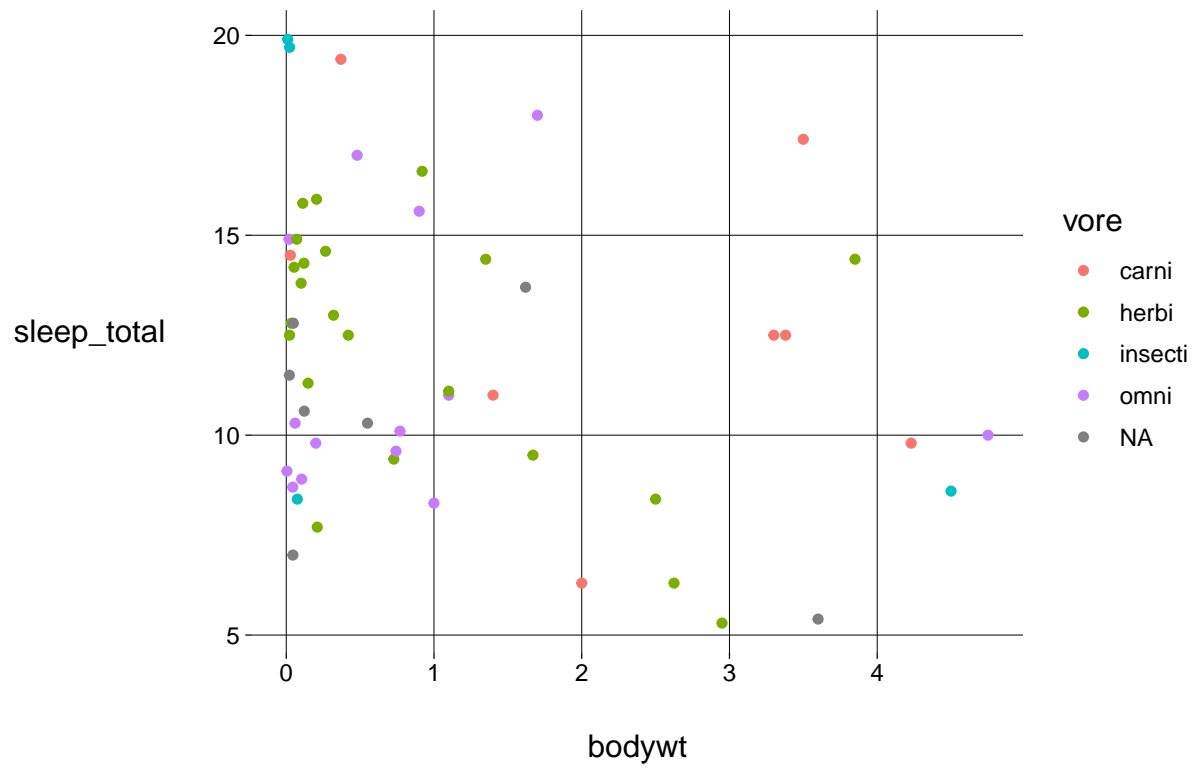


### 3.4.3

#### Incluindo a dieta

- Com a estética `color`. Observe como a legenda aparece automaticamente.

```
sono %>%  
  filter(bodywt < limite) %>%  
  ggplot(aes(x = bodywt, y = sleep_total, color = vore)) +  
    geom_point()
```

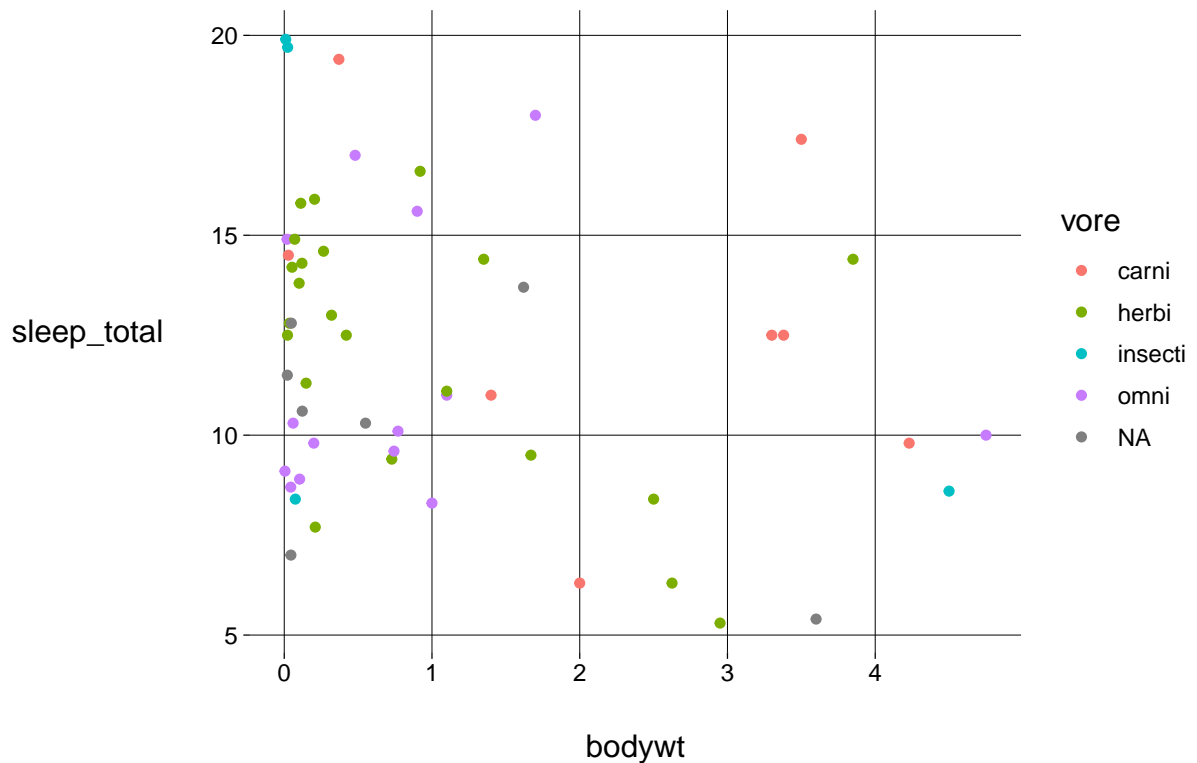


### 3.4.4

A estética pode ser especificada na `geom`

- Compare com o código anterior.

```
sono %>%  
  filter(bodywt < limite) %>%  
  ggplot() +  
    geom_point(aes(x = bodywt, y = sleep_total, color = vore))
```



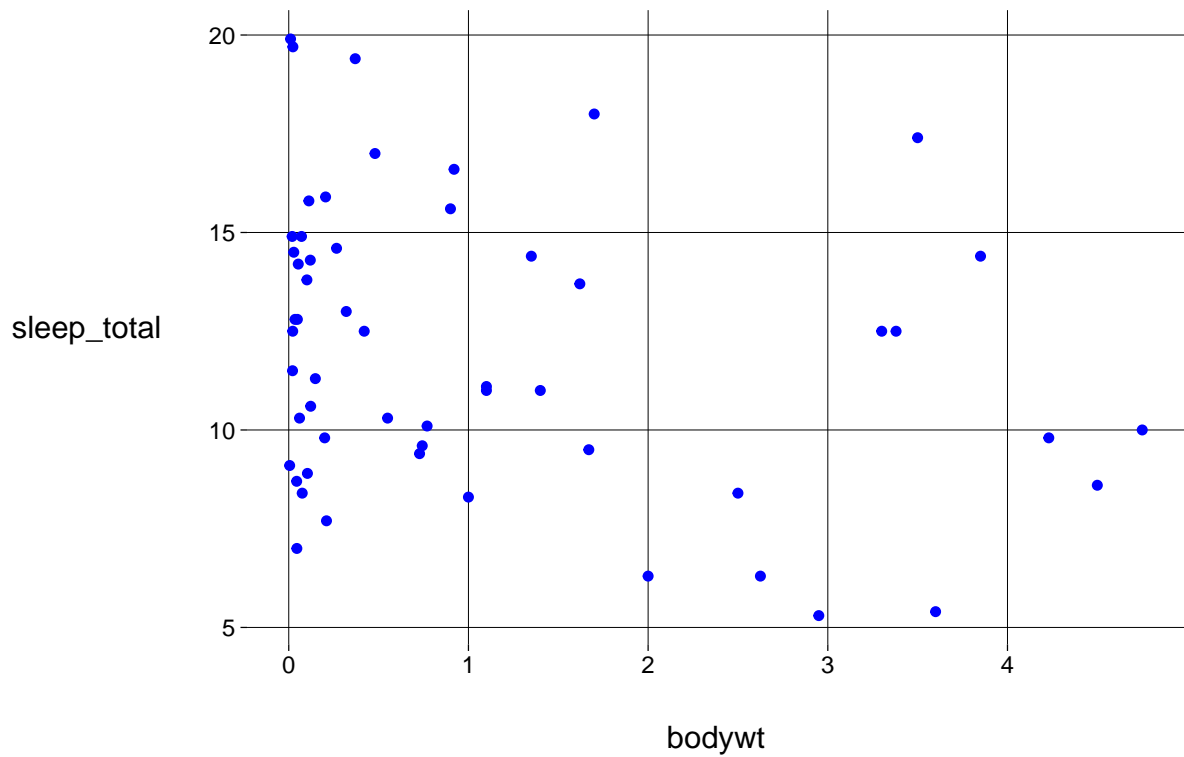
- Fazendo deste modo, a estética só vale para uma geometria. Se você acrescentar outras geometrias (linhas, por exemplo), a estética não valerá para elas.

### 3.4.5

#### Aparência fixa ou dependendo de variável?

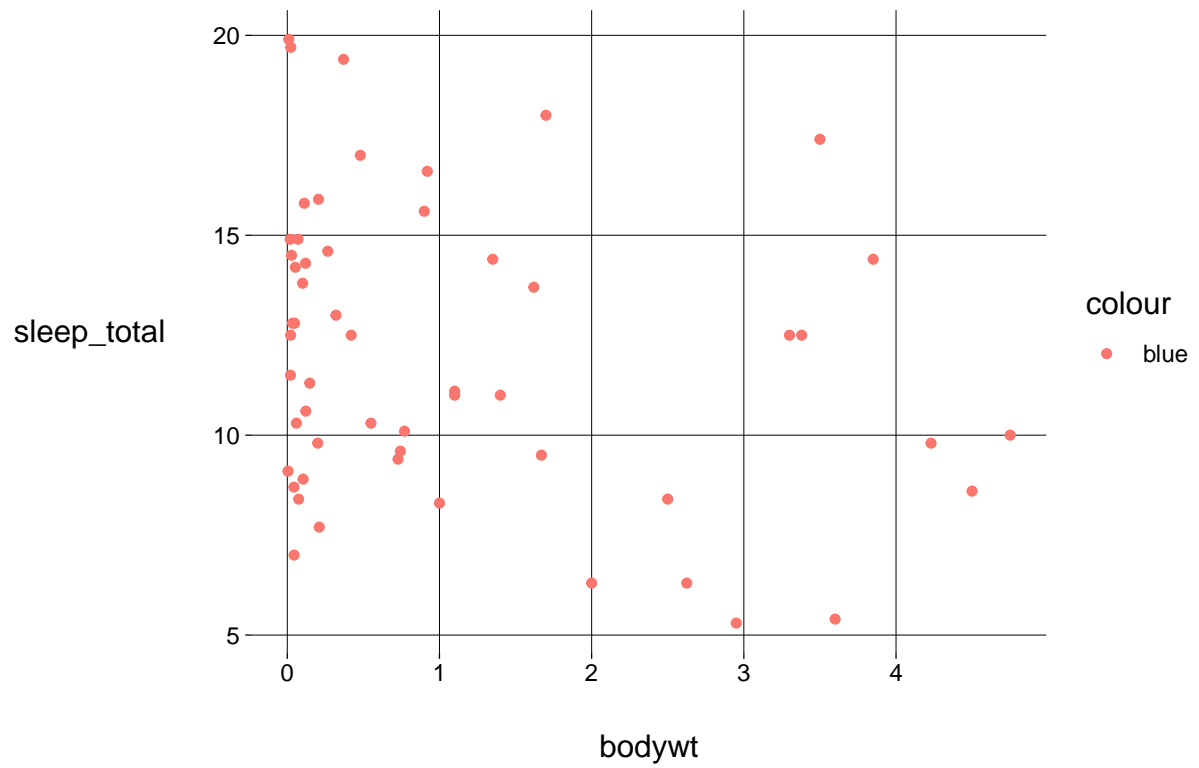
- Se for fixa, não é estética. Não representa informação.
- Se depender de variável, é estética. Representa informação.
- Compare o último *chunk* acima com:

```
sono %>%
  filter(bodywt < limite) %>%
  ggplot() +
    geom_point(aes(x = bodywt, y = sleep_total), color = 'blue')
```



- Se for uma estética, precisa estar associada a uma variável, não a um valor fixo. Um erro comum seria fazer:

```
sono %>%  
  filter(bodywt < limite) %>%  
  ggplot() +  
    geom_point(aes(x = bodywt, y = sleep_total, color = 'blue'))
```



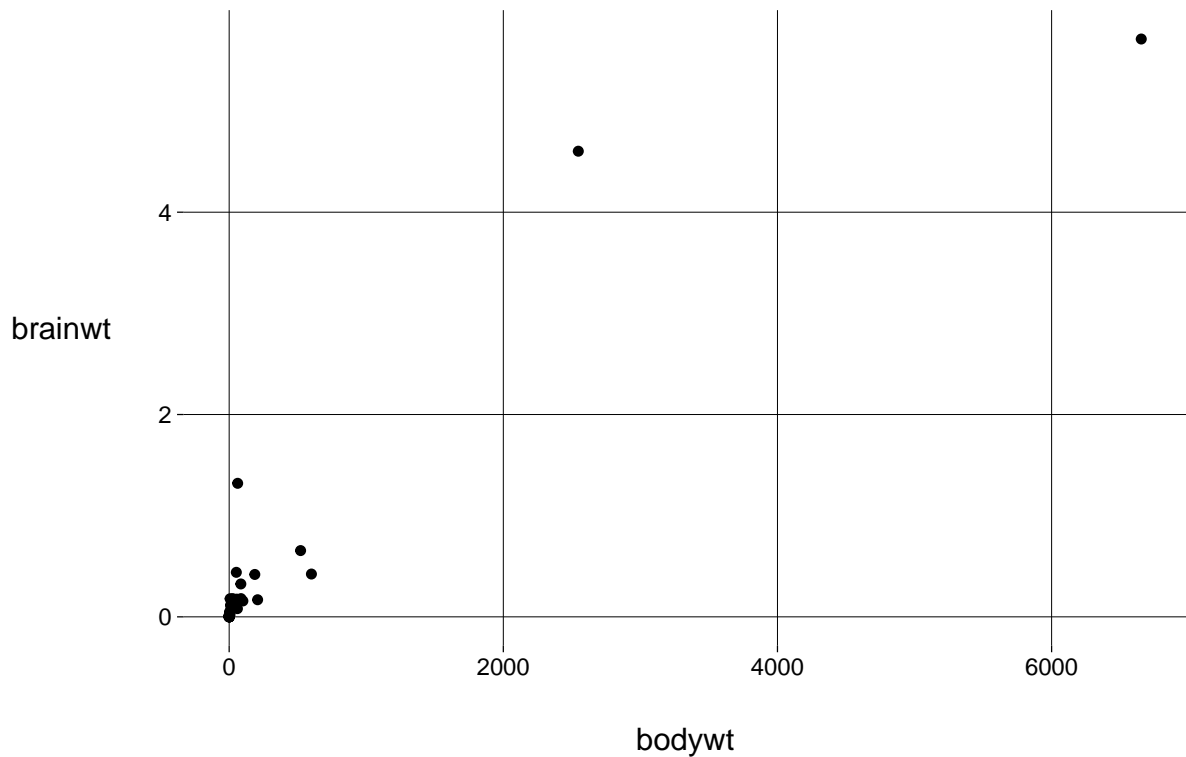
### 3.4.6

#### Uma correlação mais clara

- Peso cerebral versus peso corporal:

```
sono %>%  
  ggplot(aes(x = bodywt, y = brainwt)) +  
    geom_point()  
## Warning: Removed 27 rows containing missing values (geom_point).
```





- A mensagem de aviso (*warning*) diz que há 27 valores faltantes (NA) em bodywt ou brainwt. De fato:

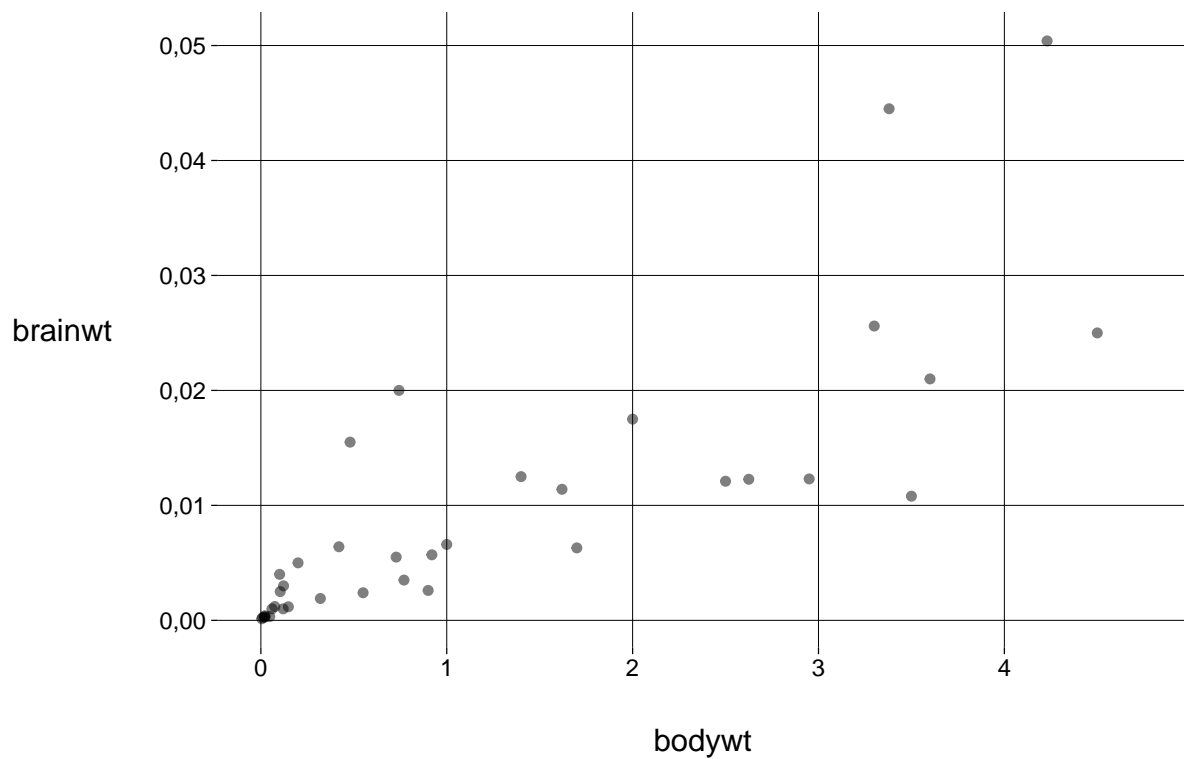
```
sono %>%
  filter(is.na(bodywt)) %>%
  count()
## # A tibble: 1 x 1
##       n
##   <int>
## 1     0
```

```
sono %>%
  filter(is.na(brainwt)) %>%
  count()
## # A tibble: 1 x 1
##       n
##   <int>
## 1    27
```

- Vamos restringir aos animais mais leves e mudar a opacidade:

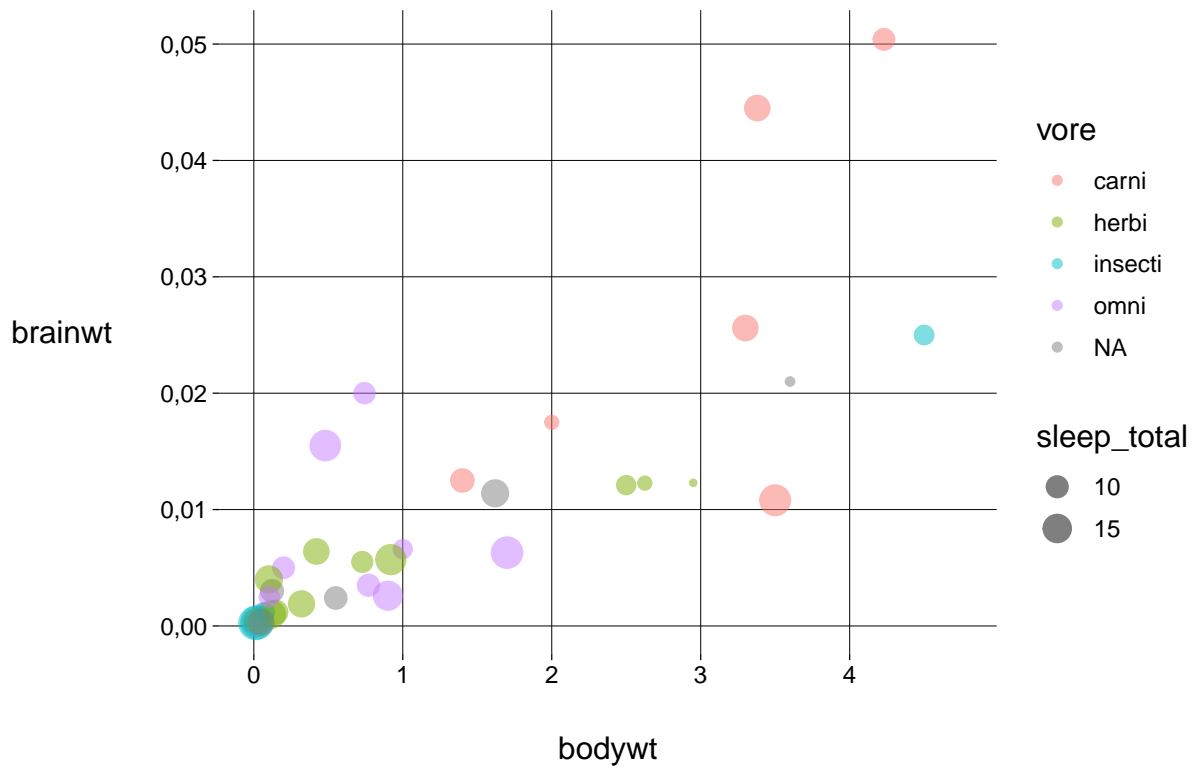
```
sono %>%
  filter(bodywt < limite) %>%
  ggplot(aes(x = bodywt, y = brainwt)) +
  geom_point(alpha = .5)
```

```
## Warning: Removed 18 rows containing missing values (geom_point).
```



- Vamos incluir horas de sono e dieta. Observe as estéticas usadas.

```
sono %>%  
  filter(bodywt < limite) %>%  
  ggplot(  
    aes(  
      x = bodywt,  
      y = brainwt,  
      size = sleep_total,  
      color = vore  
    )  
  ) +  
  geom_point(alpha = .5)  
## Warning: Removed 18 rows containing missing values (geom_point).
```



- Vamos mudar a escala dos tamanhos e incluir rótulos:

```
grafico <- sono %>%
  filter(bodywt < limite) %>%
  ggplot(
    aes(
      x = bodywt,
      y = brainwt,
      size = sleep_total,
      color = vore
    )
  ) +
  geom_point(alpha = .5) +
  scale_size(
    breaks = seq(0, 24, 4)
  ) +
  labs(
    title = 'Peso do cérebro versus peso corporal',
    subtitle = paste0(
      'para mamíferos com menos de ',
      limite,
      ' kg'
    ),
    caption = 'Fonte: dataset `msleep`',
    x = 'Peso corporal (kg)',
  )
```

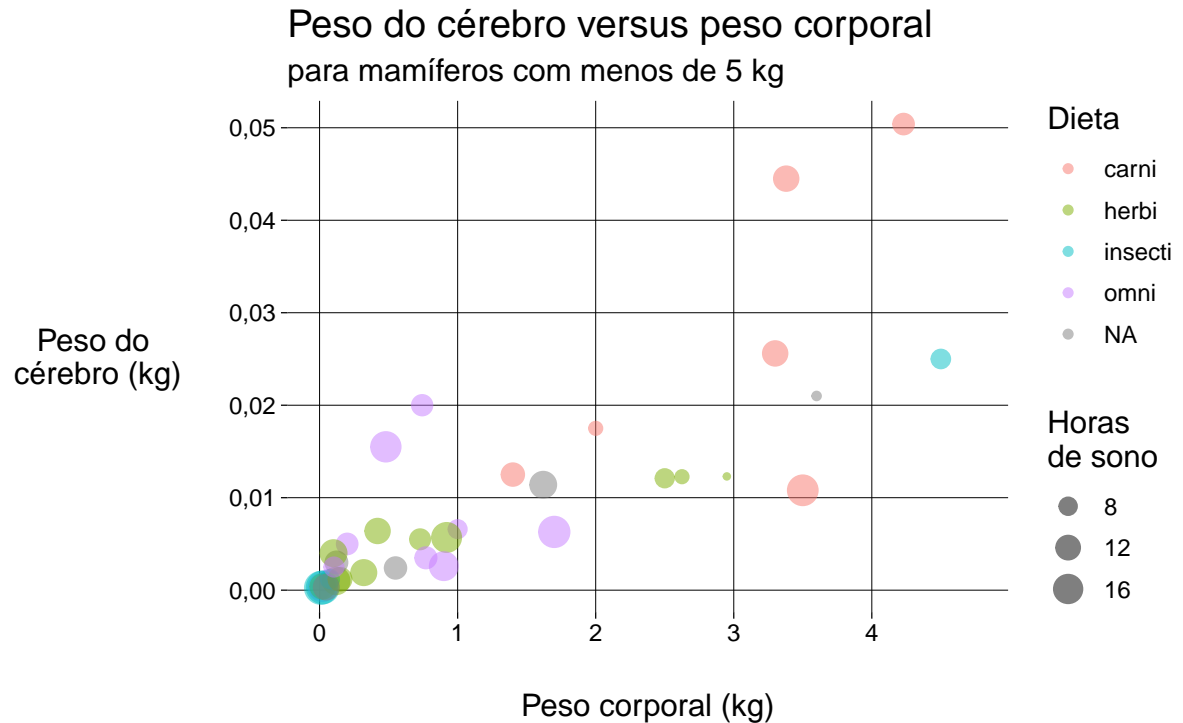
```

y = 'Peso do\n cérebro (kg)',
color = 'Dieta',
size = 'Horas\nde sono'
)

```

grafico

## Warning: Removed 18 rows containing missing values (geom\_point).



- Vamos mudar as cores usadas para a dieta, usando uma escala diferente.

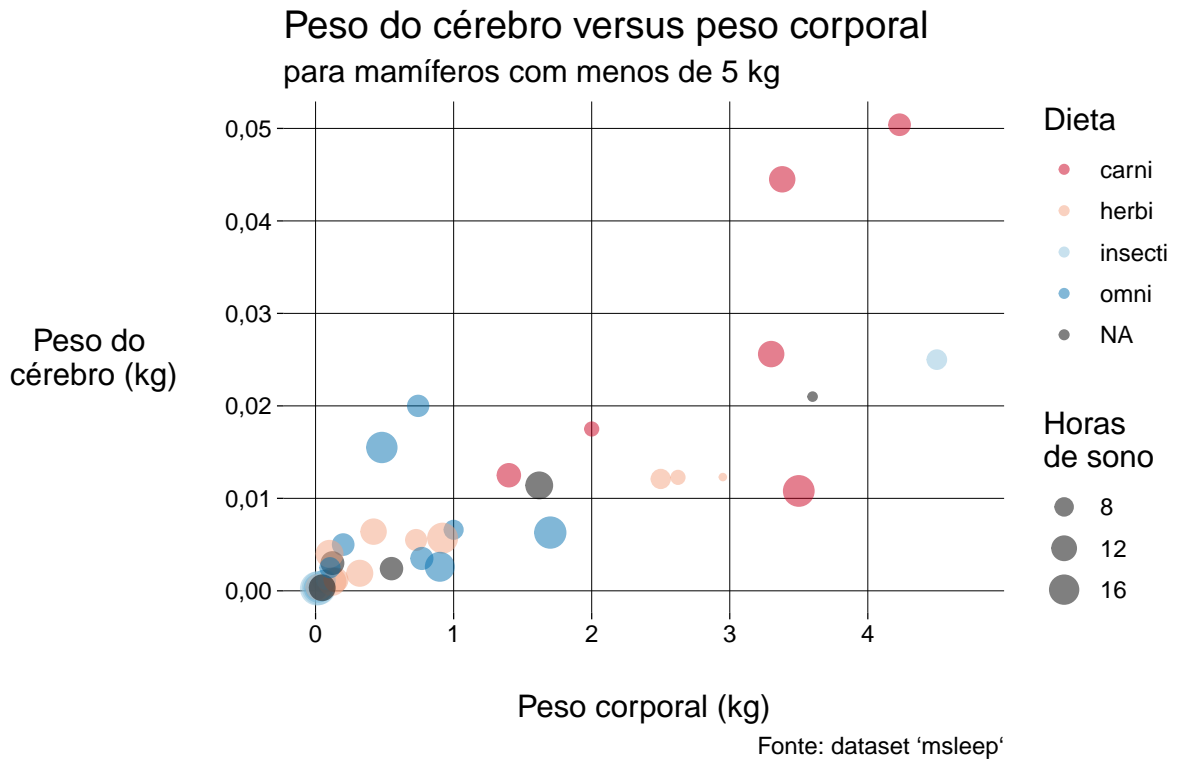
```

grafico2 <- grafico +
  scale_color_discrete(
    palette = 'RdBu',
    na.value = 'black',
    type = scale_color_brewer
  )

```

grafico2

## Warning: Removed 18 rows containing missing values (geom\_point).



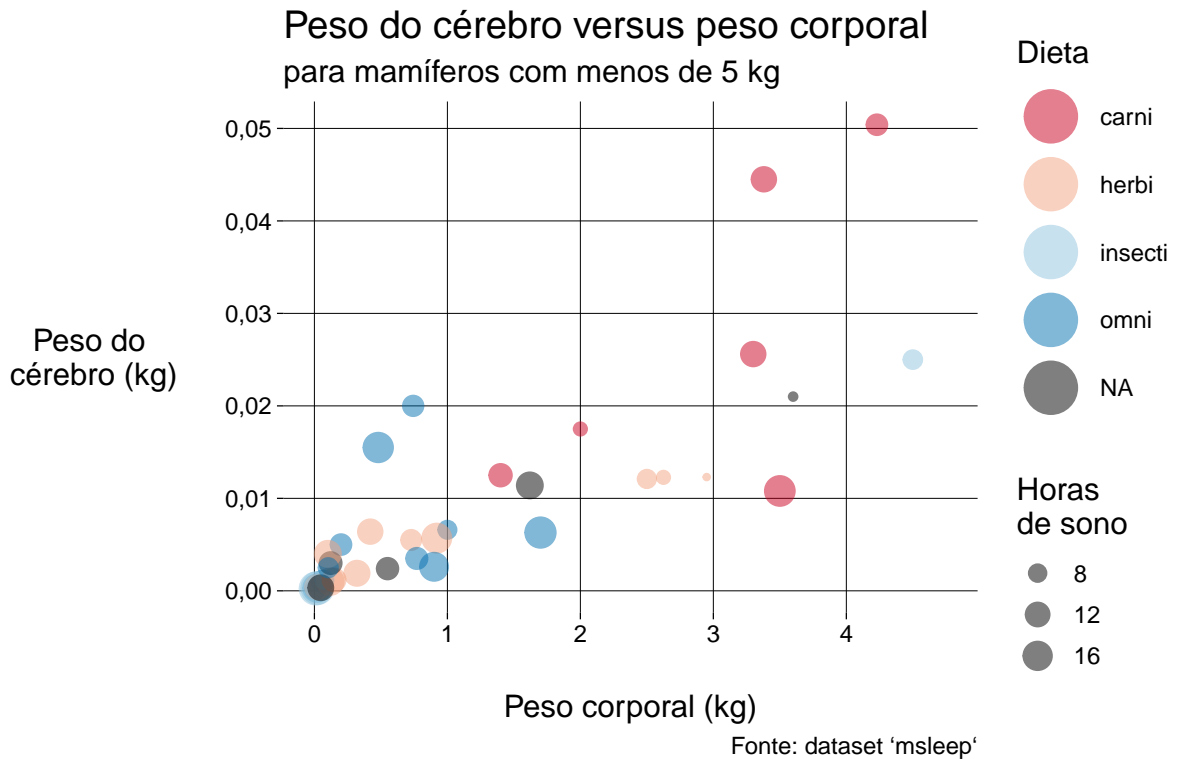
- Observe como usamos o gráfico já salvo na variável `grafico` e simplesmente acrescentamos a nova escala. Este tipo de “montagem” de gráficos `ggplot2` é bem conveniente, para evitar repetição de código.
- Um último ajuste na aparência: os pontos na legenda “Dieta” estão pequenos demais. Quase não identificamos as cores deles.

Vamos usar a função `guides` para modificar (*override*) a estética `color` — apenas na legenda, não nos pontos mostrados no gráfico, cujos tamanhos representam o número de horas de sono — tornando o tamanho maior. Leia mais sobre `override.aes` neste [link](https://ggplot2-book.org/scale-colour.html#guide_legend) (em inglês)<sup>1</sup>.

```
grafico3 <- grafico2 +
  guides(color = guide_legend(override.aes = list(size = 10)))

grafico3
## Warning: Removed 18 rows containing missing values (geom_point).
```

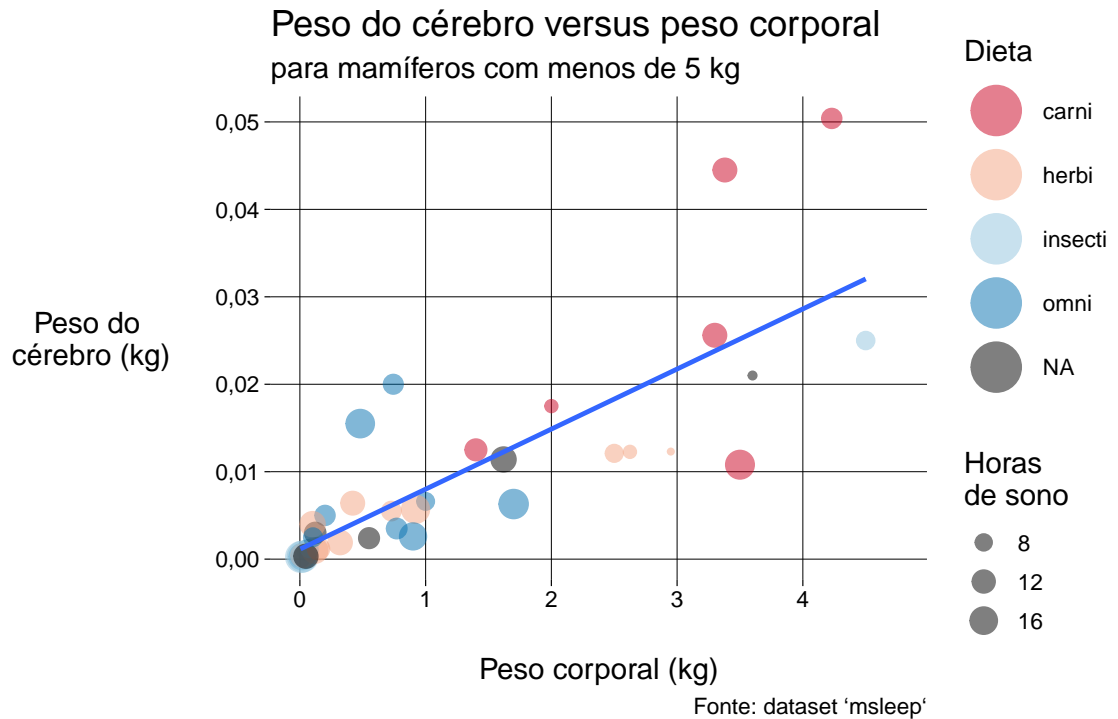
<sup>1</sup>[https://ggplot2-book.org/scale-colour.html#guide\\_legend](https://ggplot2-book.org/scale-colour.html#guide_legend)



- Agora podemos finalmente comentar sobre a informação que o gráfico mostra sobre os dados:
  - De fato, existe uma correlação entre peso cerebral e peso corporal: quanto maior o peso corporal, maior o peso cerebral. Nada surpreendente.
  - Podemos fazer o `ggplot2` traçar uma reta de regressão com a geometria `geom_smooth`. Vamos falar mais sobre correlação **em um capítulo futuro**.

```
grafico4 <- grafico3 +
  geom_smooth(
    aes(group = 1),
    show.legend = FALSE,
    method = 'lm',
    se = FALSE
  )

grafico4
## `geom_smooth()` using formula 'y ~ x'
## Warning: Removed 18 rows containing non-finite values (stat_smooth).
## Warning: Removed 18 rows containing missing values (geom_point).
```



- Todos os carnívoros têm peso corporal maior que 1kg e peso cerebral maior ou igual a 10g.
- Só um carnívoro dorme 8 horas ou menos. Qual?
- Todos os insetívoros — com exceção de um (qual?) — são muito leves e dormem muito.
- Todos os onívoros têm menos de 2kg de peso corporal e 20g ou menos de peso cerebral.

### 3.5

#### Vídeo 2

<https://youtu.be/c-LoZ9e8xWc>

### 3.6

#### Histogramas e cia.

- A idéia agora é agrupar indivíduos em classes, dependendo do valor de uma variável quantitativa.

### 3.6.1

#### Distribuições de frequência

- Vamos nos concentrar nas horas de sono.

```
sono$sleep_total
## [1] 12,1 17,0 14,4 14,9 4,0 14,4 8,7 7,0 10,1 3,0 5,3 9,4 10,0
## [14] 12,5 10,3 8,3 9,1 17,4 5,3 18,0 3,9 19,7 2,9 3,1 10,1 10,9
## [27] 14,9 12,5 9,8 1,9 2,7 6,2 6,3 8,0 9,5 3,3 19,4 10,1 14,2
## [40] 14,3 12,8 12,5 19,9 14,6 11,0 7,7 14,5 8,4 3,8 9,7 15,8 10,4
## [53] 13,5 9,4 10,3 11,0 11,5 13,7 3,5 5,6 11,1 18,1 5,4 13,0 8,7
## [66] 9,6 8,4 11,3 10,6 16,6 13,8 15,9 12,8 9,1 8,6 15,8 4,4 15,6
## [79] 8,9 5,2 6,3 12,5 9,8
```

- Antes de montar o histograma, vamos construir uma **distribuição de frequência**.
- A **amplitude** é a diferença entre o valor máximo e o valor mínimo. A função `range` não retorna a amplitude, mas sim os valores mínimo e máximo:

```
sono$sleep_total %>% range()
## [1] 1,9 19,9
```

- Vamos decidir que cada classe vai ter 2 horas. A função `cut` substitui os valores do vetor pelos nomes das classes:

```
sono$sleep_total %>%
  cut(breaks = seq(0, 20, 2), right = FALSE)
## [1] [12,14) [16,18) [14,16) [14,16) [4,6) [14,16) [8,10) [6,8)
## [9] [10,12) [2,4) [4,6) [8,10) [10,12) [12,14) [10,12) [8,10)
## [17] [8,10) [16,18) [4,6) [18,20) [2,4) [18,20) [2,4) [2,4)
## [25] [10,12) [10,12) [14,16) [12,14) [8,10) [0,2) [2,4) [6,8)
## [33] [6,8) [8,10) [8,10) [2,4) [18,20) [10,12) [14,16) [14,16)
## [41] [12,14) [12,14) [18,20) [14,16) [10,12) [6,8) [14,16) [8,10)
## [49] [2,4) [8,10) [14,16) [10,12) [12,14) [8,10) [10,12) [10,12)
## [57] [10,12) [12,14) [2,4) [4,6) [10,12) [18,20) [4,6) [12,14)
## [65] [8,10) [8,10) [8,10) [10,12) [10,12) [16,18) [12,14) [14,16)
## [73] [12,14) [8,10) [8,10) [14,16) [4,6) [14,16) [8,10) [4,6)
## [81] [6,8) [12,14) [8,10)
## 10 Levels: [0,2) [2,4) [4,6) [6,8) [8,10) [10,12) [12,14) ... [18,20)
```

- A função `table` faz a contagem dos elementos de cada classe:

```
sono$sleep_total %>%
  cut(breaks = seq(0, 20, 2), right = FALSE) %>%
  table(dnn = 'Horas de sono') %>%
  as.data.frame()
## # A tibble: 10 x 2
```



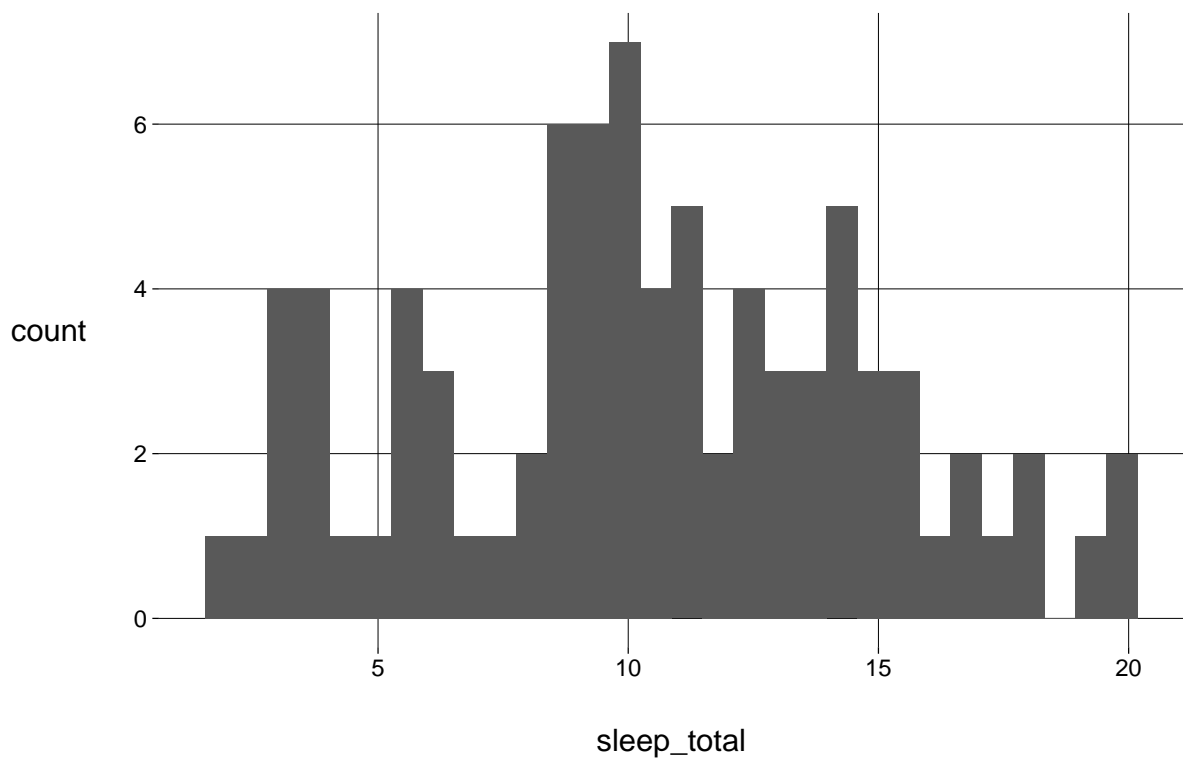
```
##   Horas.de.sono   Freq
##   <fct>         <int>
## 1 [0,2)          1
## 2 [2,4)          8
## 3 [4,6)          7
## 4 [6,8)          5
## 5 [8,10)         17
## 6 [10,12)        14
## # ... with 4 more rows
```

### 3.6.2

#### Histograma

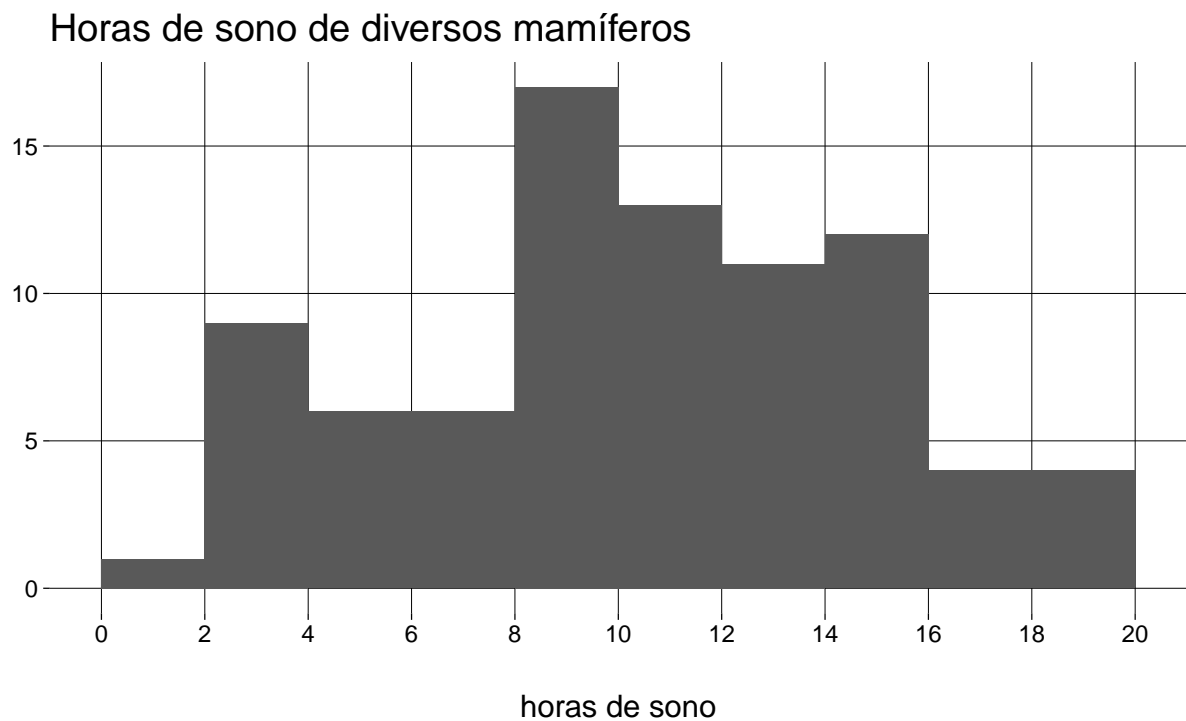
- Na verdade, o ggplot2 já faz esses cálculos para nós.
- O *default* é criar 30 classes (*bins*):

```
sono %>%
  ggplot(aes(x = sleep_total)) +
  geom_histogram()
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```



- Vamos mudar isto passando um vetor de limites das classes (*breaks*). Vamos acrescentar rótulos também:

```
sono %>%
  ggplot(aes(x = sleep_total)) +
    geom_histogram(breaks = seq(0, 20, 2)) +
    scale_x_continuous(breaks = seq(0, 20, 2)) +
    labs(
      title = 'Horas de sono de diversos mamíferos',
      x = 'horas de sono',
      y = NULL,
      caption = 'Fonte: dataset `msleep`'
    )
)
```



Fonte: dataset 'msleep'

- Nossas impressões:

- A classe que mais tem elementos é a de 8 a 10 horas.
- A distribuição é mais ou menos simétrica.
- A distribuição tem forma aproximada de sino: há poucos mamíferos com valores extremos de horas de sono; a maioria está próxima do valor médio:

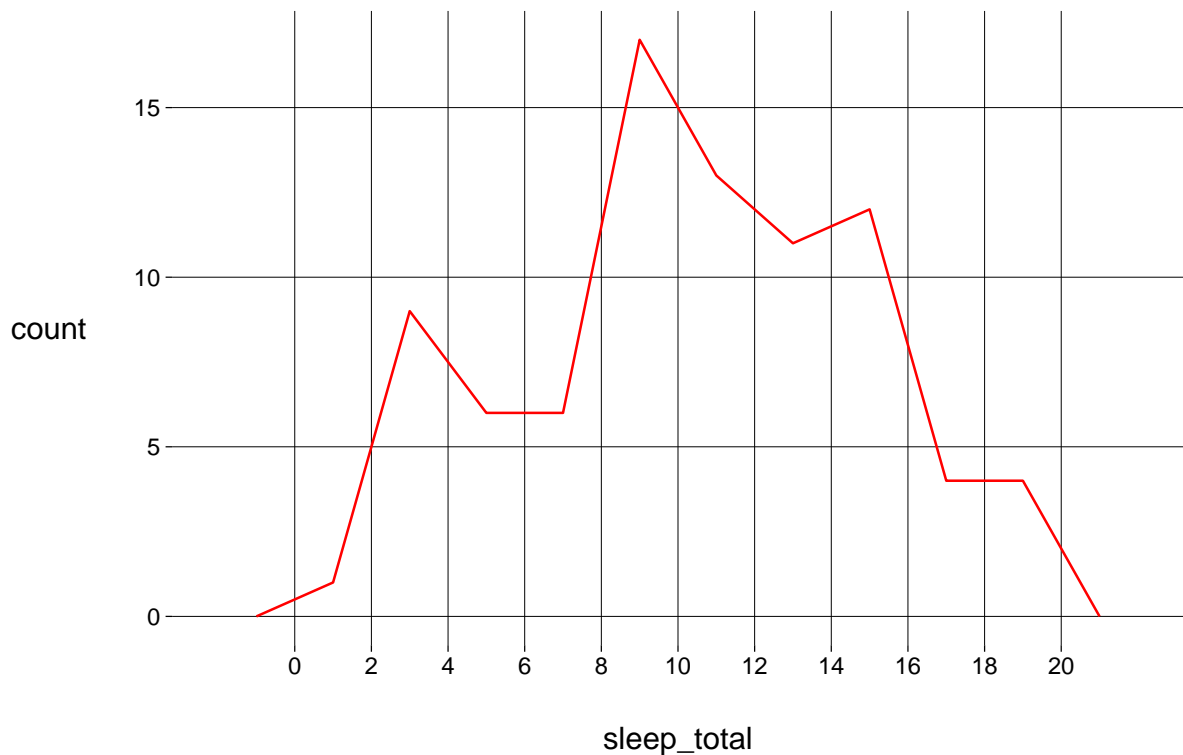
```
mean(sono$sleep_total)
## [1] 10,43373
```

### 3.6.3

#### Polígono de frequência

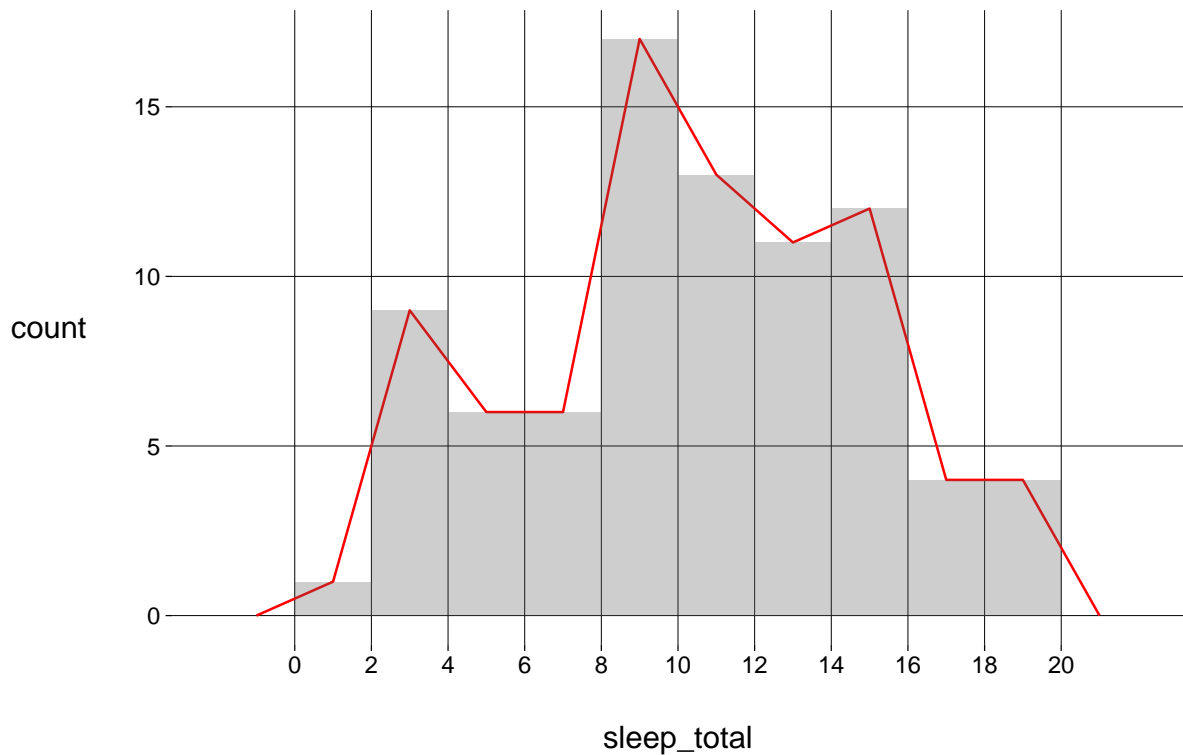
- Em vez das barras do histograma, podemos desenhar uma linha ligando seus topos.
- O resultado é um **polígono de frequência**.

```
pf <- sono %>%  
  ggplot(aes(x = sleep_total)) +  
    geom_freqpoly(breaks = seq(0, 20, 2), color = 'red') +  
    scale_x_continuous(breaks = seq(0, 20, 2))  
  
pf
```



- Vamos sobrepor o polígono de frequência ao histograma, para deixar claro o que está acontecendo:

```
pf + geom_histogram(breaks = seq(0, 20, 2), alpha = .3)
```

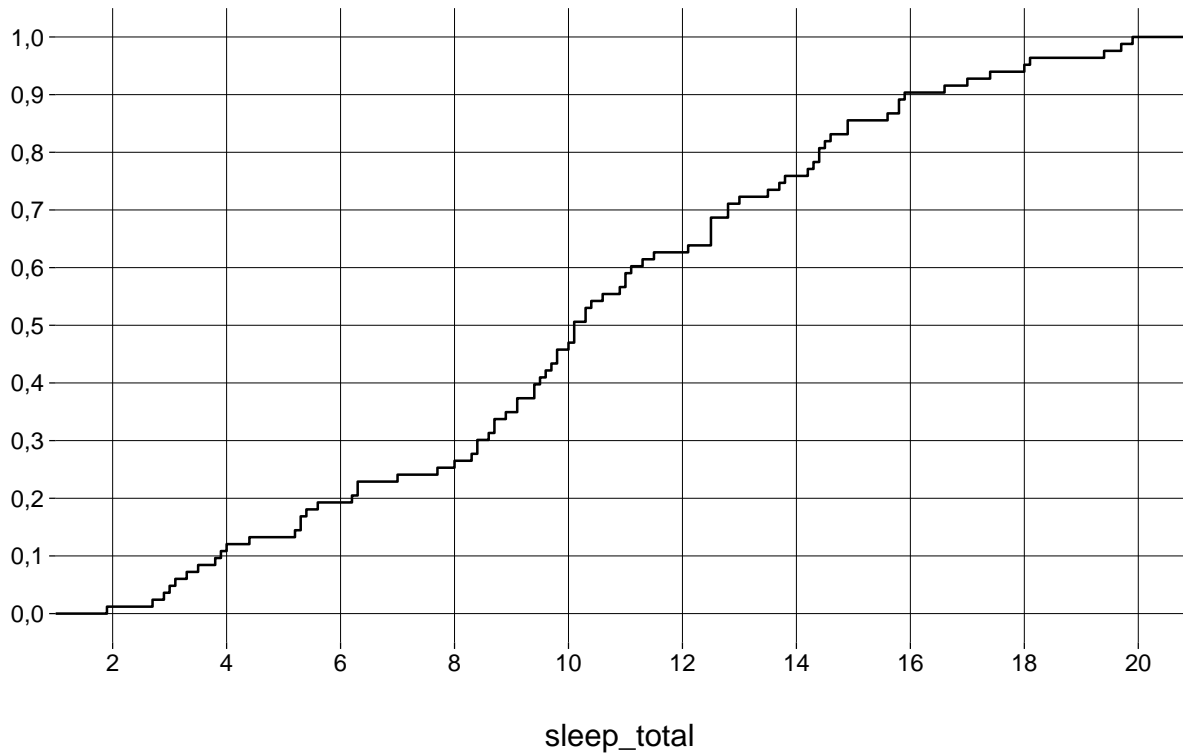


### 3.7

#### Ogiva

- A ogiva é um gráfico que mostra a **frequência acumulada**: para cada valor  $v$  da variável no eixo  $x$ , a proporção de indivíduos com valor menor ou igual a  $v$ .
- A geometria `geom_step` gera o gráfico de uma **função degrau**.
- Cada geometria está ligada a uma **stat**, um algoritmo para computar o que vai ser desenhado. Aqui, passamos para a geometria a função **`ecdf` (empirical cumulative distribution function)**, do pacote `stats`, que calcula as frequências acumuladas.

```
sono %>%
  ggplot(aes(x = sleep_total)) +
    geom_step(stat = 'ecdf') +
    scale_x_continuous(breaks = seq(0, 20, 2)) +
    scale_y_continuous(breaks = seq(0, 1, .1)) +
    labs(y = NULL)
```



- Com a ogiva, podemos obter informações difíceis de visualizar no histograma. Por exemplo:
  - Cerca de 20% dos mamíferos têm menos de 6 horas de sono.
  - Cerca de metade dos mamíferos têm menos de 10 horas de sono.
  - Cerca de 10% dos mamíferos têm mais de 16 horas de sono.

### 3.8

## Ramos e folhas

- No início dos anos 1900, quando estatísticas eram feitas à mão, Arthur Bowley criou os **diagramas de ramos e folhas**.
- Um diagrama de ramos e folhas é, basicamente, uma listagem de todos os valores de uma variável, agrupados de maneira que todos os valores de uma classe (i.e., de uma linha) têm os algarismos iniciais dentro de um intervalo.
- Para as horas de sono dos mamíferos:

```
sono$sleep_total %>%
  stem()
##
##   The decimal point is at the |
##
```

```
##    0 | 9
##    2 | 79013589
##    4 | 0423346
##    6 | 23307
##    8 | 03446779114456788
##   10 | 01113346900135
##   12 | 15555880578
##   14 | 234456996889
##   16 | 604
##   18 | 01479
```

- A primeira linha representa um indivíduo com 0,9 horas de sono.
- A penúltima linha representa 3 valores:
  - 16,6
  - 17,0
  - 17,4

### 3.9

## Exercícios



Não se esqueça de incluir títulos nos gráficos e rótulos nos eixos.

#### 3.9.1

### Peso cerebral e peso corporal

1. Observe os comandos que geraram o gráfico `grafico4`.
2. O que acontece se você retirar `aes(group = 1)` da chamada a `geom_smooth`? Explique.
3. O que acontece se você mudar `show.legend = FALSE` para `show.legend = TRUE` na chamada a `geom_smooth`? Explique.
4. O que acontece se você mudar `se = FALSE` para `se = TRUE` na chamada a `geom_smooth`? Explique.
5. Acrescente ao gráfico a camada `facet_wrap(~vore)`. O que acontece?
6. Examine o *data frame* `sono` e identifique o único insetívoro com mais de 4kg.
7. Instale o pacote `gg_repel` e acrescente ao gráfico `grafico4` (não facetado) a geometria `geom_label_repel` (consulte a ajuda) para rotular o mamífero insetívoro identificado no item anterior com o seu nome, **sem cobrir outros pontos do gráfico**. Cuidado para não alterar a legenda que já existe.

### 3.9.2

#### Peso cerebral e horas de sono

Use o *data frame* `sono` definido como

```
library(ggplot2)

sono <- msleep %>%
  select(
    name, order, genus, vore, bodywt,
    brainwt, awake, sleep_total
  )
```

1. Construa um histograma da variável `brainwt`. Escolha o número de classes que você achar melhor. O que acontece com os valores NA?
2. Descubra que função da forma `scale_x_...` usar<sup>2</sup> para fazer com que o eixo *x* tenha uma escala logarítmica. Gere um novo histograma.
3. Qual dos dois histogramas é melhor para responder a pergunta “Qual a faixa de peso cerebral que tem mais animais?” de forma satisfatória?
4. Construa um *scatter plot* de horas de sono versus peso do cérebro. Você percebe alguma correlação entre estas variáveis? Se precisar, concentre-se em um subconjunto dos dados.
5. Usando `geom_smooth` (leia a respeito<sup>3</sup>), sobreponha uma reta de regressão ao gráfico de dispersão, usando o método `lm` e sem o erro padrão (i.e., com `se = FALSE`). O que você observa? Discuta.

### 3.9.3

#### Igualdade de gênero entre furacões?

Este artigo<sup>4</sup> tenta achar uma relação entre o gênero do nome de um furacão e a quantidade de vítimas fatais provocadas por ele.

<sup>2</sup>[http://sillasgonzaga.com/material/curso\\_visualizacao/ggplot2-parte-ii.html#customizando-escalas](http://sillasgonzaga.com/material/curso_visualizacao/ggplot2-parte-ii.html#customizando-escalas)

<sup>3</sup><https://cdr.ibpad.com.br/ggplot2.html#objetos-geom%C3%A9tricos-e-tipos-de-gr%C3%A1ficos>

<sup>4</sup><https://www.pnas.org/content/111/24/8782>

Os dados estão no pacote DAAG, que deve ser instalado:

```
if (!require(DAAG))  
  install.packages("DAAG")
```

Vamos usar apenas algumas das variáveis, com nomes em português.

```
df <- hurricNamed %>%  
  as_tibble() %>%  
  transmute(  
    id = paste(Year, Name, sep = '-'),  
    nome = Name,  
    ano = Year,  
    velocidade = LF.WindsMPH * 1.8,      # convertido para km/h  
    pressao = LF.PressureMB,             # mbar  
    prejuizo = BaseDam2014 %>% round(), # milhões de dólares de 2014  
    mortes = deaths,  
    genero = mf  
  )
```

1. Crie histogramas para as seguintes variáveis, escolhendo a quantidade de barras que você achar melhor.

- velocidade
- prejuizo
- mortes

Não se esqueça de incluir títulos nos gráficos e rótulos nos eixos.

Comente os histogramas.

2. Os histogramas de prejuízos e mortes não ficaram bons. Vamos gerar histogramas transformados.

No *data frame*, crie duas novas colunas:

- logprejuizo: *logaritmo* do prejuízo (na base 10)
- logmortes: *logaritmo* do número de mortes (na base 10)

Agora, gere histogramas destas duas novas variáveis.

3. O que significa o valor do logaritmo do prejuízo na base 10?
4. O que significa o valor do logaritmo do número de mortes na base 10?
5. Por que o histograma do logaritmo do número de mortes vem com uma mensagem de aviso?
6. Por que isto não acontece com o logaritmo do prejuízo?
7. Faça um gráfico de dispersão com *pressao* no eixo *y* e *velocidade* no eixo *x*.



8. Usando `geom_smooth` (leia a respeito<sup>5</sup>), sobreponha uma reta de regressão ao gráfico, usando o método `lm` e sem o erro padrão (i.e., com `se = FALSE`). O que você observa? Discuta.
9. Faça um gráfico de dispersão com `logmortes` no eixo  $y$  e `pressao` no eixo  $x$ .
10. Usando `geom_smooth` (leia a respeito<sup>6</sup>), sobreponha uma reta de regressão ao gráfico, usando o método `lm` e sem o erro padrão (i.e., com `se = FALSE`). O que você observa? Discuta.
11. Faça um gráfico de dispersão com `logmortes` no eixo  $y$  e `pressao` no eixo  $x$ , com pontos coloridos de acordo com o gênero do nome do furacão.
12. Usando `geom_smooth` (leia a respeito<sup>7</sup>), sobreponha retas de regressão ao gráfico, uma para cada gênero, usando o método `lm` e sem o erro padrão (i.e., com `se = FALSE`). O que você observa? Discuta.



Visualizações como esta ajudam a explorar os dados, mas não servem para testar rigorosamente a hipótese de que furacões mulheres matam mais do que furacões homens.

Mais adiante no curso, vamos aprender a fazer testes mais rigorosos sobre hipóteses como esta.

---

<sup>5</sup><https://cdr.ibpad.com.br/ggplot2.html#objetos-geom%C3%A9tricos-e-tipos-de-gr%C3%A1ficos>

<sup>6</sup><https://cdr.ibpad.com.br/ggplot2.html#objetos-geom%C3%A9tricos-e-tipos-de-gr%C3%A1ficos>

<sup>7</sup><https://cdr.ibpad.com.br/ggplot2.html#objetos-geom%C3%A9tricos-e-tipos-de-gr%C3%A1ficos>

### Visualização com ggplot2 (continuação)

---



Busque mais informações sobre os pacotes `tidyverse` e `ggplot2` nas referências recomendadas.

#### 4.1

---

##### Vídeo 1

<https://youtu.be/TjgLDeIQHIc>

#### 4.2

---

##### *Boxplots*

##### 4.2.1

---

##### Conjunto de dados

- Vamos continuar a trabalhar com os dados sobre as horas de sono de alguns mamíferos:

```
sono <- msleep %>%
  select(name, vore, order, sleep_total)

sono
## # A tibble: 83 x 4
##   name                vore order      sleep_total
##   <chr>              <chr> <chr>      <dbl>
## 1 Cheetah            carni Carnivora      12.1
## 2 Owl monkey         omni  Primates       17
## 3 Mountain beaver    herbi Rodentia      14.4
## 4 Greater short-tailed shrew omni  Soricomorpha  14.9
## 5 Cow                herbi Artiodactyla    4
## 6 Three-toed sloth    herbi Pilosa      14.4
## # ... with 77 more rows
```

#### 4.2.2

#### Mediana e quartis

- Para entender *boxplots*, precisamos, antes, entender algumas medidas.
- Se tomarmos as quantidades de horas de sono de todos os animais do conjunto de dados e **classificarmos estas quantidades em ordem crescente**, vamos ter:

```
horas <- sono %>%
  pull(sleep_total) %>%
  sort()

horas
## [1] 1,9 2,7 2,9 3,0 3,1 3,3 3,5 3,8 3,9 4,0 4,4 5,2 5,3
## [14] 5,3 5,4 5,6 6,2 6,3 6,3 7,0 7,7 8,0 8,3 8,4 8,4 8,6
## [27] 8,7 8,7 8,9 9,1 9,1 9,4 9,4 9,5 9,6 9,7 9,8 9,8 10,0
## [40] 10,1 10,1 10,1 10,3 10,3 10,4 10,6 10,9 11,0 11,0 11,1 11,3 11,5
## [53] 12,1 12,5 12,5 12,5 12,5 12,8 12,8 13,0 13,5 13,7 13,8 14,2 14,3
## [66] 14,4 14,4 14,5 14,6 14,9 14,9 15,6 15,8 15,8 15,9 16,6 17,0 17,4
## [79] 18,0 18,1 19,4 19,7 19,9
```

- Quantos valores são?

```
length(horas)
## [1] 83
```

- O valor que está **bem no meio desta fila** — i.e., na posição 42 — é a **mediana**:

```
horas[ceiling(length(horas) / 2)]
## [1] 10,1
```

- Em R:

```
median(horas)
## [1] 10,1
```

Mediana e média são coisas muito diferentes.

Por acaso, neste exemplo, a média das horas é próxima da mediana:



```
mean(horas)
## [1] 10,43373
```

Isto costuma acontecer quando a distribuição dos dados é aproximadamente simétrica.

- Os **quartis** são os valores que estão nas posições  $\frac{1}{4}$ ,  $\frac{1}{2}$  e  $\frac{3}{4}$  da fila. São o **primeiro, segundo e terceiro quartis**, respectivamente.

```
horas[
  c(
    ceiling(length(horas) / 4),
    ceiling(length(horas) / 2),
    ceiling(3 * length(horas) / 4)
  )
]
## [1] 7,7 10,1 13,8
```

- **Sim, a mediana é o segundo quartil.**
- Em R, a **função quantile** generaliza esta ideia: dado um número  $q$  entre 0 e 1, o **quantil (com “N”)  $q$  é o elemento que está na posição que corresponde à fração  $q$  da fila ordenada.**

```
horas %>% quantile(c(.25, .5, .75))
## 25% 50% 75%
## 7,85 10,10 13,75
```

- Na verdade, R tem 9 algoritmos diferentes para calcular os quantis de uma amostra! Leia a ajuda da função `quantile` para conhecê-los.
- As diferenças entre nossos cálculos “à mão” e os resultados retornados por `quantile` são porque, em algumas situações, `quantile` calcula uma média ponderada entre elementos vizinhos. Por isso, `quantile` pode retornar valores que nem estão no vetor.
- Em R, a **função summary** mostra o **mínimo**, os **quartis (com “R”)**, a **média**, e o **máximo** de um vetor:

```
summary(horas)
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##      1,90    7,85   10,10   10,43   13,75   19,90
```

### 4.2.3

#### Média × mediana

- Vamos ver um exemplo simples para entender a diferença entre a média e a mediana.
- Imagine o seguinte vetor com as receitas mensais de algumas pessoas (em milhares de reais:)

```
receitas <- c(1, 2, 2, 3.5, 1, 4, 1)
```

- Eis a mediana e a média deste vetor:

```
summary(receitas)[c('Median', 'Mean')]
##      Median      Mean
## 2,000000 2,071429
```

- A mediana e a média são bem próximas.
- Imagine, agora, que adicionamos ao vetor um sujeito com receita mensal de 100 mil reais:

```
receitas <- c(1, 2, 2, 3.5, 1, 4, 1, 100)
```

- Eis a nova mediana e a nova média:

```
summary(receitas)[c('Median', 'Mean')]
##      Median      Mean
##      2,0000 14,3125
```

- O sujeito com a receita de 2 mil reais continua no meio da fila, mas a média (que é a soma de todas as receitas, dividida pelo número de indivíduos) ficou muito diferente.
- A receita do novo sujeito é um **valor discrepante**, ou, em inglês, um **outlier**.

#### Conclusão:



A **mediana é robusta**, pouco afetada por *outliers*.

A **média é pouco robusta**, muito sensível a *outliers*.

#### 4.2.4

### Intervalo interquartil (IQR) e *outliers*

- Qual fração dos elementos está entre o primeiro e o terceiro quartis?

```
length(
  horas[between(horas, quantile(horas, .25), quantile(horas, .75))]
) /
length(
  horas
)
## [1] 0,4939759
```

- Metade do total de elementos está entre o primeiro e o terceiro quartis.
- Este é o chamado intervalo interquartil (*interquartile range*, em inglês).
- No nosso vetor `horas`, os limites do IQR são

```
quantile(horas, c(.25, .75))
##      25%      75%
##  7,85 13,75
```

- O comprimento deste intervalo é calculado pela função `IQR`:

```
IQR(horas)
## [1] 5,9
```

- Valores muito abaixo do primeiro quartil podem ser considerados discrepantes (*outliers*), mas quão abaixo?
- A resposta (puramente convencional) é  $1,5 \times \text{IQR}$  abaixo do primeiro quartil.
- No nosso vetor `horas`, isto significa valores abaixo de

```
limite_inferior <- quantile(horas, .25) - 1.5 * IQR(horas)

unnamed(limite_inferior)
## [1] -1
```

- Neste caso, não há *outliers*:

```
horas[horas < limite_inferior]
## numeric(0)
```

- Da mesma forma, valores muito acima do terceiro quartil podem ser considerados discrepantes (*outliers*), mas quão acima?
- De novo, a resposta (puramente convencional) é  $1,5 \times \text{IQR}$  acima do terceiro quartil.

- No nosso vetor `horas`, isto significa valores acima de

```
limite_superior <- quantile(horas, .75) + 1.5 * IQR(horas)

unnname(limite_superior)
## [1] 22,6
```

- Neste caso, também não há *outliers*:

```
horas[horas > limite_superior]
## numeric(0)
```

- Outro exemplo: vamos tomar apenas os mamíferos onívoros:

```
onivoros <- sono %>%
  filter(vore == 'omni')

onivoros
## # A tibble: 20 x 4
##   name                vore order      sleep_total
##   <chr>              <chr> <chr>          <dbl>
## 1 Owl monkey        omni  Primates        17
## 2 Greater short-tailed shrew omni  Soricomorpha    14.9
## 3 Grivet            omni  Primates        10
## 4 Star-nosed mole    omni  Soricomorpha    10.3
## 5 African giant pouched rat omni  Rodentia         8.3
## 6 Lesser short-tailed shrew omni  Soricomorpha     9.1
## # ... with 14 more rows
```

- Vamos extrair o vetor de horas de sono:

```
horas <- onivoros %>%
  pull(sleep_total)

horas
## [1] 17,0 14,9 10,0 10,3 8,3 9,1 18,0 10,1 10,9 9,8 8,0 10,1 9,7
## [14] 9,4 11,0 8,7 9,6 9,1 15,6 8,9
```

- Vamos calcular o primeiro e terceiro quartis:

```
quartis <- horas %>%
  quantile(c(.25, .75))

quartis
## 25% 75%
## 9,100 10,925
```

- Vamos achar o IQR:

```
IQR(horas)
## [1] 1,825
```

- E os limites a partir dos quais os valores são *outliers*:

```
limites <- quartis + c(-1, 1) * 1.5 * IQR(horas)

unnname(limites)
## [1] 6,3625 13,6625
```

- Existem *outliers* inferiores?

```
onivoros %>%
  filter(sleep_total < limites[1])
## # A tibble: 0 x 4
## # ... with 4 variables: name <chr>, vore <chr>, order <chr>,
## #   sleep_total <dbl>
```

Não.

- Existem *outliers* superiores?

```
onivoros %>%
  filter(sleep_total > limites[2])
## # A tibble: 4 x 4
##   name                vore order      sleep_total
##   <chr>              <chr> <chr>         <dbl>
## 1 Owl monkey         omni  Primates         17
## 2 Greater short-tailed shrew omni  Soricomorpha    14.9
## 3 North American Opossum  omni  Didelphimorphia  18
## 4 Tenrec              omni  Afrosoricida     15.6
```

Sim! Estes animais dormem demais em comparação com os outros onívoros.

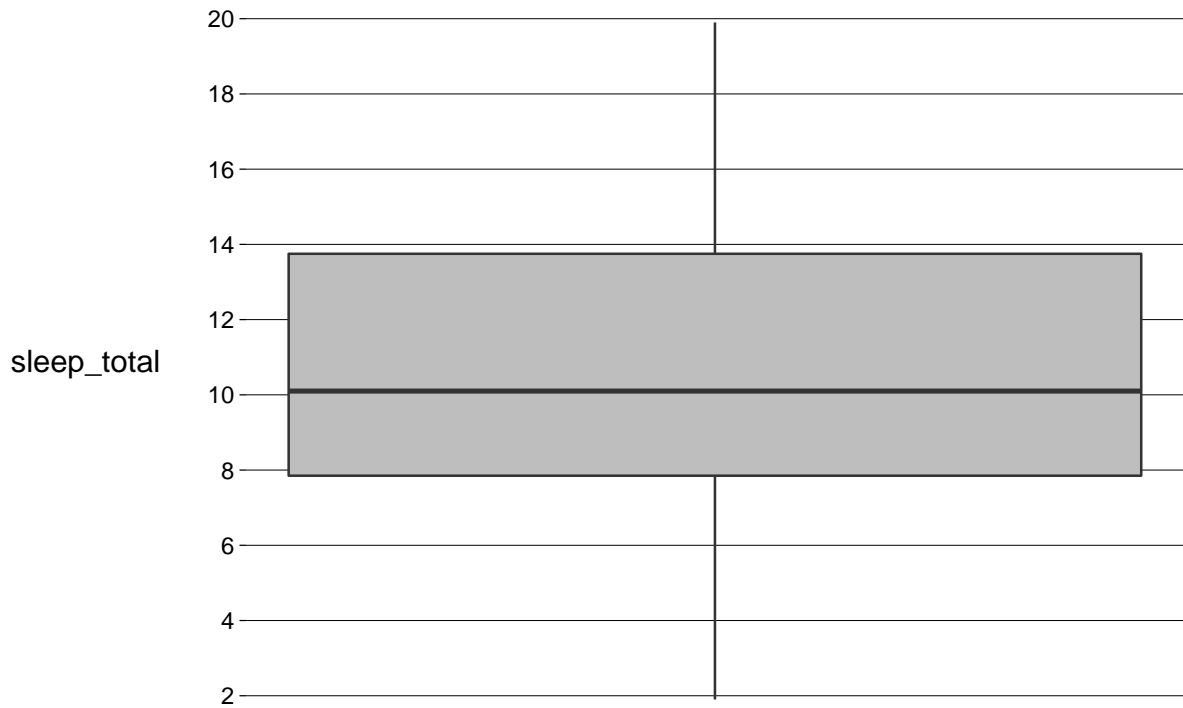
## 4.2.5

### Gerando boxplots

- Um *boxplot* é uma representação visual dos valores que calculamos acima.
- No `ggplot2`, a geometria `geom_boxplot` constrói *boxplots*:

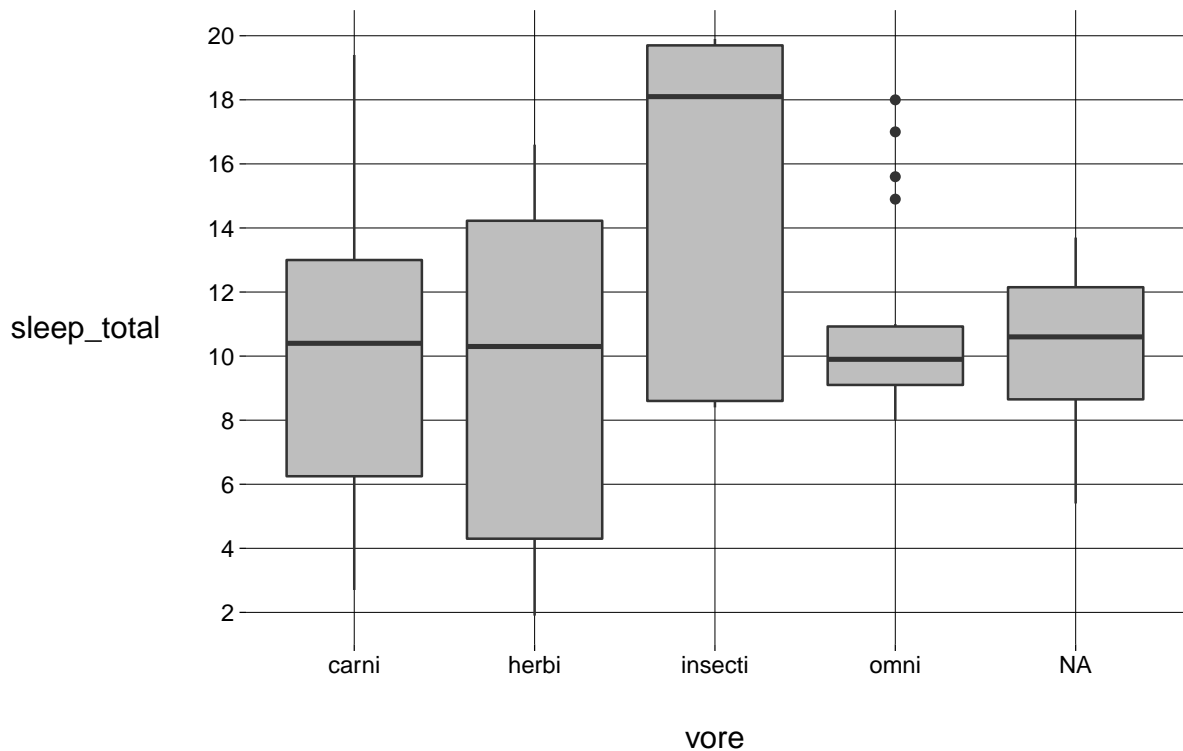
```
sono %>%
  ggplot(aes(y = sleep_total)) +
  geom_boxplot(fill = 'gray') +
  scale_x_continuous(breaks = NULL) +
  scale_y_continuous(breaks = seq(0, 20, 2))
```





- A **caixa** vai do valor do **primeiro quartil** (embaixo) até o **terceiro quartil** (em cima).
- A **linha horizontal dentro da caixa** representa o valor da **mediana**.
- As **linhas verticais** acima e abaixo da caixa (pitorescamente chamadas de “bigodes”) vão até o **limite inferior** (primeiro quartil  $- 1,5 \times \text{IQR}$ ) e até o **limite superior** (terceiro quartil  $+ 1,5 \times \text{IQR}$ ).
- Neste *boxplot*, não há *outliers*.
- Podemos usar a posição *x* para desenhar vários *boxplots*, um para cada dieta:

```
sono %>%  
  ggplot(aes(x = vore, y = sleep_total)) +  
    geom_boxplot(fill = 'gray') +  
    scale_y_continuous(breaks = seq(0, 20, 2))
```



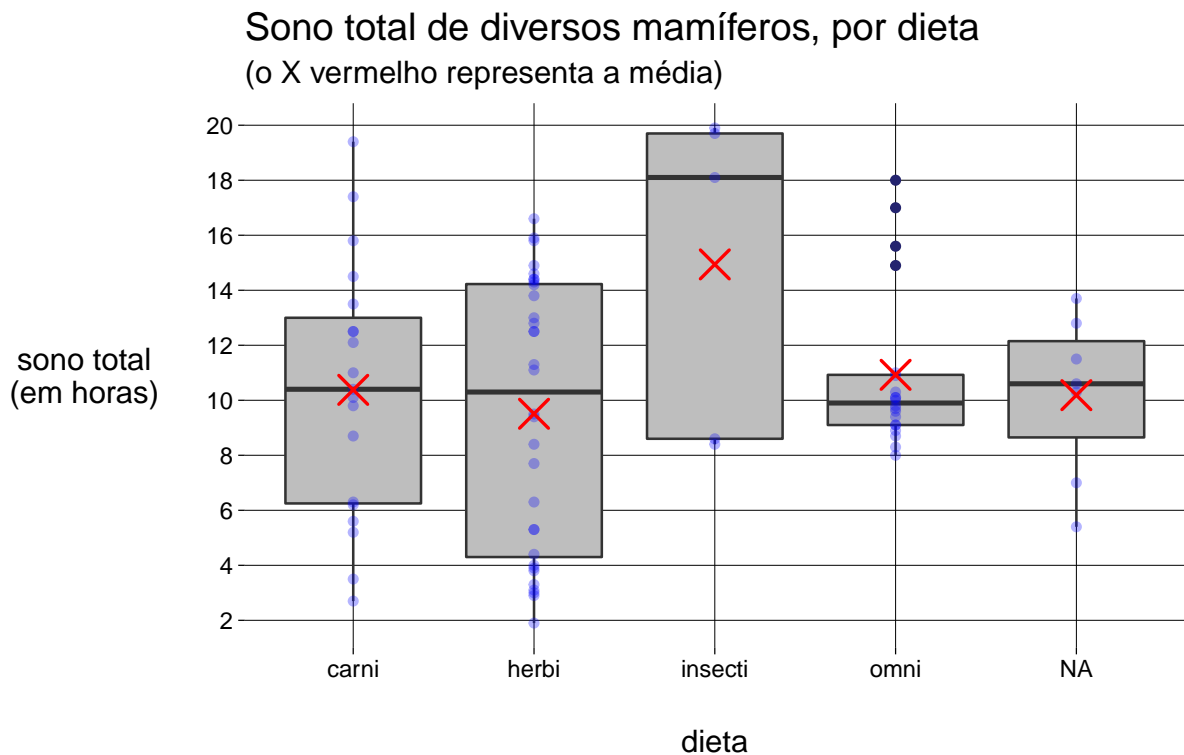
- No *boxplot* de onívoros, **os outliers aparecem como pontos isolados**, acima da caixa, além dos alcances do bigode superior (aliás, onde está bigode superior?).
- *Boxplots* lado a lado são úteis para compararmos grupos diferentes de dados.
- Veja como, com exceção dos insetívoros, as medianas dos grupos são parecidas.
- Veja como carnívoros, insetívoros e herbívoros apresentam maior variação, enquanto onívoros e animais sem dieta registrada apresentam menor variação.
- Vamos combinar, em um só gráfico
  - Os pontos representando os animais,
  - Os *boxplots*,
  - As médias (que podem estar próximas ou distantes das medianas).

```
sono %>%
  ggplot(aes(x = vore, y = sleep_total)) +
    geom_boxplot(fill = 'gray') +
    scale_y_continuous(breaks = seq(0, 20, 2)) +
    geom_point(
      color = 'blue',
      alpha = .3
    ) +
    stat_summary(
      fun = mean,
      geom = 'point',
```

```

color = 'red',
shape = 'cross',
size = 5,
stroke = 1
) +
labs(
  title = 'Sono total de diversos mamíferos, por dieta',
  subtitle = '(o X vermelho representa a média)',
  x = 'dieta',
  y = 'sono total\n(em horas)'
)

```



- Quando a caixa é longa, o IQR é grande, e os valores estão muito espalhados; é o caso dos herbívoros e insetívoros.
- Quando a caixa é curta, o IQR é pequeno, e os valores estão pouco espalhados; é o caso dos onívoros. Como o IQR é pequeno, os 4 mamíferos com mais de 14 horas de sono são *outliers*.
- Observe, ainda, como os *outliers* “puxam” a média dos onívoros para cima.

## 4.3

### Vídeo 2

<https://youtu.be/QqnOvgBXJ-s>

## 4.4

### Gráficos de barras e de colunas

#### 4.4.1

##### Dataset

```
HairEyeColor
## , , Sex = Male
##
##      Eye
## Hair   Brown Blue Hazel Green
## Black   32   11   10    3
## Brown   53   50   25   15
## Red     10   10    7    7
## Blond    3   30    5    8
##
## , , Sex = Female
##
##      Eye
## Hair   Brown Blue Hazel Green
## Black   36    9    5    2
## Brown   66   34   29   14
## Red     16    7    7    7
## Blond    4   64    5    8
```

```
df_orig <- as.data.frame(HairEyeColor) %>%
  uncount(Freq) %>%
  as_tibble()
```

```
df_orig %>% dfSummary() %>% print()
```

Variável	Estatísticas / Valores	Freqs (% de Válidos)	Faltante
Hair [factor]	1. Black	108 (18,2%)	0
	2. Brown	286 (48,3%)	(0,0%)
	3. Red	71 (12,0%)	
	4. Blond	127 (21,5%)	

Variável	Estatísticas / Valores	Freqs (% de Válidos)	Faltante
Eye [factor]	1. Brown	220 (37,2%)	0
	2. Blue	215 (36,3%)	(0,0%)
	3. Hazel	93 (15,7%)	
	4. Green	64 (10,8%)	
Sex [factor]	1. Male	279 (47,1%)	0
	2. Female	313 (52,9%)	(0,0%)

```
cabelo <- c(
  'Brown' = 'Castanhos',
  'Blond' = 'Louros',
  'Black' = 'Pretos',
  'Red' = 'Ruivos'
)

olhos <- c(
  'Brown' = 'Castanhos',
  'Blue' = 'Azuis',
  'Hazel' = 'Avelã',
  'Green' = 'Verdes'
)

sexo <- c(
  'Male' = 'Homem',
  'Female' = 'Mulher'
)

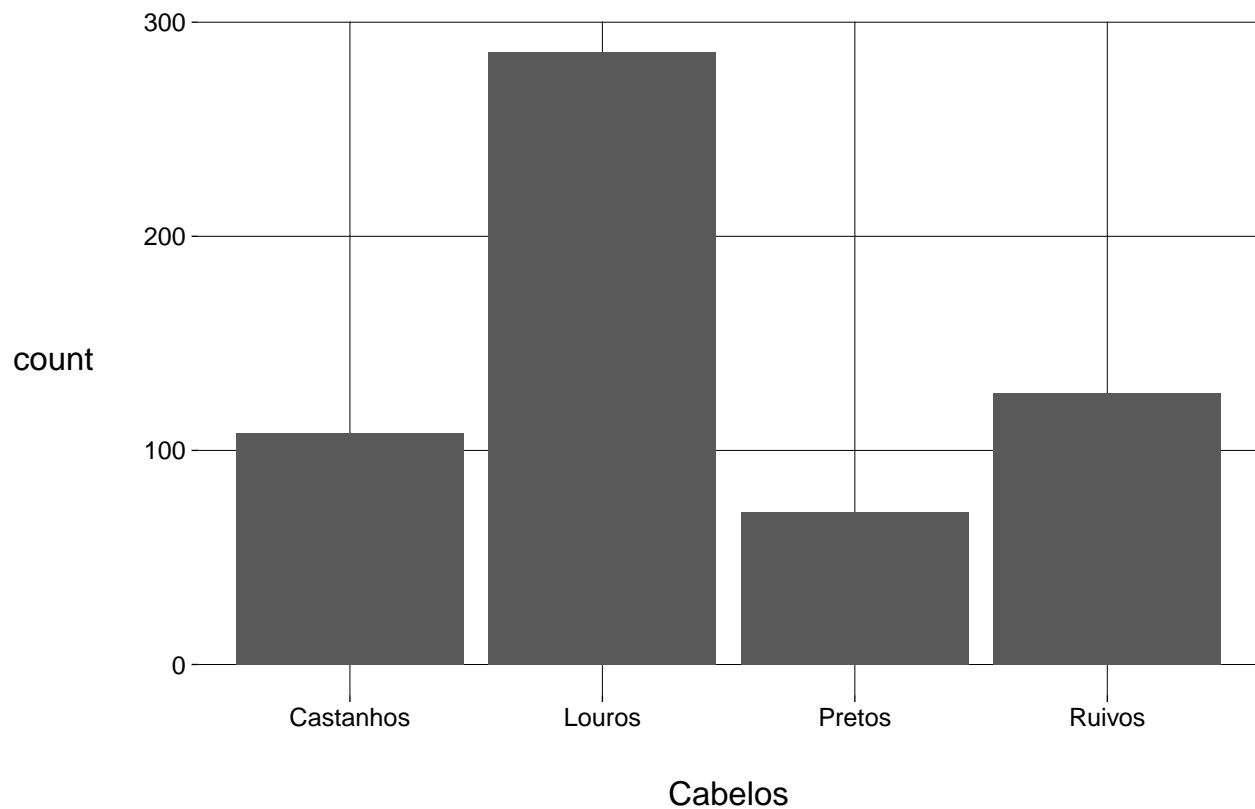
df <- df_orig %>%
  transmute(
    Cabelos = cabelo[Hair],
    Olhos = olhos[Eye],
    Sexo = sexo[Sex]
  )
```

```
df %>% dfSummary() %>% print()
```

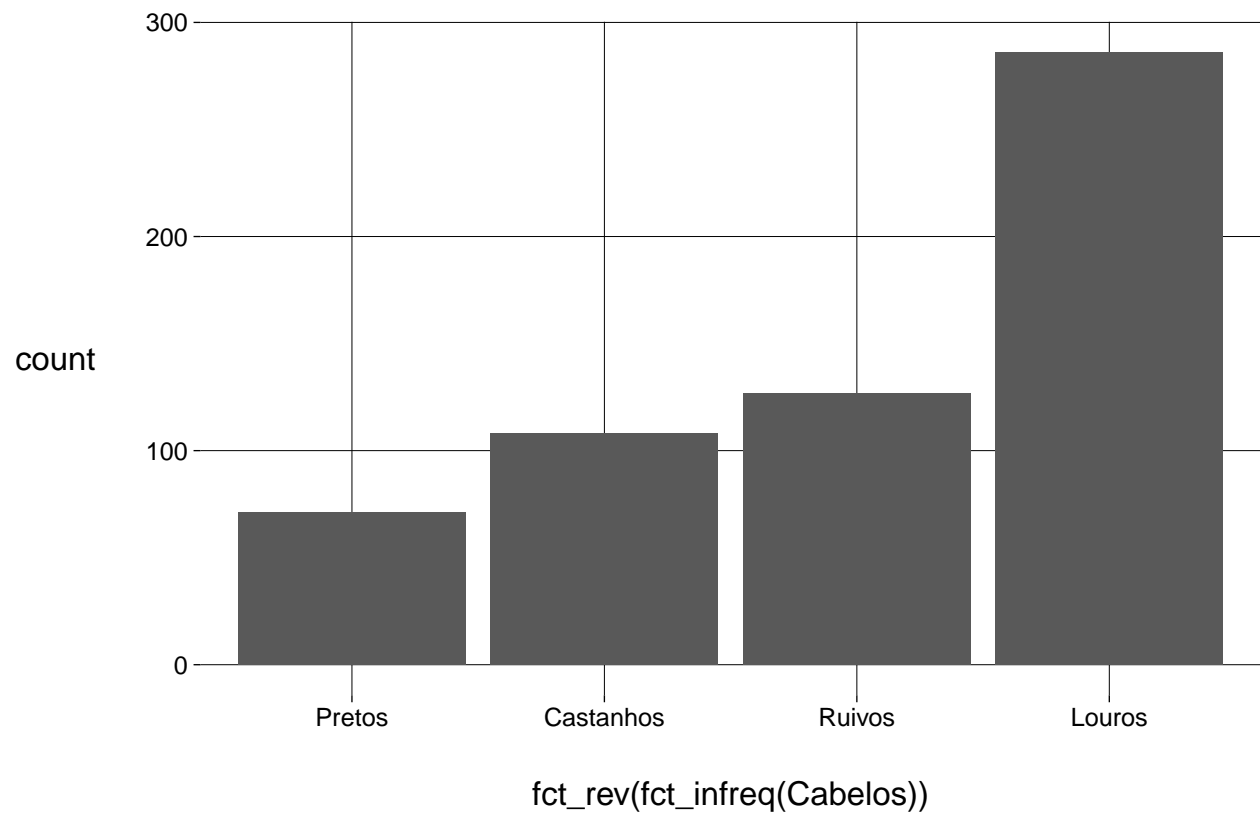
Variável	Estatísticas / Valores	Freqs (% de Válidos)	Faltante
Cabelos [character]	1. Castanhos	108 (18,2%)	0
	2. Louros	286 (48,3%)	(0,0%)
	3. Pretos	71 (12,0%)	
	4. Ruivos	127 (21,5%)	
Olhos [character]	1. Avelã	93 (15,7%)	0
	2. Azuis	215 (36,3%)	(0,0%)
	3. Castanhos	220 (37,2%)	
	4. Verdes	64 (10,8%)	

Variável	Estatísticas / Valores	Freqs (% de Válidos)	Faltante
Sexo [character]	1. Homem	279 (47,1%)	0
	2. Mulher	313 (52,9%)	(0,0%)

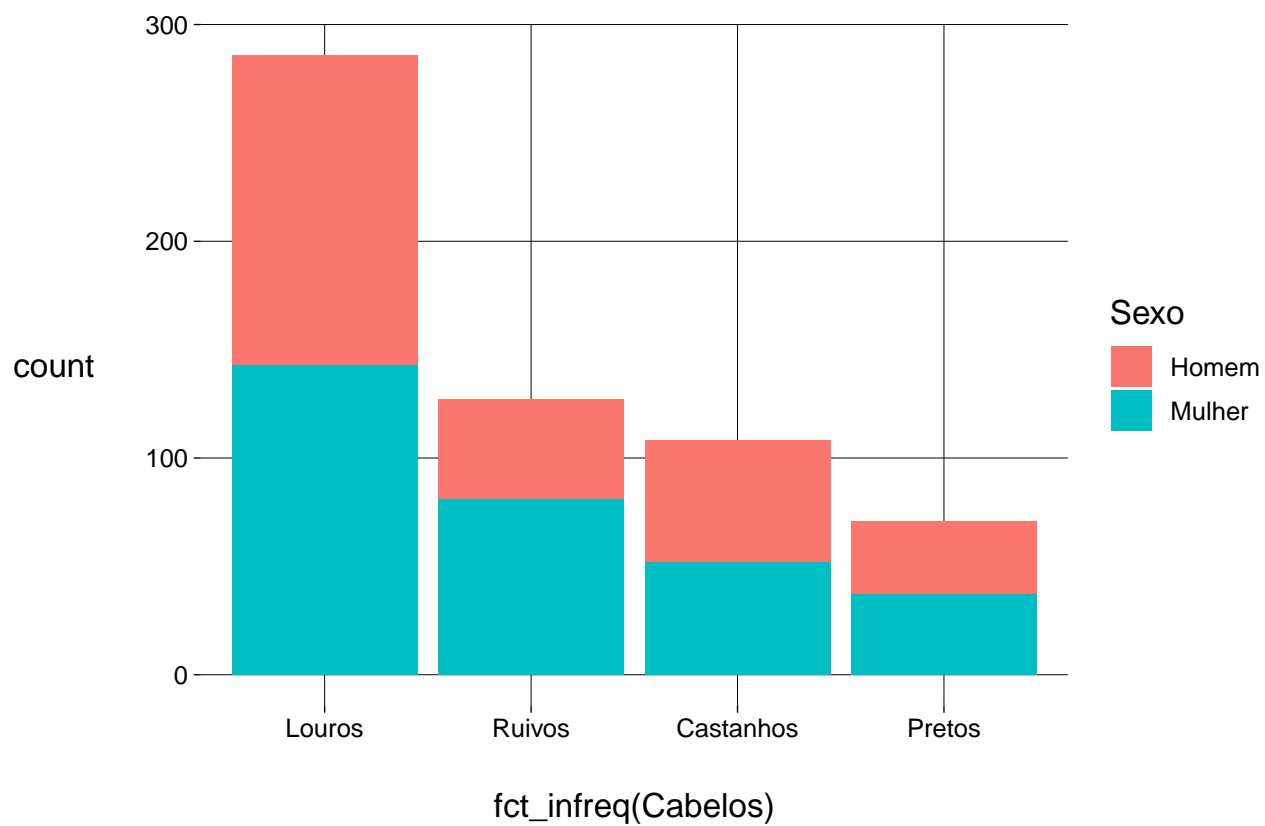
```
df %>%
  ggplot(aes(x = Cabelos)) +
    geom_bar()
```



```
df %>%
  ggplot(aes(x = fct_rev(fct_infreq(Cabelos)))) +
    geom_bar()
```

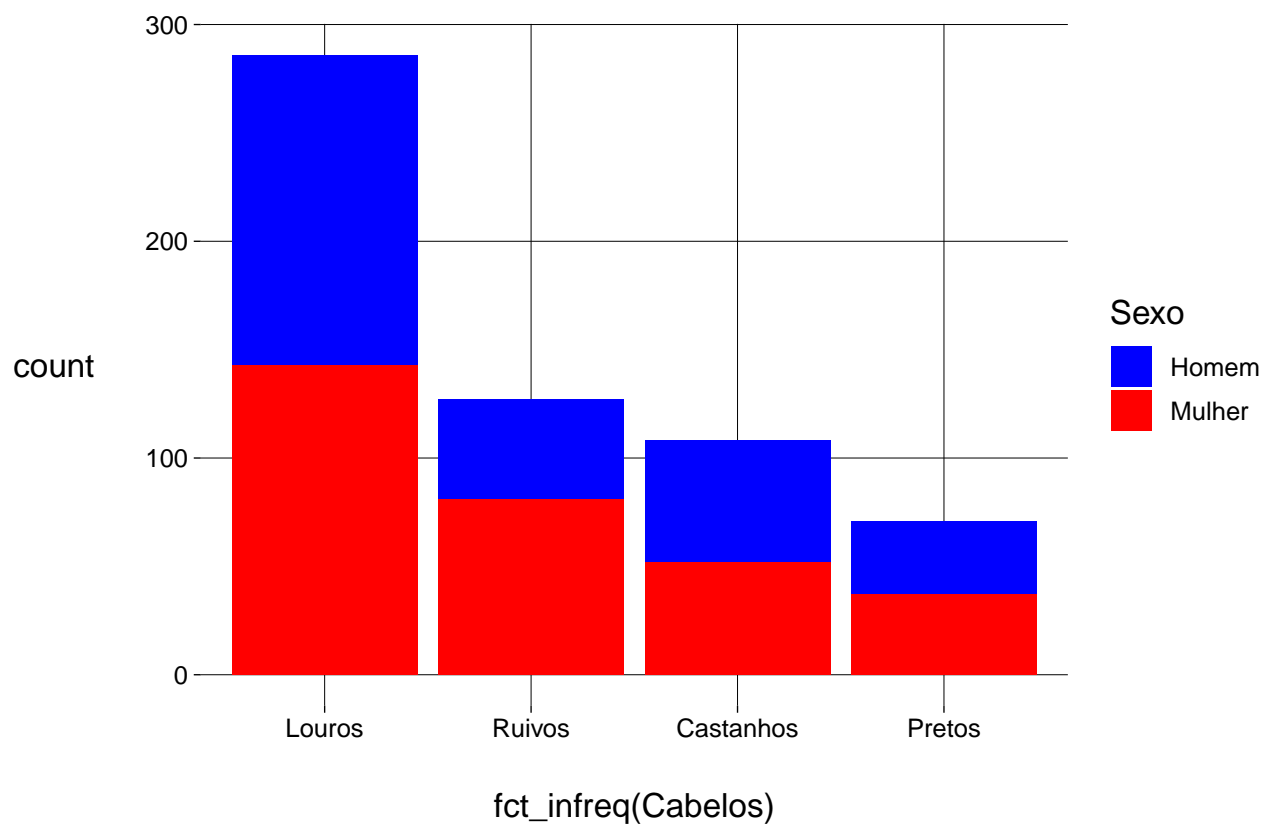


```
df %>%  
  ggplot(aes(x = fct_infreq(Cabelos), fill = Sexo)) +  
  geom_bar()
```

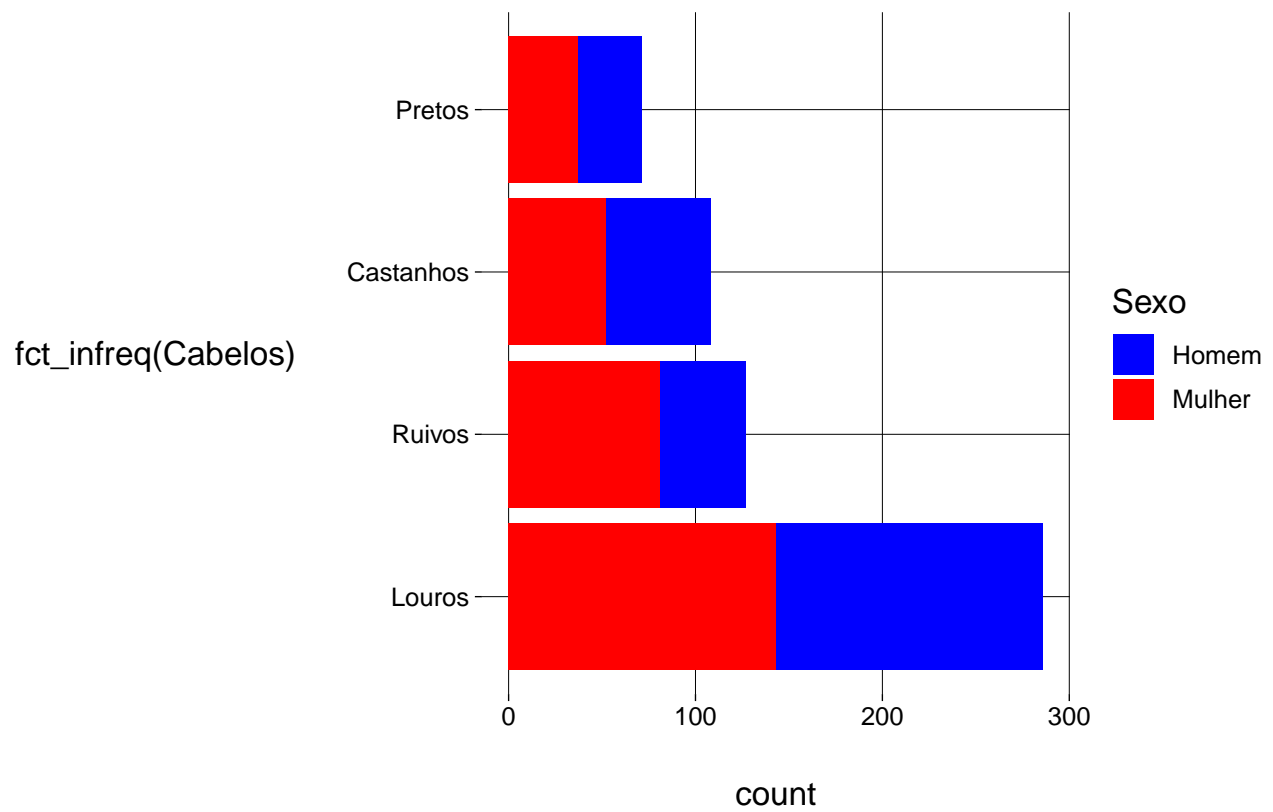


```
df %>%  
  ggplot(aes(x = fct_infreq(Cabelos), fill = Sexo)) +  
    geom_bar() +  
    scale_fill_discrete(type = c('blue', 'red'))
```

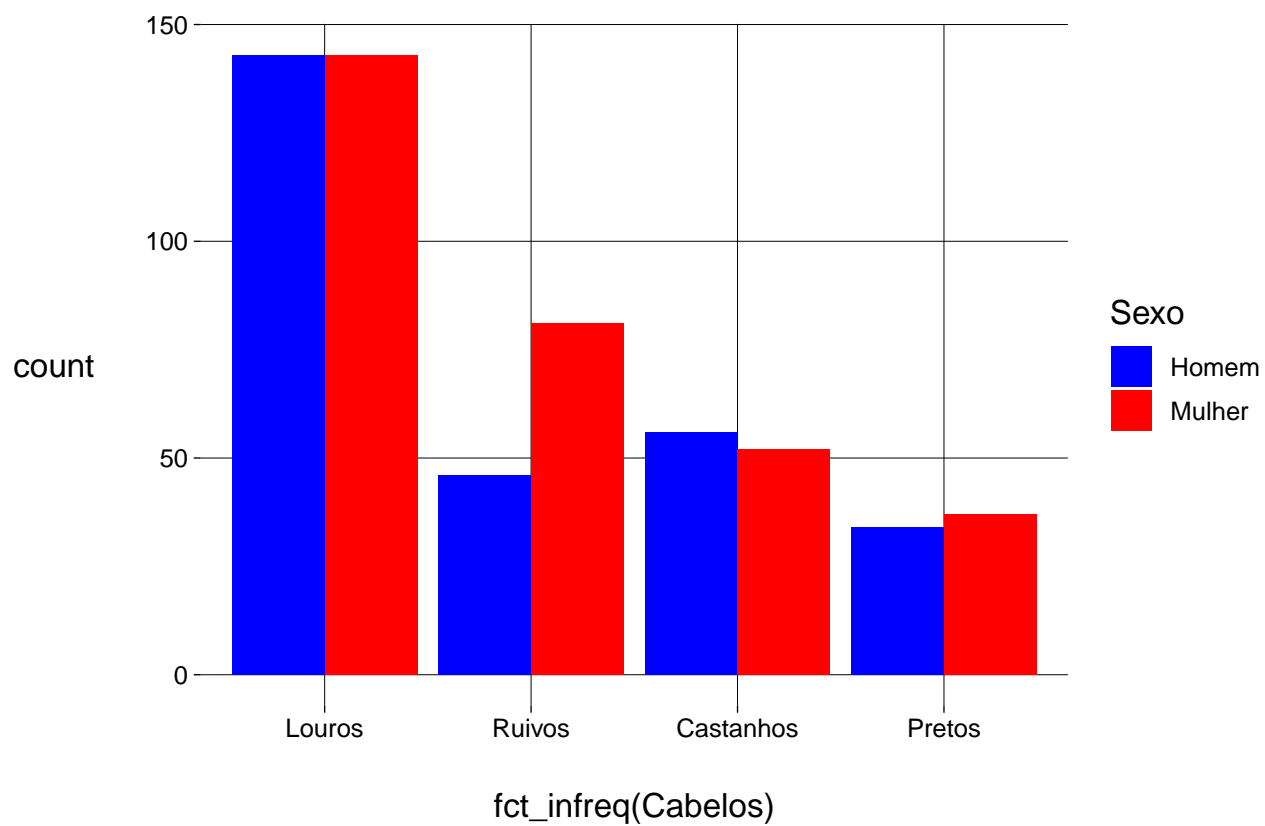




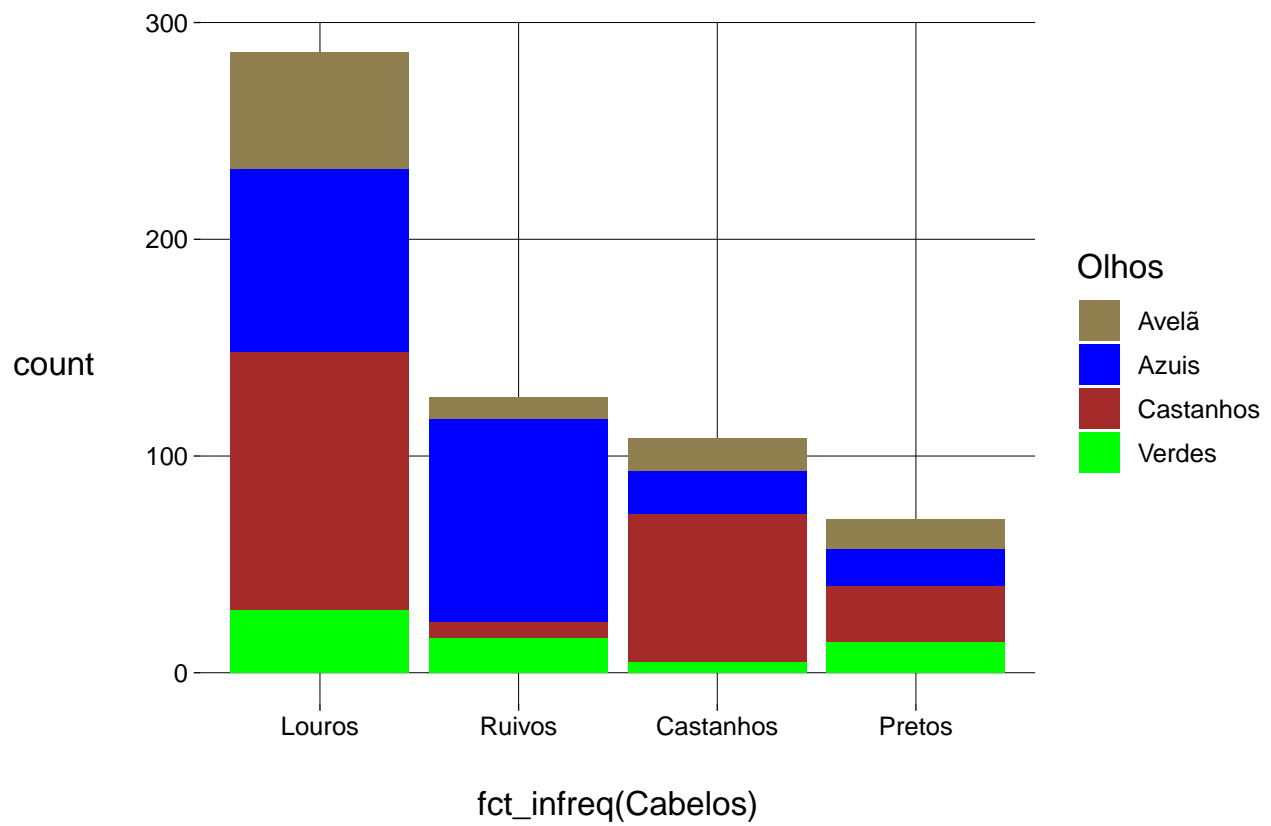
```
df %>%  
  ggplot(aes(x = fct_infreq(Cabelos), fill = Sexo)) +  
    geom_bar() +  
    scale_fill_discrete(type = c('blue', 'red')) +  
    coord_flip()
```



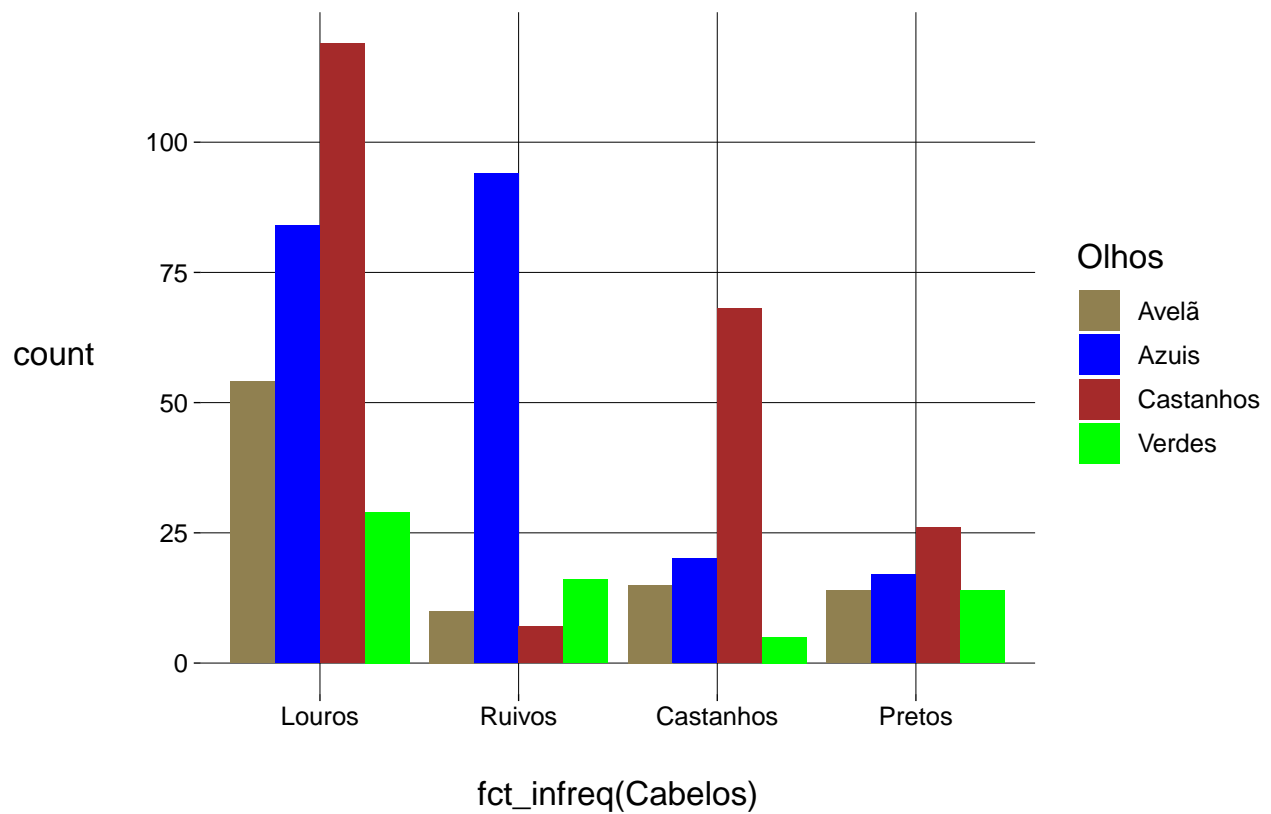
```
df %>%  
  ggplot(aes(x = fct_infreq(Cabelos), fill = Sexo)) +  
  geom_bar(position = 'dodge') +  
  scale_fill_discrete(type = c('blue', 'red'))
```



```
df %>%  
  ggplot(aes(x = fct_infreq(Cabelos), fill = Olhos)) +  
    geom_bar() +  
    scale_fill_discrete(type = c('#908050', 'blue', 'brown', 'green'))
```

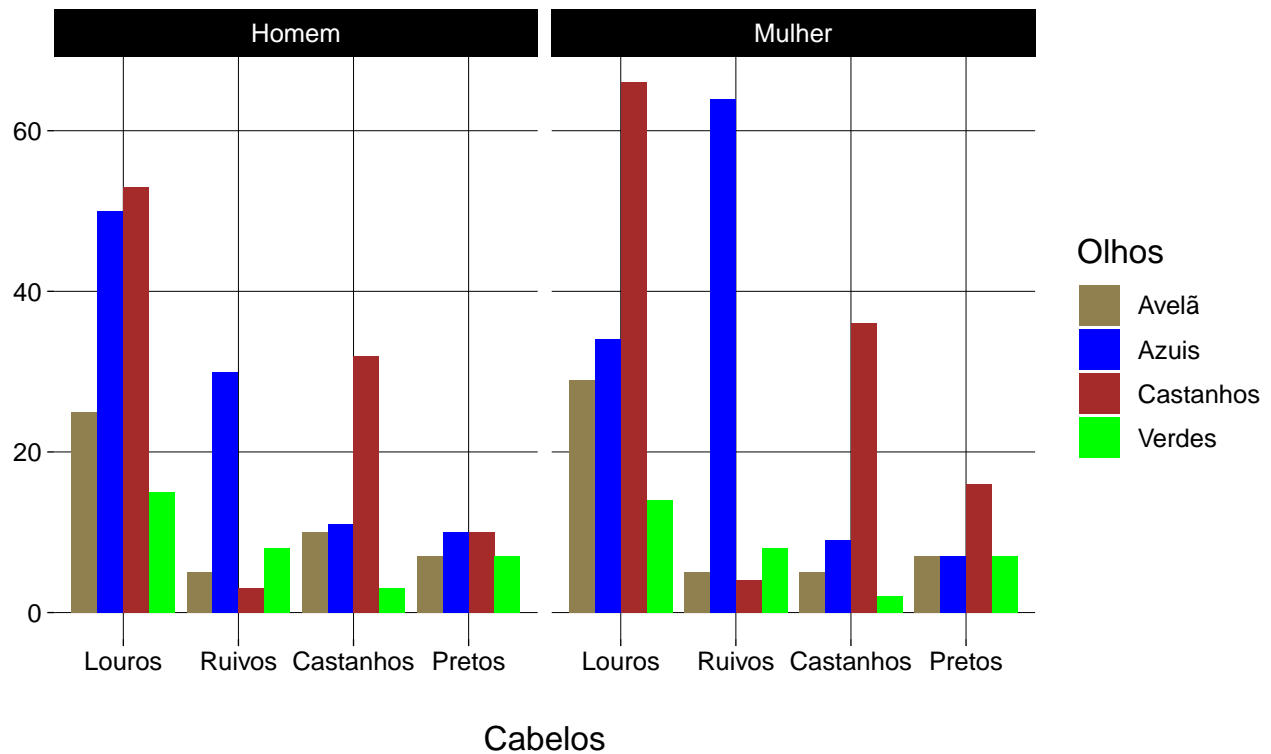


```
df %>%  
  ggplot(aes(x = fct_infreq(Cabelos), fill = Olhos)) +  
  geom_bar(position = 'dodge') +  
  scale_fill_discrete(type = c('#908050', 'blue', 'brown', 'green'))
```



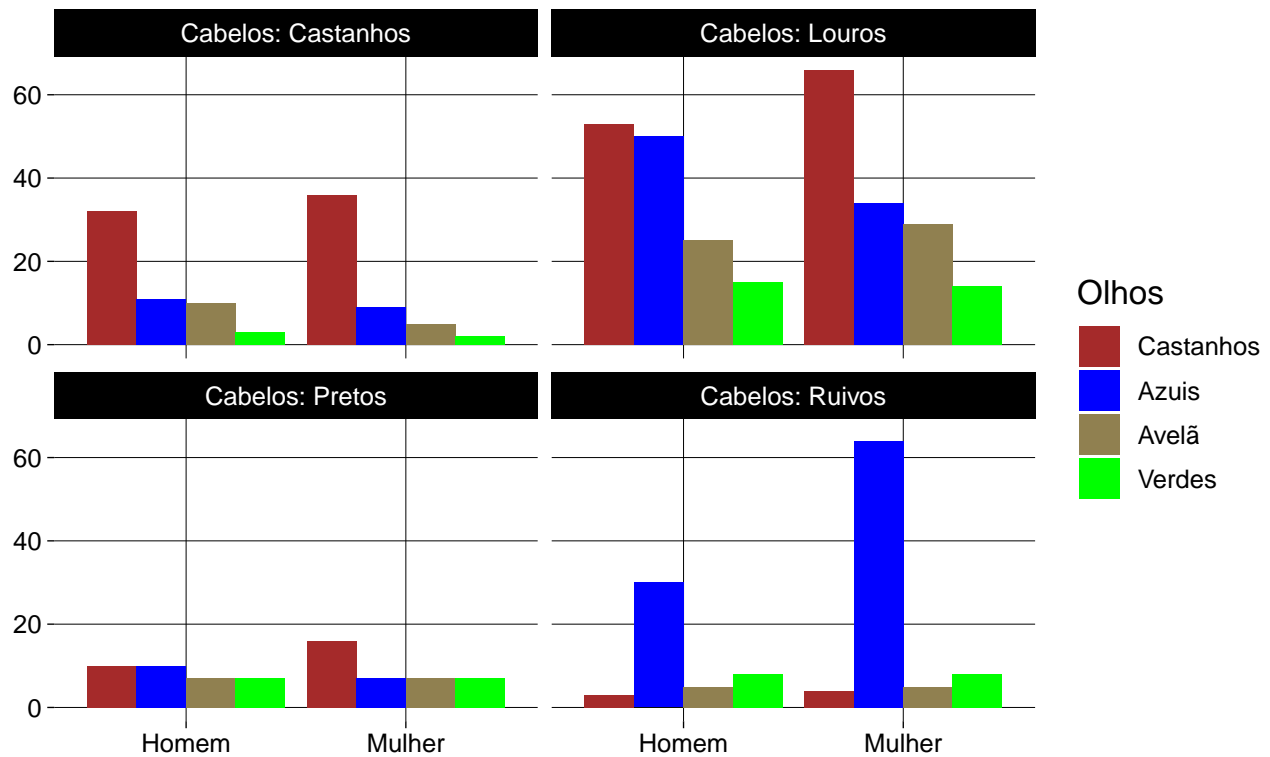
```
df %>%
  ggplot(aes(x = fct_infreq(Cabelos), fill = Olhos)) +
    geom_bar(position = 'dodge') +
    scale_fill_discrete(type = c('#908050', 'blue', 'brown', 'green')) +
    facet_wrap(~Sexo) +
    labs(
      title = 'Cores de cabelos e olhos por sexo',
      y = NULL,
      x = 'Cabelos'
    )
)
```

## Cores de cabelos e olhos por sexo



```
df %>%
  ggplot(aes(x = Sexo, fill = fct_infreq(Olhos))) +
  geom_bar(position = 'dodge') +
  facet_wrap(~Cabelos, labeller = label_both) +
  scale_fill_discrete(type = c('brown', 'blue', '#908050', 'green')) +
  labs(
    x = NULL,
    y = NULL,
    fill = 'Olhos',
    title = 'Cor dos olhos e sexo por cor dos cabelos'
  )
```

## Cor dos olhos e sexo por cor dos cabelos



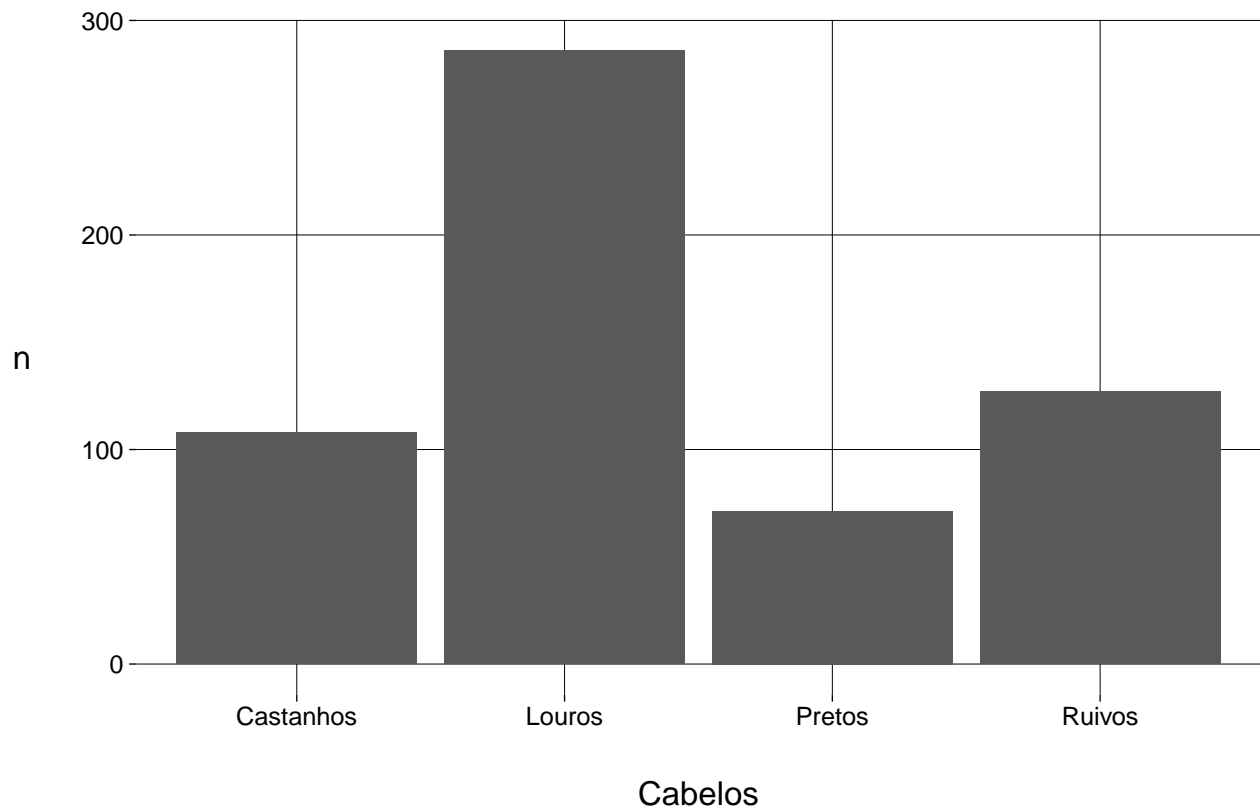
### 4.4.2

**Dataframe já contendo os totais**

```
df_tot <- df %>%
  count(Sexo, Cabelos, Olhos)

df_tot
## # A tibble: 32 x 4
##   Sexo  Cabelos  Olhos      n
##   <chr> <chr>    <chr>  <int>
## 1 Homem Castanhos Avelã    10
## 2 Homem Castanhos Azuis    11
## 3 Homem Castanhos Castanhos  32
## 4 Homem Castanhos Verdes     3
## 5 Homem Louros    Avelã    25
## 6 Homem Louros    Azuis    50
## # ... with 26 more rows

df_tot %>%
  ggplot(aes(x = Cabelos, y = n)) +
  geom_col()
```



## 4.5

### Gráficos de linha e séries temporais

#### 4.5.1

##### Dataset

WorldPhones

```
##      N.Amer Europe Asia S.Amer Oceania Africa Mid.Amer
## 1951 45939 21574 2876  1815  1646    89    555
## 1956 60423 29990 4708  2568  2366  1411    733
## 1957 64721 32510 5230  2695  2526  1546    773
## 1958 68484 35218 6662  2845  2691  1663    836
## 1959 71799 37598 6856  3000  2868  1769    911
## 1960 76036 40341 8220  3145  3054  1905   1008
## 1961 79831 43173 9053  3338  3224  2005   1076
```

```
fones <- WorldPhones %>%
  as_tibble(rownames = 'Ano') %>%
  mutate(Ano = as.numeric(Ano))
```

fones

```
## # A tibble: 7 x 8
```

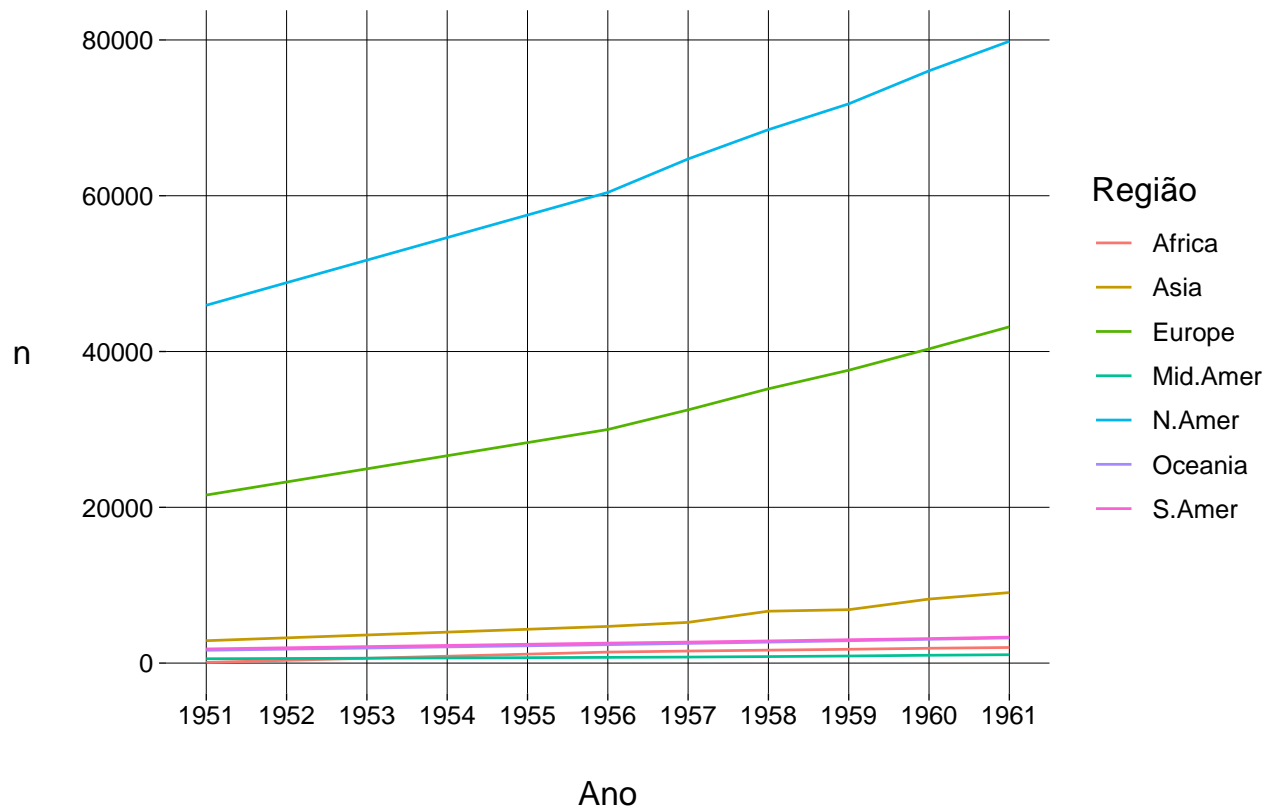


```
##      Ano N.Amer Europe  Asia S.Amer Oceania Africa Mid.Amer
##      <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
## 1  1951  45939  21574  2876   1815   1646    89   555
## 2  1956  60423  29990  4708   2568   2366   1411   733
## 3  1957  64721  32510  5230   2695   2526   1546   773
## 4  1958  68484  35218  6662   2845   2691   1663   836
## 5  1959  71799  37598  6856   3000   2868   1769   911
## 6  1960  76036  40341  8220   3145   3054   1905  1008
## # ... with 1 more row
```

```
fones_long <- fones %>%
  pivot_longer(
    cols = -Ano,
    names_to = 'Região',
    values_to = 'n'
  )

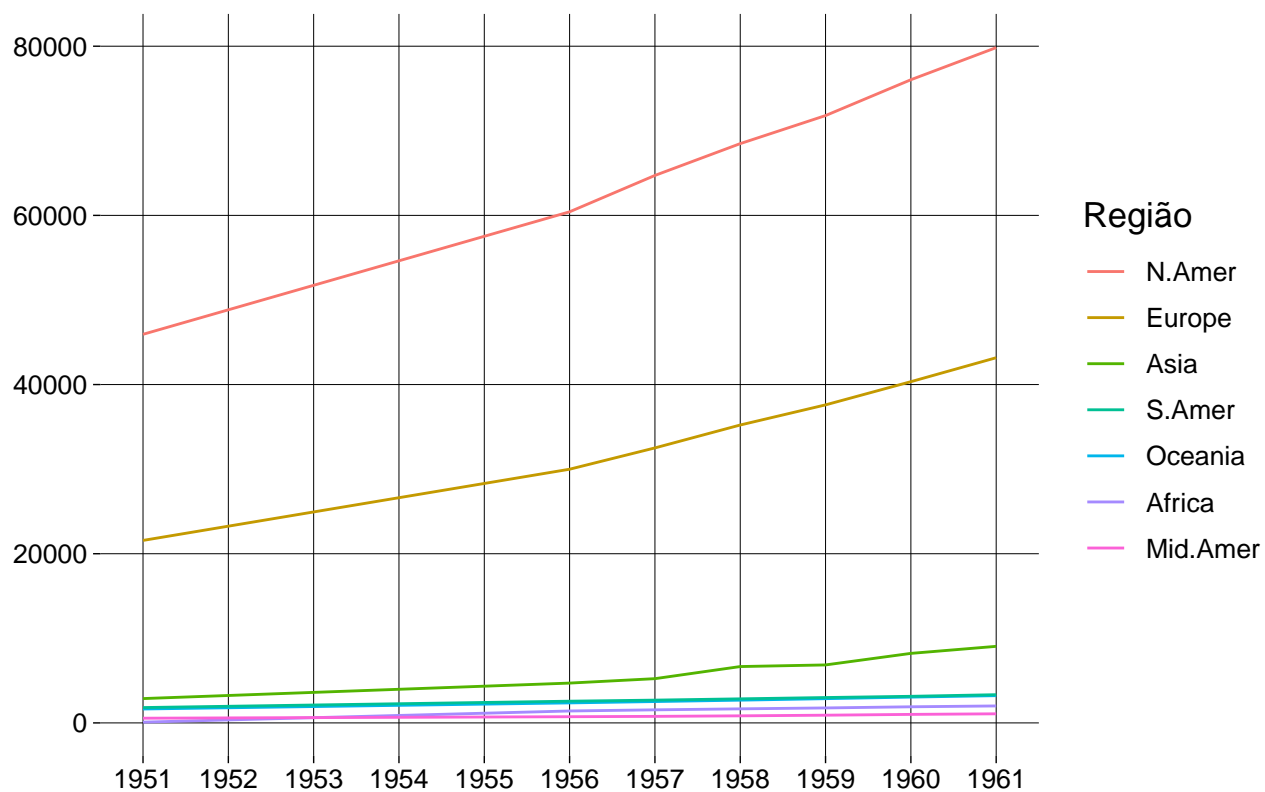
fones_long
## # A tibble: 49 x 3
##      Ano Região      n
##      <dbl> <chr>   <dbl>
## 1  1951 N.Amer  45939
## 2  1951 Europe  21574
## 3  1951 Asia    2876
## 4  1951 S.Amer  1815
## 5  1951 Oceania 1646
## 6  1951 Africa    89
## # ... with 43 more rows
```

```
fones_long %>%
  ggplot(aes(x = Ano, y = n, group = Região, color = Região)) +
  geom_line() +
  scale_x_continuous(breaks = 1951:1961)
```

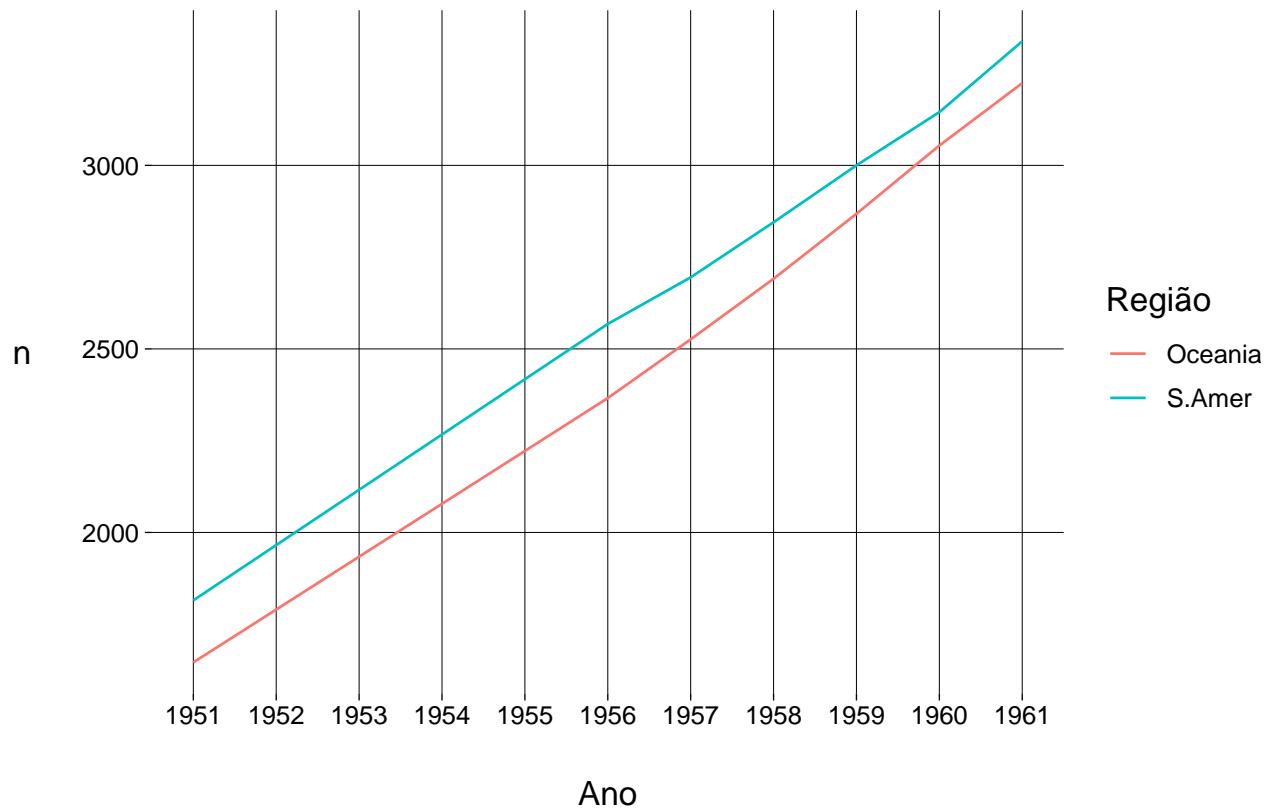


```
fones_long %>%
  ggplot(
    aes(
      x = Ano,
      y = n,
      group = Região,
      color = fct_rev(fct_reorder(Região, n))
    )
  ) +
  geom_line() +
  scale_x_continuous(breaks = 1951:1961) +
  labs(
    color = 'Região',
    y = '',
    x = NULL,
    title = 'Quantidade de aparelhos de telefone por ano, por região'
  )
)
```

## Quantidade de aparelhos de telefone por ano, por região



```
fones_long %>%  
  filter(Região %in% c('S.Amer', 'Oceania')) %>%  
  ggplot(aes(x = Ano, y = n, group = Região, color = Região)) +  
    geom_line() +  
    scale_x_continuous(breaks = 1951:1961)
```



```
library(tsibble)  
?`tsibble-package`
```

## 4.6

### Referências sobre visualização e R



Busque mais informações sobre os pacotes `tidyverse` e `ggplot2` nas referências recomendadas.

## CAPÍTULO 5

---

### Medidas

---