

# Probabilidade e Estatística com R

Fernando Náufel

(versão de 12/04/2022)

---

## 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 . . . . .	7
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 Introdução ao tidyverse</b>	<b>12</b>
3.1 Criando uma <i>tibble</i> . . . . .	12
3.2 Operador de <i>pipe</i> ( <i>%&gt;%</i> ) . . . . .	14
3.3 Formato <i>tidy</i> . . . . .	15
3.4 Manipulando os dados . . . . .	20
3.5 Exercícios . . . . .	34
<b>4 Visualização com ggplot2</b>	<b>37</b>
4.1 Vídeo 1 . . . . .	37
4.2 Componentes de um gráfico ggplot2 . . . . .	37
4.3 Conjunto de dados . . . . .	40
4.4 Gráficos de dispersão ( <i>scatter plots</i> ) . . . . .	46
4.5 Vídeo 2 . . . . .	63
4.6 Histogramas e cia. . . . .	63
4.7 Ogiva . . . . .	68
4.8 Ramos e folhas . . . . .	69

4.9	Personalização do tema	70
4.10	Exercícios	73
<b>5</b>	<b>Visualização com ggplot2 (continuação)</b>	<b>77</b>
5.1	Vídeo 1	77
5.2	<i>Boxplots</i>	77
5.3	Vídeo 2	87
5.4	Gráficos de barras e de colunas	87
5.5	Gráficos de linha e séries temporais	101
5.6	Exercícios	106
5.7	Referências sobre visualização e R	106
<b>6</b>	<b>Medidas</b>	<b>107</b>
6.1	Vídeo	107
6.2	Medidas de centralidade	107
6.3	Formas de uma distribuição	115
6.4	Re-expressão	121
6.5	Medidas de posição	122
6.6	Medidas de dispersão	123
6.7	Coeficiente de variação	130
6.8	Escores-padrão	132
6.9	Teorema de Tchebychev	134

---

## Apresentação

---

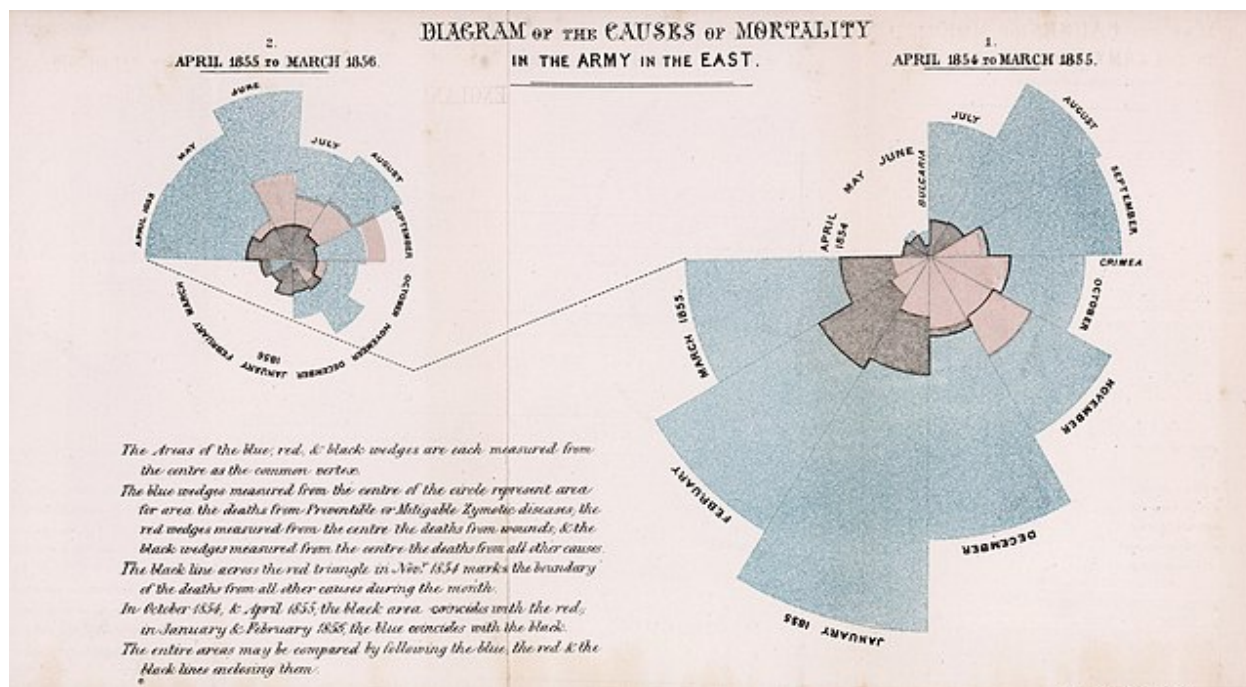
### Atenção



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](#), 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.
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](#).

---

## 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 Grolemund, 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](#)?



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?

## 1.3

### Vídeo 2

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





LIMITAÇÕES DE ESTUDOS COM CEGAMENTO

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

## 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), 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 (`...`), tecla *Enter*.
13. Se nada funcionar, tecla *ESC*.
14. Para sair do `swirl()`, digite `bye()` no *prompt* do R.
15. Para voltar para os exercícios, digite

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

---

### Introdução ao tidyverse

---



Busque mais informações sobre os pacotes que compõem o tidyverse [nas referências recomendadas](#).

### 3.1

---

#### Criando uma *tibble*

- Uma *tibble* é uma tabela retangular.
- Cada coluna é um vetor:

```
cores <- tibble(  
  pessoa = c('João', 'Maria', 'Pedro', 'Ana'),  
  'cor favorita' = c('azul', 'rosa', 'preto', 'branco')  
)
```

```
cores
```

```
## # A tibble: 4 x 2  
##   pessoa `cor favorita`  
##   <chr>   <chr>  
## 1 João   azul  
## 2 Maria  rosa
```

```
## 3 Pedro  preto
## 4 Ana    branco
```

- Isto é um pouco diferente da maneira como estamos acostumados a ver tabelas (como uma coleção de linhas, em vez de uma coleção de colunas).
- A função `tribble` permite a entrada de forma mais natural, linha a linha. **Lembre-se de usar `~` antes dos nomes das colunas.**

```
cores <- tribble(
  ~pessoa, ~'cor favorita',
  "João",   "azul",
  "Maria",  "rosa",
  "Pedro",  "preto",
  "Ana",    "branco"
)

cores
```

```
## # A tibble: 4 x 2
##   pessoa `cor favorita`
##   <chr>   <chr>
## 1 João   azul
## 2 Maria  rosa
## 3 Pedro  preto
## 4 Ana    branco
```



Mesmo que você crie uma *tibble* linha a linha, o R vai continuar tratando sua *tibble* como uma coleção de colunas.

É importante lembrar disto para entender a forma como R manipula estas tabelas.

- Se uma coluna não puder ser armazenada em um vetor, a coluna será uma lista (com vetores como elementos):

```
cores <- tibble(
  pessoa = c('João', 'Maria', 'Pedro', 'Ana'),
  'cor favorita' = list(
    c('azul', 'roxo'),
    c('rosa', 'magenta'),
    NA,
    'branco'
  )
)

cores
```

```
## # A tibble: 4 x 2
```

```
##  pessoa `cor favorita`
##  <chr>  <list>
## 1 João   <chr [2]>
## 2 Maria  <chr [2]>
## 3 Pedro  <lgl [1]>
## 4 Ana    <chr [1]>
```

- Use `View()` para examinar interativamente o conteúdo de uma coluna-lista:

```
cores %>% View()
```

## 3.2

### Operador de *pipe* (`%>%`)

- O tidyverse inclui o pacote `magrittr`, que contém o operador `%>%`, chamado *pipe*.<sup>1</sup>
- A ideia é facilitar a leitura de **composições de funções**. O código

```
y <- h(g(f(x)))
```

pode ser escrito como

```
y <- x %>% f() %>% g() %>% h()
```

- Esta segunda versão é mais fiel à ordem em que as operações acontecem.
- Na verdade, R tem um operador de **atribuição para a direita**, mas poucas pessoas recomendam usá-lo:

```
x %>% f() %>% g() %>% h() -> y
```

- Se `f`, `g` e `h` forem funções de um argumento só, os parênteses podem ser omitidos:

```
y <- x %>% f %>% g %>% h
```

- Se a função `f` tiver outros argumentos, escreva-os normalmente na chamada a `f`:

```
y <- x %>% mean(na.rm = TRUE)
```

- O *pipe* `EXP %>% f(...)` sempre insere o resultado da expressão `EXP` como o **primeiro argumento da função `f`**.
- Se você precisar que o resultado da expressão `EXP` seja inserido em outra posição na lista de argumentos de `f`, use um ponto “.” para isso:

<sup>1</sup>Por que o nome do pacote e o nome do operador formam um trocadilho?

```
x %>% consultar(df, .)
```

equivale a

```
consultar(df, x)
```

### 3.3

## Formato *tidy*

- Nossa última versão da *tibble* `cores` é um pouco mais complexa do que deveria ser:

```
cores
```

```
## # A tibble: 4 x 2
##   pessoa `cor favorita`
##   <chr>   <list>
## 1 João   <chr [2]>
## 2 Maria <chr [2]>
## 3 Pedro <lgl [1]>
## 4 Ana   <chr [1]>
```

- O formato *tidy* exige que
  1. Cada linha da *tibble* corresponda a uma observação sobre um indivíduo,
  2. Cada coluna corresponda a uma variável observada, e
  3. Cada célula contenha um valor da variável.
- Na *tibble* `cores`, a primeira e a segunda exigências são satisfeitas, mas a terceira não, pois algumas células contêm valores múltiplos.
- A *tibble* não está no formato *tidy*.
- Podemos “extrair” estes vetores “aninhados” usando o comando `unnest`, do pacote `tidyr`:

```
cores <- cores %>%
  unnest(`cor favorita`)

cores
```

```
## # A tibble: 6 x 2
##   pessoa `cor favorita`
##   <chr>   <chr>
## 1 João   azul
## 2 João   roxo
## 3 Maria  rosa
## 4 Maria  magenta
```



```
## 5 Pedro <NA>
## 6 Ana branco
```

- A maioria das funções do tidyverse exige que as *tibbles* estejam neste formato *tidy*.
- Um exemplo mais complexo é o *dataset* billboard, com as seguintes colunas (para cada música que estava no *top 100* da Billboard no ano de 2000):
  - Nome do artista ou banda;
  - Nome da música;
  - Data em que a música entrou no *top 100* da Billboard;
  - Para cada uma das 76 semanas seguintes, a posição da música no *top 100*.

```
billboard %>% glimpse()
```

```
## Rows: 317
## Columns: 79
## $ artist      <chr> "2 Pac", "2Ge+her", "3 Doors Down", "3 Doors Dow~
## $ track       <chr> "Baby Don't Cry (Keep...", "The Hardest Part Of ~
## $ date.entered <date> 2000-02-26, 2000-09-02, 2000-04-08, 2000-10-21,~
## $ wk1         <dbl> 87, 91, 81, 76, 57, 51, 97, 84, 59, 76, 84, 57, ~
## $ wk2         <dbl> 82, 87, 70, 76, 34, 39, 97, 62, 53, 76, 84, 47, ~
## $ wk3         <dbl> 72, 92, 68, 72, 25, 34, 96, 51, 38, 74, 75, 45, ~
## $ wk4         <dbl> 77, NA, 67, 69, 17, 26, 95, 41, 28, 69, 73, 29, ~
## $ wk5         <dbl> 87, NA, 66, 67, 17, 26, 100, 38, 21, 68, 73, 23,~
## $ wk6         <dbl> 94, NA, 57, 65, 31, 19, NA, 35, 18, 67, 69, 18, ~
## $ wk7         <dbl> 99, NA, 54, 55, 36, 2, NA, 35, 16, 61, 68, 11, 2~
## $ wk8         <dbl> NA, NA, 53, 59, 49, 2, NA, 38, 14, 58, 65, 9, 17~
## $ wk9         <dbl> NA, NA, 51, 62, 53, 3, NA, 38, 12, 57, 73, 9, 17~
## $ wk10        <dbl> NA, NA, 51, 61, 57, 6, NA, 36, 10, 59, 83, 11, 1~
## $ wk11        <dbl> NA, NA, 51, 61, 64, 7, NA, 37, 9, 66, 92, 1, 17,~
## $ wk12        <dbl> NA, NA, 51, 59, 70, 22, NA, 37, 8, 68, NA, 1, 3,~
## $ wk13        <dbl> NA, NA, 47, 61, 75, 29, NA, 38, 6, 61, NA, 1, 3,~
## $ wk14        <dbl> NA, NA, 44, 66, 76, 36, NA, 49, 1, 67, NA, 1, 7,~
## $ wk15        <dbl> NA, NA, 38, 72, 78, 47, NA, 61, 2, 59, NA, 4, 10~
## $ wk16        <dbl> NA, NA, 28, 76, 85, 67, NA, 63, 2, 63, NA, 8, 17~
## $ wk17        <dbl> NA, NA, 22, 75, 92, 66, NA, 62, 2, 67, NA, 12, 2~
## $ wk18        <dbl> NA, NA, 18, 67, 96, 84, NA, 67, 2, 71, NA, 22, 2~
## $ wk19        <dbl> NA, NA, 18, 73, NA, 93, NA, 83, 3, 79, NA, 23, 2~
## $ wk20        <dbl> NA, NA, 14, 70, NA, 94, NA, 86, 4, 89, NA, 43, 4~
## $ wk21        <dbl> NA, NA, 12, NA, NA, NA, NA, NA, 5, NA, NA, 44, 4~
## $ wk22        <dbl> NA, NA, 7, NA, NA, NA, NA, NA, 5, NA, NA, NA, 50~
## $ wk23        <dbl> NA, NA, 6, NA, NA, NA, NA, NA, 6, NA, NA, NA, NA~
## $ wk24        <dbl> NA, NA, 6, NA, NA, NA, NA, NA, 9, NA, NA, NA, NA~
## $ wk25        <dbl> NA, NA, 6, NA, NA, NA, NA, NA, 13, NA, NA, NA, N~
## $ wk26        <dbl> NA, NA, 5, NA, NA, NA, NA, NA, 14, NA, NA, NA, N~
## $ wk27        <dbl> NA, NA, 5, NA, NA, NA, NA, NA, 16, NA, NA, NA, N~
```

[illegible]

```
## $ wk76          <dbl> NA, NA, NA, NA, NA, NA, NA, NA, NA, NA, NA, NA, ~
```

- Vamos renomear as colunas:

```
bb <- billboard %>%  
  rename(  
    artista = artist,  
    musica = track,  
    entrou = date.entered  
  )
```

```
bb %>% head()
```

```
## # A tibble: 6 x 79  
##   artista  musica entrou      wk1    wk2    wk3    wk4    wk5    wk6    wk7  
##   <chr>    <chr> <date>    <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>  
## 1 2 Pac    Baby ~ 2000-02-26    87    82    72    77    87    94    99  
## 2 2Ge+her  The H~ 2000-09-02    91    87    92    NA    NA    NA    NA  
## 3 3 Doors ~ Krypt~ 2000-04-08    81    70    68    67    66    57    54  
## 4 3 Doors ~ Loser 2000-10-21    76    76    72    69    67    65    55  
## 5 504 Boyz Wobbl~ 2000-04-15    57    34    25    17    17    31    36  
## 6 98^0     Give ~ 2000-08-19    51    39    34    26    26    19    2  
## # ... with 69 more variables: wk8 <dbl>, wk9 <dbl>, wk10 <dbl>,  
## #   wk11 <dbl>, wk12 <dbl>, wk13 <dbl>, wk14 <dbl>, wk15 <dbl>,  
## #   wk16 <dbl>, wk17 <dbl>, wk18 <dbl>, wk19 <dbl>, wk20 <dbl>,  
## #   wk21 <dbl>, wk22 <dbl>, wk23 <dbl>, wk24 <dbl>, wk25 <dbl>,  
## #   wk26 <dbl>, wk27 <dbl>, wk28 <dbl>, wk29 <dbl>, wk30 <dbl>,  
## #   wk31 <dbl>, wk32 <dbl>, wk33 <dbl>, wk34 <dbl>, wk35 <dbl>,  
## #   wk36 <dbl>, wk37 <dbl>, wk38 <dbl>, wk39 <dbl>, wk40 <dbl>, ...
```

- O que é uma observação neste conjunto de dados?  
A posição, em uma semana, de uma música que esteve no *top 100* da *Billboard* durante o ano 2000.
- Quais são as variáveis que qualificam cada observação?
  - O artista,
  - O título da música,
  - A posição da música no *top 100* da *Billboard* em cada uma das 76 semanas depois que ela entrou na lista.
- Este último item é complexo, e o criador da *tibble* decidiu criar uma coluna por semana.
- Uma decisão ruim, pois existe informação embutida nos nomes das colunas. A coluna `wk68` corresponde à posição da música na semana 68 após ela entrar na lista, mas o número da semana só aparece no nome da coluna!
- Isto nunca deve acontecer. A informação deve sempre estar nas células.

- Vamos simplificar as coisas criando duas colunas:
  - semana, com o número da semana; perceba que esta informação vem dos nomes das colunas,
  - pos, com a posição da música naquela semana; esta informação vem das células.
- A *tibble*, que antes era larga, **vai ser mais estreita e mais longa.**
- A função `pivot_longer`, do pacote `tidyr`, vai fazer o trabalho — inclusive extraíndo os números das semanas dos nomes das colunas:

```
bb_tidy <- bb %>%
  pivot_longer(
    wk1:wk76,
    names_to = 'semana',
    names_prefix = 'wk',
    names_transform = list(
      semana = as.integer
    ),
    values_to = 'pos'
  )

bb_tidy
```

```
## # A tibble: 24.092 x 5
##   artista musica      entrou  semana  pos
##   <chr>   <chr>      <date>    <int> <dbl>
## 1 2 Pac    Baby Don't Cry (Keep... 2000-02-26      1      87
## 2 2 Pac    Baby Don't Cry (Keep... 2000-02-26      2      82
## 3 2 Pac    Baby Don't Cry (Keep... 2000-02-26      3      72
## 4 2 Pac    Baby Don't Cry (Keep... 2000-02-26      4      77
## 5 2 Pac    Baby Don't Cry (Keep... 2000-02-26      5      87
## 6 2 Pac    Baby Don't Cry (Keep... 2000-02-26      6      94
## # ... with 24.086 more rows
```

- O R só mostra, por *default*, as 1000 primeiras linhas de uma *tibble*.
- Na verdade, o número de linhas da tabela original era

```
bb %>% nrow()
```

```
## [1] 317
```

- O número de linhas, depois de `pivot_longer`, ficou:

```
bb_tidy %>% nrow()
```

```
## [1] 24092
```

- Existem linhas onde `pos` tem o valor `NA`. São resultado da organização original dos dados, onde o `NA` indicava que a música não estava no *top 100* naquela semana.
- No novo formato, a ausência da linha com aquele número de semana já basta para indicar isto. Então, vamos eliminar as linhas onde `pos` é `NA`.
- A função `filter` **mantém** as linhas que **satisfazem** a condição dada; por isso, a condição é “`pos` não é `NA`”:

```
bb_tidy <- bb_tidy %>%
  filter(!is.na(pos))
```

```
bb_tidy
```

```
## # A tibble: 5.307 x 5
##   artista musica      entrou   semana   pos
##   <chr>    <chr>      <date>   <int> <dbl>
## 1 2 Pac    Baby Don't Cry (Keep... 2000-02-26     1     87
## 2 2 Pac    Baby Don't Cry (Keep... 2000-02-26     2     82
## 3 2 Pac    Baby Don't Cry (Keep... 2000-02-26     3     72
## 4 2 Pac    Baby Don't Cry (Keep... 2000-02-26     4     77
## 5 2 Pac    Baby Don't Cry (Keep... 2000-02-26     5     87
## 6 2 Pac    Baby Don't Cry (Keep... 2000-02-26     6     94
## # ... with 5.301 more rows
```

- O número de linhas ficou

```
bb_tidy %>% nrow()
```

```
## [1] 5307
```

### 3.3.1

#### Exercícios

- Todas as semanas deste conjunto de dados são do ano 2000?
- Qual é o tipo do **primeiro** argumento da função `filter()`?

## 3.4

### Manipulando os dados

#### 3.4.1

## Criando novas colunas: mutate, transmute

- O *data frame*<sup>2</sup> `cars` tem dados (de 1920!) sobre as distâncias de frenagem (em pés) de um carro viajando a diversas velocidades (em milhas por hora):

```
cars
```

```
## # A tibble: 50 x 2
##   speed  dist
##   <dbl> <dbl>
## 1     4     2
## 2     4    10
## 3     7     4
## 4     7    22
## 5     8    16
## 6     9    10
## # ... with 44 more rows
```

- Vamos criar colunas novas com os valores convertidos para km/h e metros; além disso, uma coluna com a taxa de frenagem:

```
cars %>%
  mutate(
    velocidade = speed * 1.6,
    distancia = dist * .33,
    taxa = velocidade / distancia
  )
```

```
## # A tibble: 50 x 5
##   speed  dist velocidade distancia  taxa
##   <dbl> <dbl>      <dbl>      <dbl> <dbl>
## 1     4     2        6.4        0.66  9.70
## 2     4    10        6.4        3.3   1.94
## 3     7     4       11.2        1.32  8.48
## 4     7    22       11.2        7.26  1.54
## 5     8    16       12.8        5.28  2.42
## 6     9    10       14.4        3.3   4.36
## # ... with 44 more rows
```

- Perceba que as colunas antigas continuam lá. Se quiser manter apenas as colunas novas, use `transmute`:

```
cars %>%
  transmute(
    velocidade = speed * 1.6,
```

<sup>2</sup>Considere *data frame* como sinônimo de *tibble*. Na verdade, *tibbles* formam um superconjunto de *data frames*: todo *data frame* é uma *tibble*, mas nem toda *tibble* é um *data frame*.

```

    distancia = dist * .33,
    taxa = velocidade / distancia
  )

```

```

## # A tibble: 50 x 3
##   velocidade distancia taxa
##   <dbl>      <dbl> <dbl>
## 1      6.4      0.66  9.70
## 2      6.4      3.3   1.94
## 3     11.2      1.32  8.48
## 4     11.2      7.26  1.54
## 5     12.8      5.28  2.42
## 6     14.4      3.3   4.36
## # ... with 44 more rows

```

- Ou use o argumento `.keep` de `mutate` para escolher com mais precisão. Veja a ajuda de `mutate`.

### 3.4.2

**Selecionando colunas:** `select`, `distinct`, `pull`

- Vamos voltar à nossa *tibble* dos *top 100* da *Billboard*.
- Para ver só a coluna de artistas:

```

bb_tidy %>%
  select(artista)

```

```

## # A tibble: 5.307 x 1
##   artista
##   <chr>
## 1 2 Pac
## 2 2 Pac
## 3 2 Pac
## 4 2 Pac
## 5 2 Pac
## 6 2 Pac
## # ... with 5.301 more rows

```

- Para eliminar as repetições:

```

bb_tidy %>%
  select(artista) %>%
  distinct()

```

```

## # A tibble: 228 x 1
##   artista
##   <chr>
## 1 2 Pac

```

```
## 2 2Ge+her
## 3 3 Doors Down
## 4 504 Boyz
## 5 98^0
## 6 A*Teens
## # ... with 222 more rows
```

- Para ver artistas e músicas:

```
bb_tidy %>%
  select(artista, musica) %>%
  distinct()
```

```
## # A tibble: 317 x 2
##   artista      musica
##   <chr>      <chr>
## 1 2 Pac      Baby Don't Cry (Keep...
## 2 2Ge+her    The Hardest Part Of ...
## 3 3 Doors Down Kryptonite
## 4 3 Doors Down Loser
## 5 504 Boyz    Wobble Wobble
## 6 98^0       Give Me Just One Nig...
## # ... with 311 more rows
```

- Para especificar colunas **a não mostrar**, use o sinal de menos “-”:

```
bb_tidy %>%
  select(-c(entrou, semana, pos))
```

```
## # A tibble: 5.307 x 2
##   artista musica
##   <chr>    <chr>
## 1 2 Pac    Baby Don't Cry (Keep...
## 2 2 Pac    Baby Don't Cry (Keep...
## 3 2 Pac    Baby Don't Cry (Keep...
## 4 2 Pac    Baby Don't Cry (Keep...
## 5 2 Pac    Baby Don't Cry (Keep...
## 6 2 Pac    Baby Don't Cry (Keep...
## # ... with 5.301 more rows
```

- Para **extrair uma coluna na forma de vetor** (unique é uma função do R base, aplicável a vetores):

```
bb_tidy %>%
  pull(artista) %>%
  unique()
```

```
## [1] "2 Pac" "2Ge+her"
## [3] "3 Doors Down" "504 Boyz"
```



##	[5]	"98~0"	"A*Teens"
##	[7]	"Aaliyah"	"Adams, Yolanda"
##	[9]	"Adkins, Trace"	"Aguilera, Christina"
##	[11]	"Alice DeeJay"	"Allan, Gary"
##	[13]	"Amber"	"Anastacia"
##	[15]	"Anthony, Marc"	"Avant"
##	[17]	"BBMak"	"Backstreet Boys, The"
##	[19]	"Badu, Erkyah"	"Baha Men"
##	[21]	"Barenaked Ladies"	"Beenie Man"
##	[23]	"Before Dark"	"Bega, Lou"
##	[25]	"Big Punisher"	"Black Rob"
##	[27]	"Black, Clint"	"Blaque"
##	[29]	"Blige, Mary J."	"Blink-182"
##	[31]	"Bloodhound Gang"	"Bon Jovi"
##	[33]	"Braxton, Toni"	"Brock, Chad"
##	[35]	"Brooks & Dunn"	"Brooks, Garth"
##	[37]	"Byrd, Tracy"	"Cagle, Chris"
##	[39]	"Cam'ron"	"Carey, Mariah"
##	[41]	"Carter, Aaron"	"Carter, Torrey"
##	[43]	"Changing Faces"	"Chesney, Kenny"
##	[45]	"Clark Family Experience"	"Clark, Terri"
##	[47]	"Common"	"Counting Crows"
##	[49]	"Creed"	"Cyrus, Billy Ray"
##	[51]	"D'Angelo"	"DMX"
##	[53]	"Da Brat"	"Davidson, Clay"
##	[55]	"De La Soul"	"Destiny's Child"
##	[57]	"Diffie, Joe"	"Dion, Celine"
##	[59]	"Dixie Chicks, The"	"Dr. Dre"
##	[61]	"Drama"	"Dream"
##	[63]	"Eastsidaz, The"	"Eiffel 65"
##	[65]	"Elliott, Missy \"Misdemeanor\""	"Eminem"
##	[67]	"En Vogue"	"Estefan, Gloria"
##	[69]	"Evans, Sara"	"Eve"
##	[71]	"Everclear"	"Fabian, Lara"
##	[73]	"Fatboy Slim"	"Filter"
##	[75]	"Foo Fighters"	"Fragma"
##	[77]	"Funkmaster Flex"	"Ghostface Killah"
##	[79]	"Gill, Vince"	"Gilman, Billy"
##	[81]	"Ginuwine"	"Goo Goo Dolls"
##	[83]	"Gray, Macy"	"Griggs, Andy"
##	[85]	"Guy"	"Hanson"
##	[87]	"Hart, Beth"	"Heatherly, Eric"
##	[89]	"Henley, Don"	"Herndon, Ty"
##	[91]	"Hill, Faith"	"Hoku"
##	[93]	"Hollister, Dave"	"Hot Boys"
##	[95]	"Houston, Whitney"	"IMx"
##	[97]	"Ice Cube"	"Ideal"
##	[99]	"Iglesias, Enrique"	"J-Shin"

## [101] "Ja Rule"	"Jackson, Alan"
## [103] "Jagged Edge"	"Janet"
## [105] "Jay-Z"	"Jean, Wyclef"
## [107] "Joe"	"John, Elton"
## [109] "Jones, Donell"	"Jordan, Montell"
## [111] "Juvenile"	"Kandi"
## [113] "Keith, Toby"	"Kelis"
## [115] "Kenny G"	"Kid Rock"
## [117] "Kravitz, Lenny"	"Kumbia Kings"
## [119] "LFO"	"LL Cool J"
## [121] "Larrieux, Amel"	"Lawrence, Tracy"
## [123] "Levert, Gerald"	"Lil Bow Wow"
## [125] "Lil Wayne"	"Lil' Kim"
## [127] "Lil' Mo"	"Lil' Zane"
## [129] "Limp Bizkit"	"Lonestar"
## [131] "Lopez, Jennifer"	"Loveless, Patty"
## [133] "Lox"	"Lucy Pearl"
## [135] "Ludacris"	"M2M"
## [137] "Madison Avenue"	"Madonna"
## [139] "Martin, Ricky"	"Mary Mary"
## [141] "Master P"	"McBride, Martina"
## [143] "McEntire, Reba"	"McGraw, Tim"
## [145] "McKnight, Brian"	"Messina, Jo Dee"
## [147] "Metallica"	"Montgomery Gentry"
## [149] "Montgomery, John Michael"	"Moore, Chante"
## [151] "Moore, Mandy"	"Mumba, Samantha"
## [153] "MusiQ"	"Mya"
## [155] "Mystikal"	"N'Sync"
## [157] "Nas"	"Nelly"
## [159] "Next"	"Nine Days"
## [161] "No Doubt"	"Nu Flavor"
## [163] "Offspring, The"	"Paisley, Brad"
## [165] "Papa Roach"	"Pearl Jam"
## [167] "Pink"	"Price, Kelly"
## [169] "Profyle"	"Puff Daddy"
## [171] "Q-Tip"	"R.E.M."
## [173] "Rascal Flatts"	"Raye, Collin"
## [175] "Red Hot Chili Peppers"	"Rimes, LeAnn"
## [177] "Rogers, Kenny"	"Ruff Endz"
## [179] "Sammie"	"Santana"
## [181] "Savage Garden"	"SheDaisy"
## [183] "Sheist, Shade"	"Shyne"
## [185] "Simpson, Jessica"	"Sisqo"
## [187] "Sister Hazel"	"Smash Mouth"
## [189] "Smith, Will"	"Son By Four"
## [191] "Sonique"	"SoulDecision"
## [193] "Spears, Britney"	"Spencer, Tracie"
## [195] "Splender"	"Sting"

## [197]	"Stone Temple Pilots"	"Stone, Angie"
## [199]	"Strait, George"	"Sugar Ray"
## [201]	"TLC"	"Tamar"
## [203]	"Tamia"	"Third Eye Blind"
## [205]	"Thomas, Carl"	"Tippin, Aaron"
## [207]	"Train"	"Trick Daddy"
## [209]	"Trina"	"Tritt, Travis"
## [211]	"Tuesday"	"Urban, Keith"
## [213]	"Usher"	"Vassar, Phil"
## [215]	"Vertical Horizon"	"Vitamin C"
## [217]	"Walker, Clay"	"Wallflowers, The"
## [219]	"Westlife"	"Williams, Robbie"
## [221]	"Wills, Mark"	"Worley, Darryl"
## [223]	"Wright, Chely"	"Yankee Grey"
## [225]	"Yearwood, Trisha"	"Ying Yang Twins"
## [227]	"Zombie Nation"	"matchbox twenty"

### 3.4.3

Filtrando linhas: `filter`, `slice`

- Apenas as músicas da Britney Spears:

```
bb_tidy %>%
  filter(artista == 'Spears, Britney')
```

```
## # A tibble: 51 x 5
##   artista      musica      entrou      semana  pos
##   <chr>        <chr>        <date>      <int> <dbl>
## 1 Spears, Britney From The Bottom Of M... 2000-01-29      1     76
## 2 Spears, Britney From The Bottom Of M... 2000-01-29      2     59
## 3 Spears, Britney From The Bottom Of M... 2000-01-29      3     52
## 4 Spears, Britney From The Bottom Of M... 2000-01-29      4     52
## 5 Spears, Britney From The Bottom Of M... 2000-01-29      5     14
## 6 Spears, Britney From The Bottom Of M... 2000-01-29      6     14
## # ... with 45 more rows
```

- Apenas músicas que chegaram à posição 1, sem mostrar a coluna pos:

```
bb_tidy %>%
  filter(pos == 1) %>%
  select(-pos)
```

```
## # A tibble: 55 x 4
##   artista      musica      entrou      semana
##   <chr>        <chr>        <date>      <int>
## 1 Aaliyah      Try Again    2000-03-18      14
## 2 Aguilera, Christina Come On Over Baby (A... 2000-08-05      11
```

```
## 3 Aguilera, Christina Come On Over Baby (A... 2000-08-05      12
## 4 Aguilera, Christina Come On Over Baby (A... 2000-08-05      13
## 5 Aguilera, Christina Come On Over Baby (A... 2000-08-05      14
## 6 Aguilera, Christina What A Girl Wants      1999-11-27      8
## # ... with 49 more rows
```

- Apenas músicas que chegaram à posição 1 em menos de 10 semanas, mostrando apenas artista e música:

```
bb_tidy %>%
  filter(pos == 1, semana < 10) %>%
  distinct(artista, musica)
```

```
## # A tibble: 5 x 2
##   artista      musica
##   <chr>        <chr>
## 1 Aguilera, Christina What A Girl Wants
## 2 Destiny's Child Independent Women Pa...
## 3 Madonna      Music
## 4 Santana      Maria, Maria
## 5 Sisco        Incomplete
```

- As funções da família `slice` filtram linhas de diversas maneiras.
- De acordo com seus índices (números de linha):

```
bb_tidy %>%
  slice(c(1, 1000, 5000))
```

```
## # A tibble: 3 x 5
##   artista      musica      entrou      semana  pos
##   <chr>        <chr>      <date>      <int> <dbl>
## 1 2 Pac        Baby Don't Cry (Keep~ 2000-02-26      1    87
## 2 Clark Family Experience Meanwhile Back At Th~ 2000-11-18      3    81
## 3 Vassar, Phil Carlene      2000-03-04      3    64
```

```
bb_tidy %>%
  slice_head(n = 4)
```

```
## # A tibble: 4 x 5
##   artista musica      entrou      semana  pos
##   <chr>   <chr>      <date>      <int> <dbl>
## 1 2 Pac   Baby Don't Cry (Keep... 2000-02-26      1    87
## 2 2 Pac   Baby Don't Cry (Keep... 2000-02-26      2    82
## 3 2 Pac   Baby Don't Cry (Keep... 2000-02-26      3    72
## 4 2 Pac   Baby Don't Cry (Keep... 2000-02-26      4    77
```

```
bb_tidy %>%
  slice_tail(n = 4)
```

```
## # A tibble: 4 x 5
##   artista      musica entrou   semana   pos
##   <chr>        <chr>   <date>    <int> <dbl>
## 1 matchbox twenty Bent    2000-04-29    36    37
## 2 matchbox twenty Bent    2000-04-29    37    38
## 3 matchbox twenty Bent    2000-04-29    38    38
## 4 matchbox twenty Bent    2000-04-29    39    48
```

- De acordo com a ordenação de uma coluna ou de uma função das colunas:

```
bb_tidy %>%
  slice_min(pos)
```

```
## # A tibble: 55 x 5
##   artista      musica      entrou   semana   pos
##   <chr>        <chr>    <date>    <int> <dbl>
## 1 Aaliyah      Try Again 2000-03-18    14    1
## 2 Aguilera, Christina Come On Over Baby (A... 2000-08-05    11    1
## 3 Aguilera, Christina Come On Over Baby (A... 2000-08-05    12    1
## 4 Aguilera, Christina Come On Over Baby (A... 2000-08-05    13    1
## 5 Aguilera, Christina Come On Over Baby (A... 2000-08-05    14    1
## 6 Aguilera, Christina What A Girl Wants      1999-11-27     8    1
## # ... with 49 more rows
```

```
bb_tidy %>%
  slice_max(semana)
```

```
## # A tibble: 1 x 5
##   artista musica      entrou   semana   pos
##   <chr>   <chr>   <date>    <int> <dbl>
## 1 Creed  Higher 1999-09-11    65    49
```

- Aleatoriamente, criando uma amostra:

```
bb_tidy %>%
  slice_sample(n = 5)
```

```
## # A tibble: 5 x 5
##   artista      musica      entrou   semana   pos
##   <chr>        <chr>    <date>    <int> <dbl>
## 1 Third Eye Blind Never Let You Go 2000-01-22     6    22
## 2 Houston, Whitney My Love Is Your Love 1999-09-04     2    68
## 3 Goo Goo Dolls Broadway          2000-04-22     7    29
## 4 Nelly          (Hot S**t) Country G... 2000-04-29     8    24
## 5 Beenie Man     Girls Dem Sugar 2000-10-21     2    72
```

- Veja a ajuda de `slice` para saber mais sobre estas funções. Por exemplo:
  - `slice_min` e `slice_max` podem considerar ou não empates.
  - Você pode especificar uma proporção de linhas (usando `prop`) em vez da quantidade de linhas (`n`).
  - Você pode fazer amostragem com reposição, ou com probabilidades diferentes para cada linha.

### 3.4.4

#### Ordenando linhas: `arrange`

- Por título, sem repetições:

```
bb_tidy %>%
  select(musica) %>%
  distinct() %>%
  arrange(musica)
```

```
## # A tibble: 316 x 1
##   musica
##   <chr>
## 1 (Hot S**t) Country G...
## 2 3 Little Words
## 3 911
## 4 A Country Boy Can Su...
## 5 A Little Gasoline
## 6 A Puro Dolor (Purest...
## # ... with 310 more rows
```

- Por título, sem repetições, em ordem inversa:

```
bb_tidy %>%
  select(musica) %>%
  distinct() %>%
  arrange(desc(musica))
```

```
## # A tibble: 316 x 1
##   musica
##   <chr>
## 1 Your Everything
## 2 You're A God
## 3 You'll Always Be Lov...
## 4 You Won't Be Lonely ...
## 5 You Should've Told M...
## 6 You Sang To Me
## # ... with 310 more rows
```

### 3.4.5

#### Contando linhas: `count`

- Quantas semanas cada artista ficou nos *top 100*? Duas músicas na mesma semana contam como duas semanas.

```
bb_tidy %>%  
  count(artista, sort = TRUE)
```

```
## # A tibble: 228 x 2  
##   artista      n  
##   <chr>      <int>  
## 1 Creed      104  
## 2 Lonestar    95  
## 3 Destiny's Child 92  
## 4 N'Sync      74  
## 5 Sisqo       74  
## 6 3 Doors Down 73  
## # ... with 222 more rows
```

- Quantas semanas cada música ficou nos *top 100*?

```
bb_tidy %>%  
  count(musica, sort = TRUE)
```

```
## # A tibble: 316 x 2  
##   musica      n  
##   <chr>      <int>  
## 1 Higher     57  
## 2 Amazed     55  
## 3 Breathe    53  
## 4 Kryptonite 53  
## 5 With Arms Wide Open 47  
## 6 I Wanna Know 44  
## # ... with 310 more rows
```

- Se houve músicas com o mesmo nome, mas de artistas diferentes, o código acima está errado. O certo é

```
bb_tidy %>%  
  count(musica, artista, sort = TRUE)
```

```
## # A tibble: 317 x 3  
##   musica      artista      n  
##   <chr>      <chr>      <int>  
## 1 Higher     Creed        57  
## 2 Amazed     Lonestar     55  
## 3 Breathe    Hill, Faith  53
```

```
## 4 Kryptonite          3 Doors Down    53
## 5 With Arms Wide Open Creed          47
## 6 I Wanna Know       Joe              44
## # ... with 311 more rows
```

De fato, há uma diferença de uma linha.

### 3.4.5.1

#### Exercício

- Ache o título da música que tem dois artistas diferentes.

**Sugestão:** conte por música e artista primeiro, depois só por música.

### 3.4.6

#### Agrupando linhas: `group_by` e `summarize`

- Qual foi a melhor posição que cada artista alcançou?

```
bb_tidy %>%
  group_by(artista) %>%
  summarize(melhor = min(pos)) %>%
  arrange(melhor)
```

```
## # A tibble: 228 x 2
##   artista          melhor
##   <chr>            <dbl>
## 1 Aaliyah          1
## 2 Aguilera, Christina 1
## 3 Carey, Mariah     1
## 4 Creed            1
## 5 Destiny's Child   1
## 6 Iglesias, Enrique 1
## # ... with 222 more rows
```

- Qual foi a melhor posição que cada música alcançou?

```
bb_tidy %>%
  group_by(artista, musica) %>%
  summarize(melhor = min(pos)) %>%
  arrange(melhor)
```

```
## `summarise()` has grouped output by 'artista'. You can override using the `.groups`
## argument.
```

```
## # A tibble: 317 x 3
##   artista          musica          melhor
##   <chr>            <chr>            <dbl>
```



```
## 1 Aaliyah          Try Again          1
## 2 Aguilera, Christina Come On Over Baby (A... 1
## 3 Aguilera, Christina What A Girl Wants      1
## 4 Carey, Mariah    Thank God I Found Yo... 1
## 5 Creed            With Arms Wide Open    1
## 6 Destiny's Child Independent Women Pa... 1
## # ... with 311 more rows
```

- Quando usamos `summarize`, só o agrupamento **mais interno** é desfeito. Isto significa que **o resultado acima ainda está agrupado por artista**.
- Quantas semanas cada artista ficou na posição 1?

A função `n()` é uma maneira conveniente de **obter o número de linhas de um grupo** (ou, se não houver agrupamento, de toda a *tibble*); mas `n()` só pode ser chamada **em certos contextos**, como `summarize()` ou `mutate()`.

```
bb_tidy %>%
  filter(pos == 1) %>%
  group_by(artista) %>%
  summarize(semanas = n()) %>%
  arrange(desc(semanas))
```

```
## # A tibble: 15 x 2
##   artista          semanas
##   <chr>           <int>
## 1 Destiny's Child      14
## 2 Santana             10
## 3 Aguilera, Christina   6
## 4 Madonna             4
## 5 Savage Garden        4
## 6 Iglesias, Enrique     3
## # ... with 9 more rows
```

- Perceba que `count`, que vimos mais acima, faz agrupamentos do mesmo modo:

```
bb_tidy %>%
  filter(pos == 1) %>%
  count(artista, sort = TRUE)
```

```
## # A tibble: 15 x 2
##   artista          n
##   <chr>           <int>
## 1 Destiny's Child      14
## 2 Santana             10
## 3 Aguilera, Christina   6
## 4 Madonna             4
## 5 Savage Garden        4
## 6 Iglesias, Enrique     3
```

```
## # ... with 9 more rows
```

- Uma pergunta diferente: quais são os artistas cujas músicas apareceram no *top 100* mais tempo depois do lançamento da música?

```
bb_tidy %>%
  group_by(artista) %>%
  summarize(semanas = max(semana)) %>%
  arrange(desc(semanas))
```

```
## # A tibble: 228 x 2
##   artista      semanas
##   <chr>         <int>
## 1 Creed          65
## 2 Lonestar       64
## 3 3 Doors Down   53
## 4 Hill, Faith    53
## 5 Joe            44
## 6 Vertical Horizon 41
## # ... with 222 more rows
```

- Qual a posição média de cada música? Lembre-se de que eliminamos as linhas com NA; logo, **a média vai ser sobre a quantidade de semanas em que a música esteve na lista.**

```
media1 <- bb_tidy %>%
  group_by(artista, musica) %>%
  summarize(media = mean(pos), .groups = 'drop') %>%
  arrange(media)
```

```
media1
```

```
## # A tibble: 317 x 3
##   artista      musica      media
##   <chr>         <chr>    <dbl>
## 1 "Santana"      Maria, Maria    10.5
## 2 "Madonna"      Music          13.5
## 3 "N'Sync"       Bye Bye Bye     14.3
## 4 "Elliott, Missy \"Misdemeanor\"" Hot Boyz        14.3
## 5 "Destiny's Child" Independent Women Pa... 14.8
## 6 "Iglesias, Enrique" Be With You     15.8
## # ... with 311 more rows
```

- E se quisermos **a média sobre o número de semanas desde a entrada da música até a última semana** em que a música apareceu na lista?

```
media2 <- bb_tidy %>%
```

```
group_by(artista, musica) %>%
  summarize(media = sum(pos)/max(semana), .groups = 'drop') %>%
  arrange(media)
```

```
media2
```

```
## # A tibble: 317 x 3
##   artista                musica                media
##   <chr>                  <chr>                <dbl>
## 1 "Santana"              Maria, Maria           10.5
## 2 "Madonna"             Music                 13.5
## 3 "N'Sync"              Bye Bye Bye           14.3
## 4 "Elliott, Missy \"Misdemeanor\"" Hot Boyz              14.3
## 5 "Destiny's Child"     Independent Women Pa... 14.8
## 6 "Iglesias, Enrique"   Be With You           15.8
## # ... with 311 more rows
```

As primeiras linhas são iguais, mas os resultados são diferentes:

```
identical(media1, media2)
```

```
## [1] FALSE
```

## 3.5

### Exercícios

1. Vamos trabalhar com um conjunto de dados sobre super-heróis.

Carregue o tidyverse com o comando

```
library(tidyverse)
```

Execute o seguinte comando para ler os dados para uma *tibble*:

```
arquivo <- paste0(
  'https://github.com/fnaufel/',
  'probestr/raw/master/data/',
  'heroes_information.csv'
)

herois_info <- read_csv(
  arquivo,
  na = c(' ', '-', 'NA')
) %>%
  # Eliminar a primeira coluna (números de série)
  select(-1) %>%
  # Renomear colunas restantes
```

```

rename(
  nome = name,
  sexo = Gender,
  olhos = 'Eye color',
  raça = Race,
  cabelos = 'Hair color',
  altura = Height,
  editora = Publisher,
  pele = 'Skin color',
  lado = Alignment,
  peso = Weight
)

```

2. Quantas linhas tem a *tibble*?
3. Existem heróis que aparecem em mais de uma linha?
4. Quantas editoras diferentes existem na *tibble*? Liste-as em ordem decrescente de quantidade de heróis.
5. Vamos colocar todas as editores menores em uma classe só.

Na coluna *editora*, substitua

- ‘Marvel Comics’ por ‘Marvel’,
- ‘DC Comics’ por ‘DC’, e
- todas as outras editoras pelo termo ‘Outras’.

**Dica:** use a função `case_when()`, do tidyverse.

6. Confira, novamente, a quantidade de valores diferentes na coluna *editora*.
7. Existem heróis sem informação de editora. Quantos? Quais são?
8. Altere novamente a coluna *editora*, colocando o valor ‘Outras’ para os heróis sem informação de editora. Use a função `if_else()` (com *underscore*, não a função `ifelse`).
9. Confira, mais uma vez, a quantidade de valores diferentes na coluna *editora*.
10. Existem heróis sem informação de sexo? Quantos? Para estes heróis, coloque o valor ‘Desconhecido’ na coluna *sexo*.
11. Qual a altura mínima? Qual a altura máxima? Substitua as alturas negativas por `NA`.
12. Qual o peso mínimo? Qual o peso máximo? Substitua os pesos negativos por `NA`.
13. Qual é o peso médio de todos os heróis? Ignore os valores `NA`.
14. Qual é a altura média de todos os heróis? Ignore os valores `NA`.
15. Qual é a altura média dos heróis, por editora? Ignore os valores `NA`.
16. Quais são os 3 heróis mais altos de cada sexo?
17. Quais são as 3 cores de olhos mais comuns para cada sexo?
18. Liste, por editora, as quantidades de heróis do bem, do mal, e neutros.

19. Quantas raças diferentes existem?

20. Qual a quantidade de raças diferentes de cada editora?

21. **DESAFIO:** Liste as raças que só pertencem a uma única editora.

Existem várias maneiras de fazer isto. Experimente várias, até achar uma que seja mais elegante.

---

### Visualização com ggplot2

---



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

#### 4.1

---

##### Vídeo 1

<https://youtu.be/OBpNjqIIyhI>

#### 4.2

---

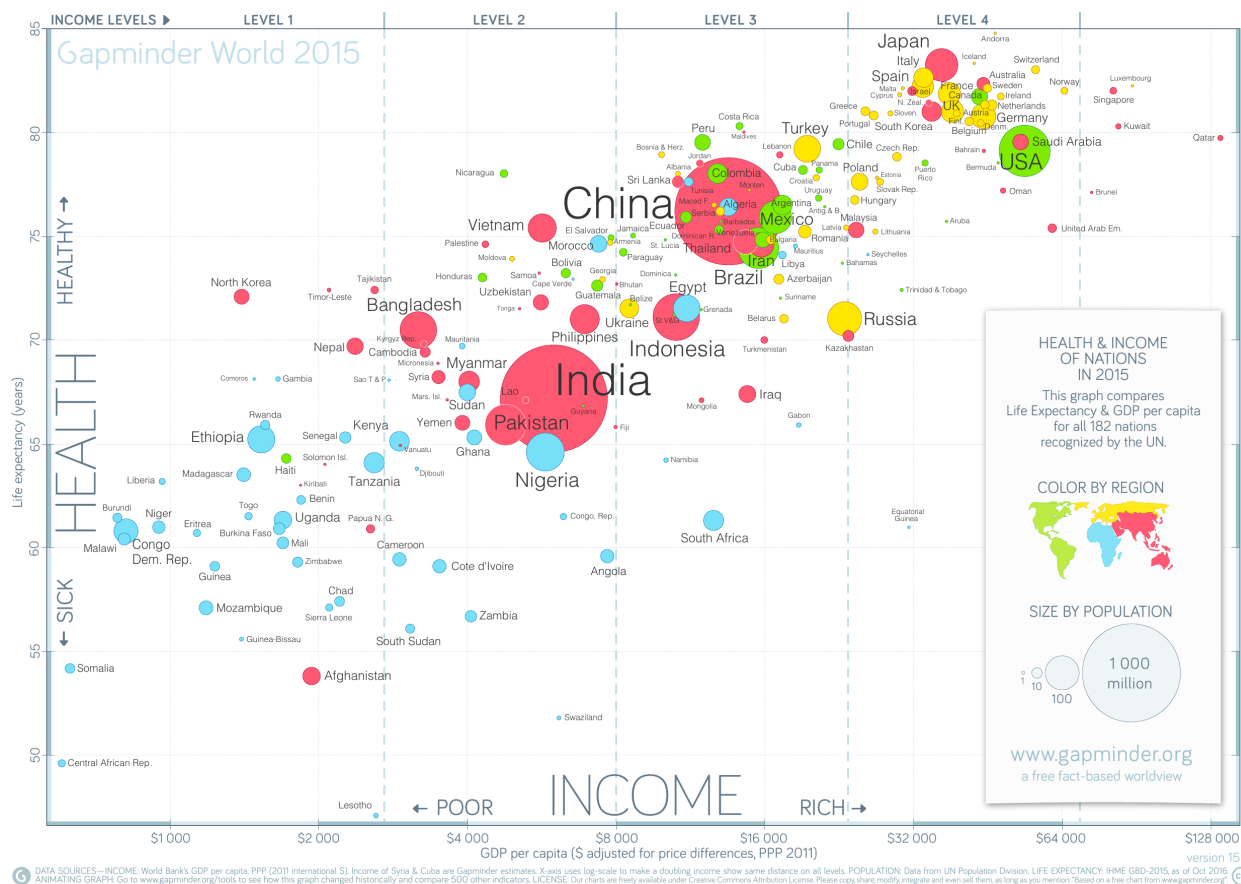
### Componentes de um gráfico ggplot2

#### 4.2.1

---

##### 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.

#### 4.2.2

---

##### 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.

#### 4.2.3

---

##### 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.

#### 4.2.4

---

##### 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.

## 4.3

### 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      Acin~ carni Carn~ lc             12.1      NA
## 2 Owl monkey   Aotus omni Prim~ <NA>          17        1.8
## 3 Mountain beaver Aplo~ herbi Rode~ nt             14.4      2.4
## 4 Greater short-t~ Blar~ omni Sori~ lc             14.9      2.3
## 5 Cow          Bos   herbi Arti~ domesticated      4        0.7
## 6 Three-toed sloth Brad~ herbi Pilo~ <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   Acin~ carn~ Carn~ lc             12.1        NA
## 2 Owl monkey Aotus omni  Prim~ <NA>          17          1.8
## 3 Mountain beaver Aplo~ herbi Rode~ nt             14.4        2.4
## 4 Greater short-t~ Blar~ omni  Sori~ lc             14.9        2.3
## 5 Cow       Bos   herbi Arti~ domesticated    4          0.7
## 6 Three-toed sloth Brad~ herbi Pilo~ <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	1. African elephant	1 ( 1,2%)	0
[character]	2. African giant pouched rat	1 ( 1,2%)	(0,0%)
	3. African striped mouse	1 ( 1,2%)	
	4. Arctic fox	1 ( 1,2%)	
	5. Arctic ground squirrel	1 ( 1,2%)	
	6. Asian elephant	1 ( 1,2%)	
	7. Baboon	1 ( 1,2%)	
	8. Big brown bat	1 ( 1,2%)	
	9. Bottle-nosed dolphin	1 ( 1,2%)	
	10. Brazilian tapir	1 ( 1,2%)	
	[ 73 outros ]	73 (88,0%)	
genus	1. Panthera	3 ( 3,6%)	0
[character]	2. Spermophilus	3 ( 3,6%)	(0,0%)
	3. Equus	2 ( 2,4%)	
	4. Vulpes	2 ( 2,4%)	
	5. Acinonyx	1 ( 1,2%)	
	6. Aotus	1 ( 1,2%)	
	7. Aplodontia	1 ( 1,2%)	
	8. Blarina	1 ( 1,2%)	
	9. Bos	1 ( 1,2%)	
	10. Bradypus	1 ( 1,2%)	
	[ 67 outros ]	67 (80,7%)	
vore	1. carni	19 (25,0%)	7
[character]	2. herbi	32 (42,1%)	(8,4%)
	3. insecti	5 ( 6,6%)	
	4. omni	20 (26,3%)	
order	1. Rodentia	22 (26,5%)	0
[character]	2. Carnivora	12 (14,5%)	(0,0%)
	3. Primates	12 (14,5%)	
	4. Artiodactyla	6 ( 7,2%)	
	5. Soricomorpha	5 ( 6,0%)	
	6. Cetacea	3 ( 3,6%)	
	7. Hyracoidea	3 ( 3,6%)	
	8. Perissodactyla	3 ( 3,6%)	
	9. Chiroptera	2 ( 2,4%)	
	10. Cingulata	2 ( 2,4%)	
	[ 9 outros ]	13 (15,7%)	
conservation	1. cd	2 ( 3,7%)	29
[character]	2. domesticated	10 (18,5%)	(34,9%)
	3. en	4 ( 7,4%)	
	4. lc	27 (50,0%)	
	5. nt	4 ( 7,4%)	
	6. vu	7 (13,0%)	

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. Sperophilus 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%)

## 4.4

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

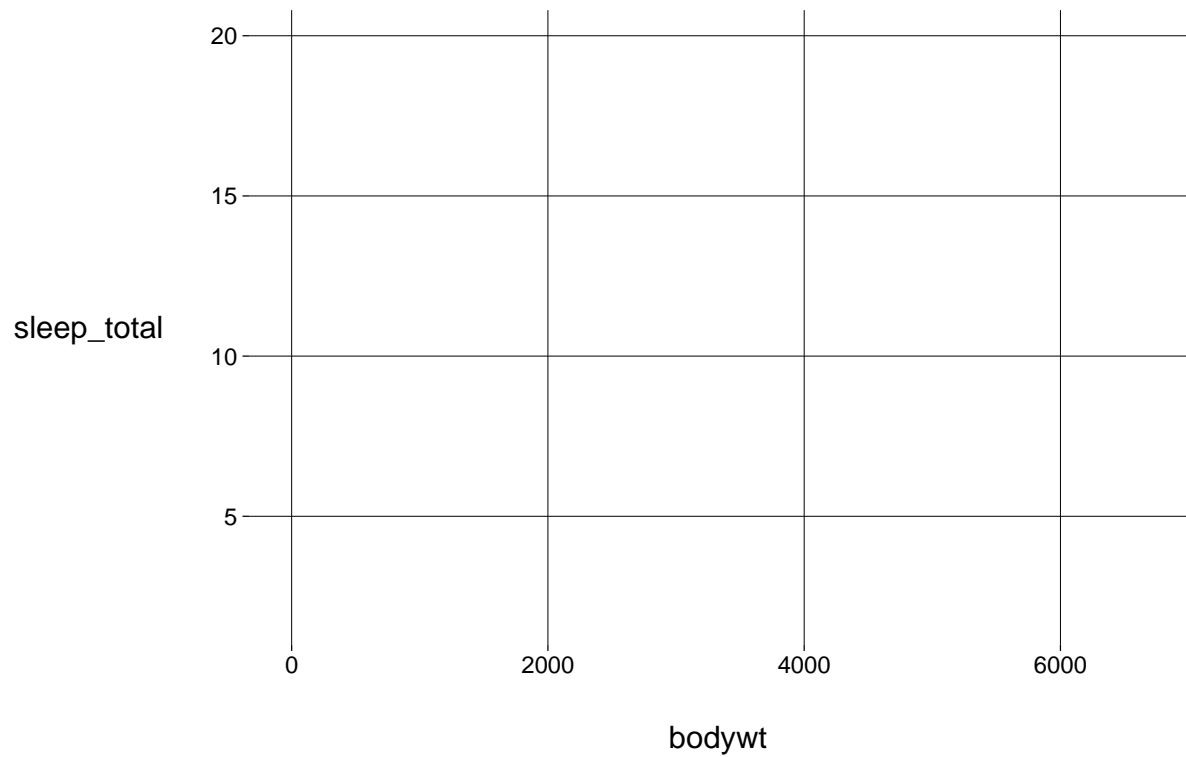
- Servem para visualizar a *relação* entre **duas variáveis quantitativas**.
- **Essa relação não é necessariamente de causa e efeito**.
- Isto é, a variável do eixo horizontal não determina, necessariamente, os valores da variável do eixo vertical.
- Pense em **associação**, **correlação**, não em causalidade.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

#### 4.4.1

#### Horas de sono e peso corporal

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

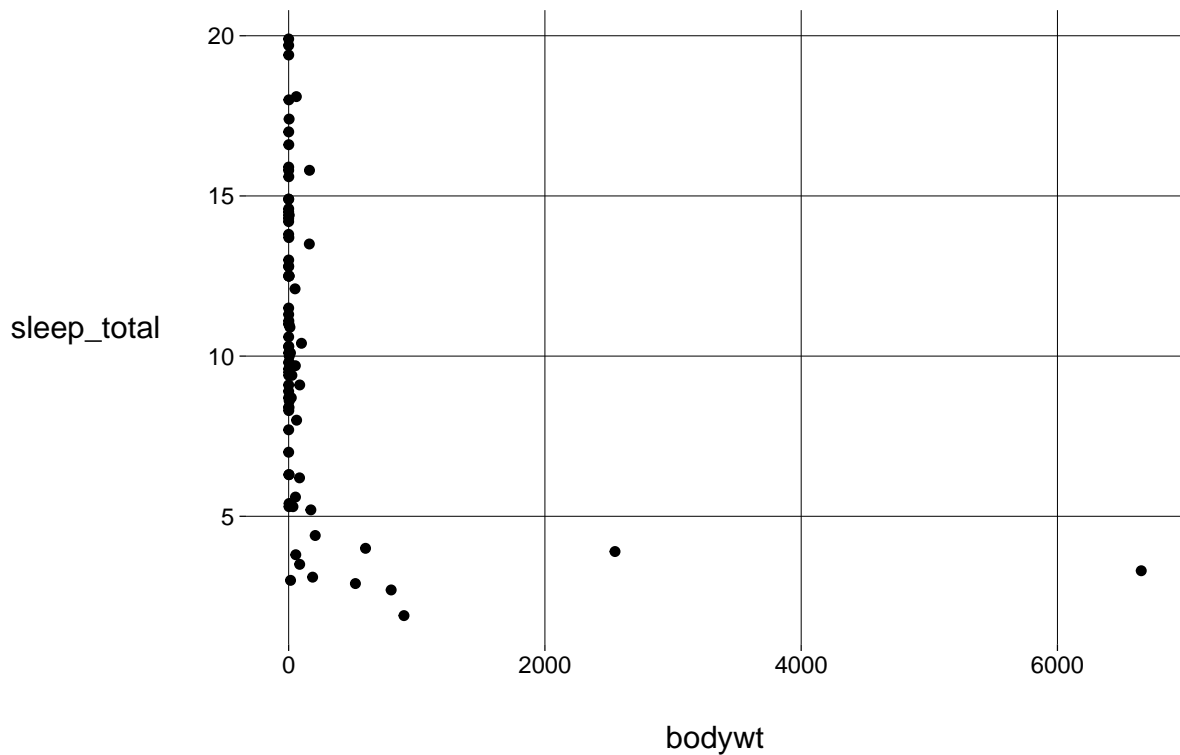
```
sono %>%
  ggplot(aes(x = bodywt, y = sleep_total))
```



- 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?

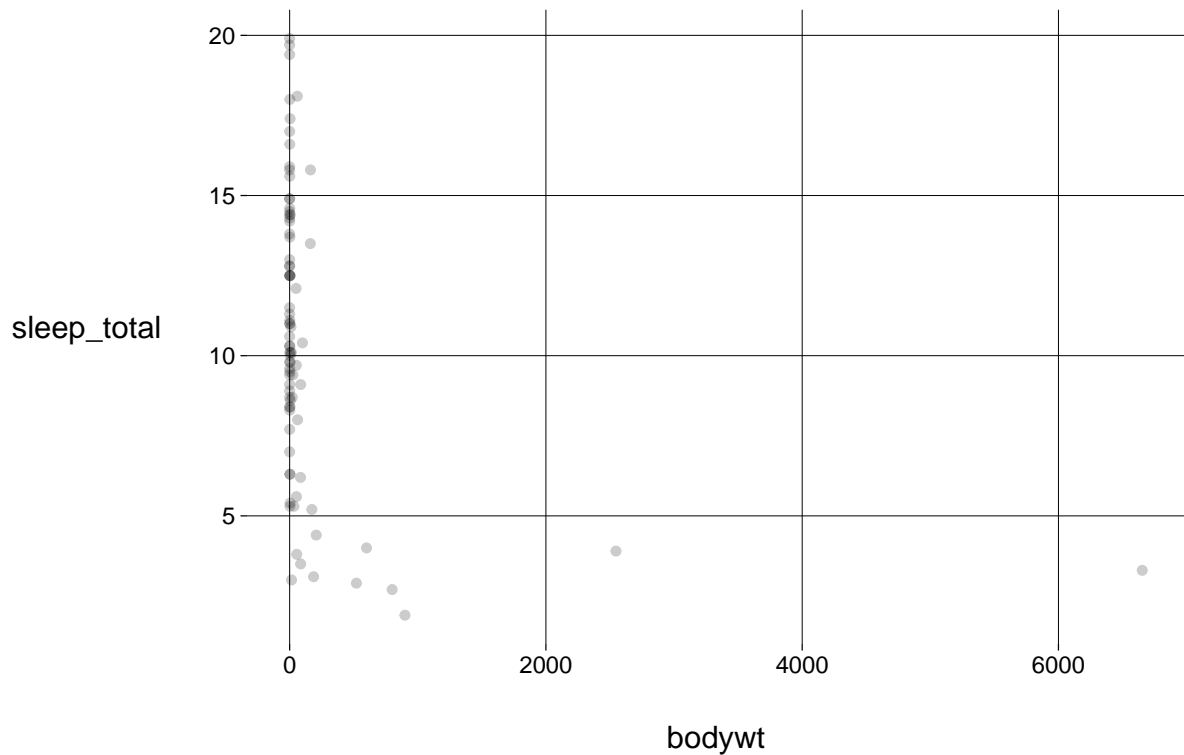
```
sono %>%
  filter(bodywt > 250) %>%
  select(name, bodywt) %>%
  arrange(bodywt)
```

```
## # 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.

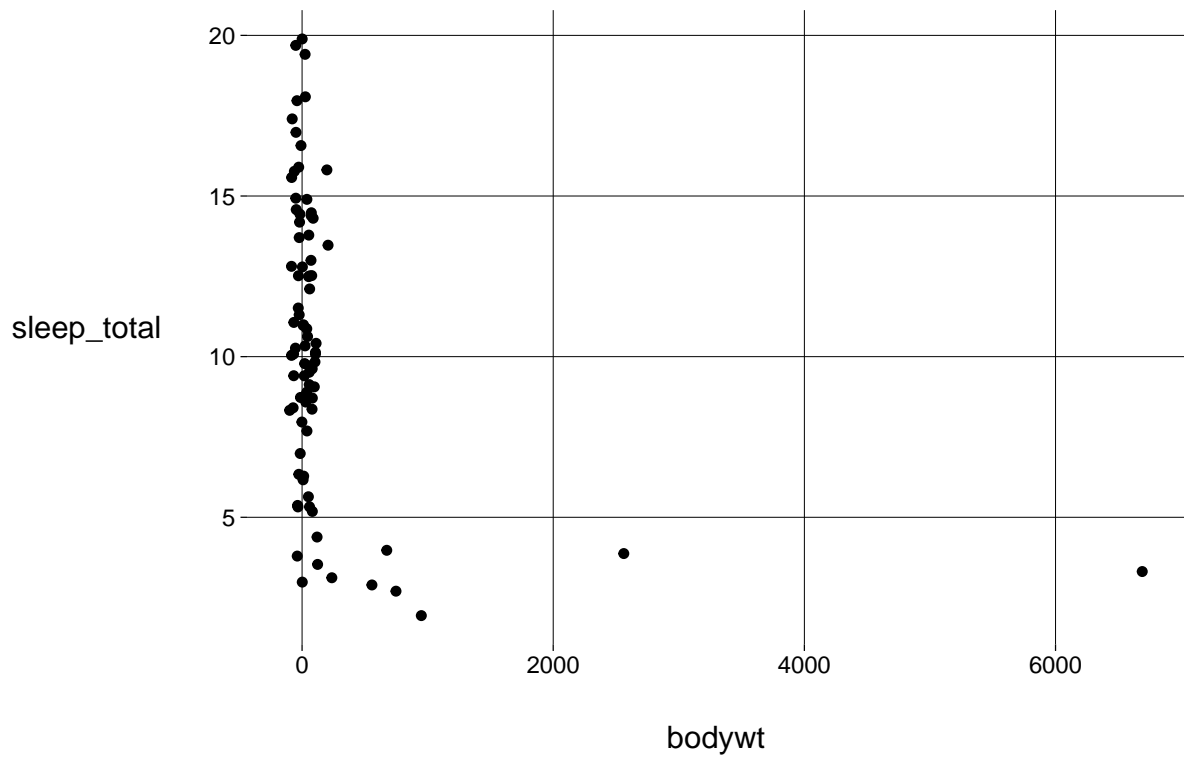
- A primeira delas é **alterando a opacidade dos pontos**. Isto é um ajuste na geometria apenas, pois a opacidade, aqui, não representa informação nenhuma.

```
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 “tremer”. As posições dos pontos são ligeiramente perturbadas, para evitar colisões. Perde-mos precisão, mas a visualização fica melhor.

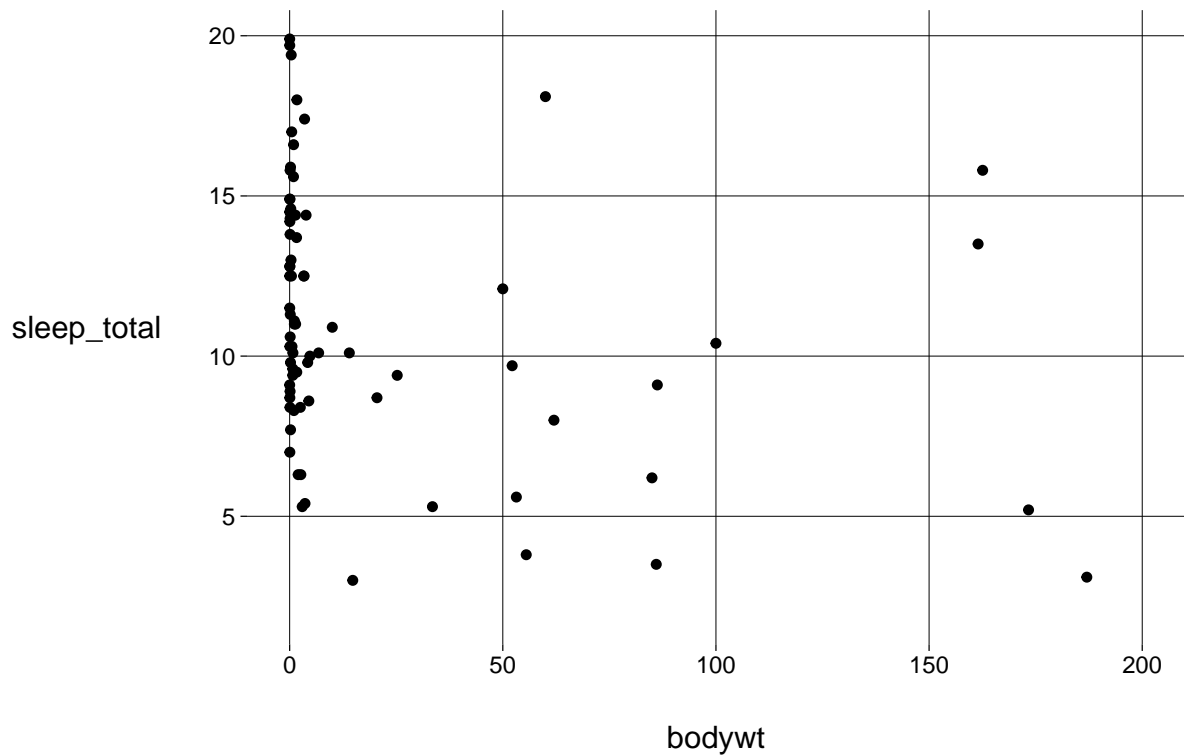
```
sono %>%  
  ggplot(aes(x = bodywt, y = sleep_total)) +  
    geom_jitter(width = 100)
```



- 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.

#### 4.4.2

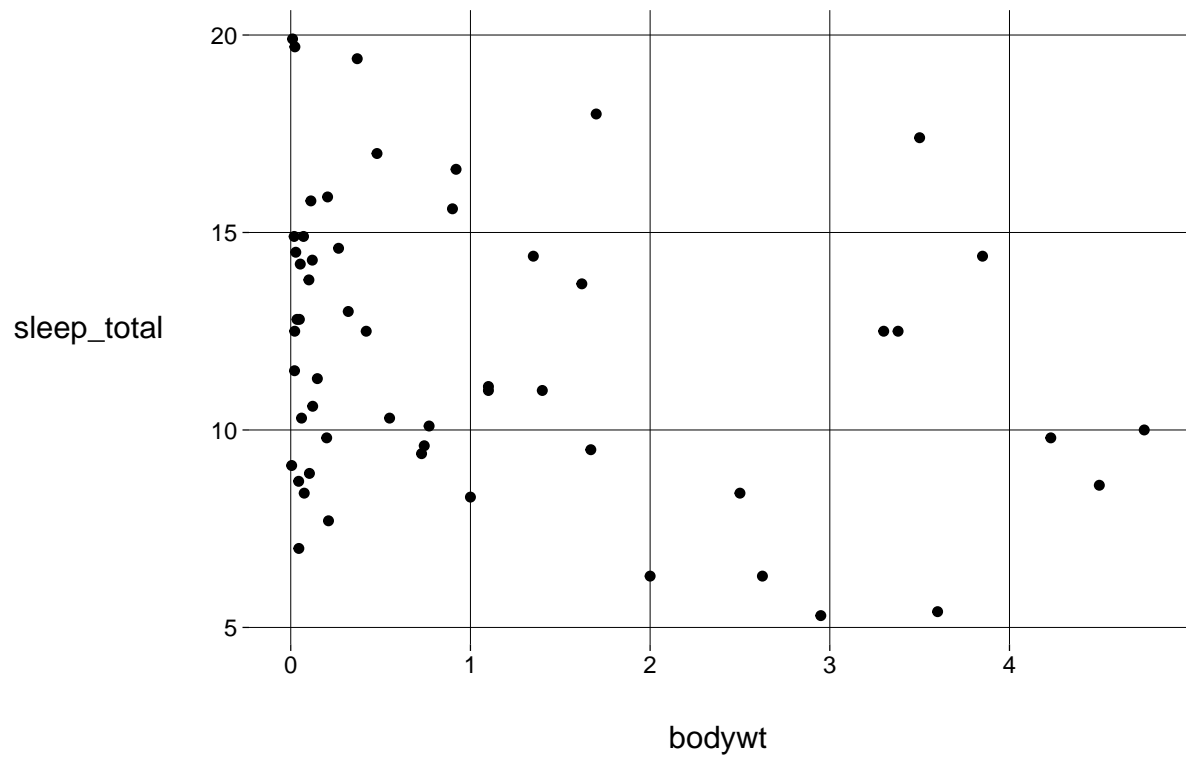
##### Horas de sono e peso corporal para animais pequenos

- 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()
```

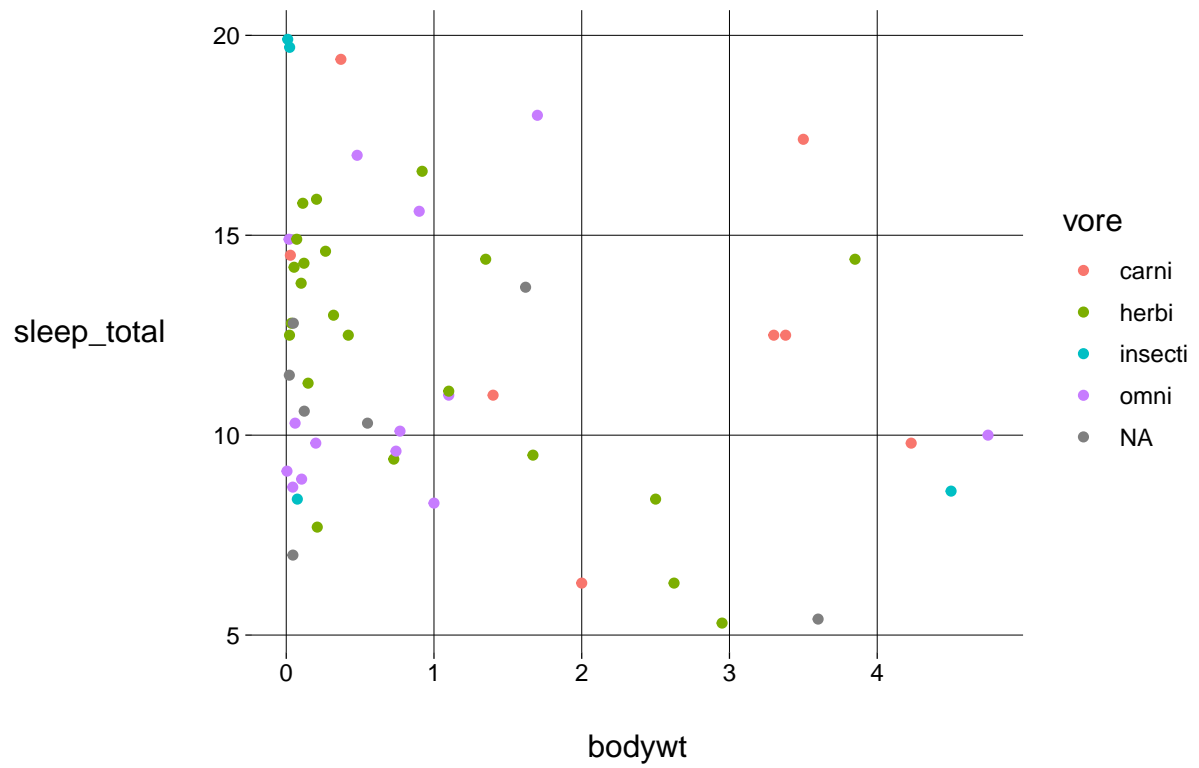


#### 4.4.3

##### 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()
```

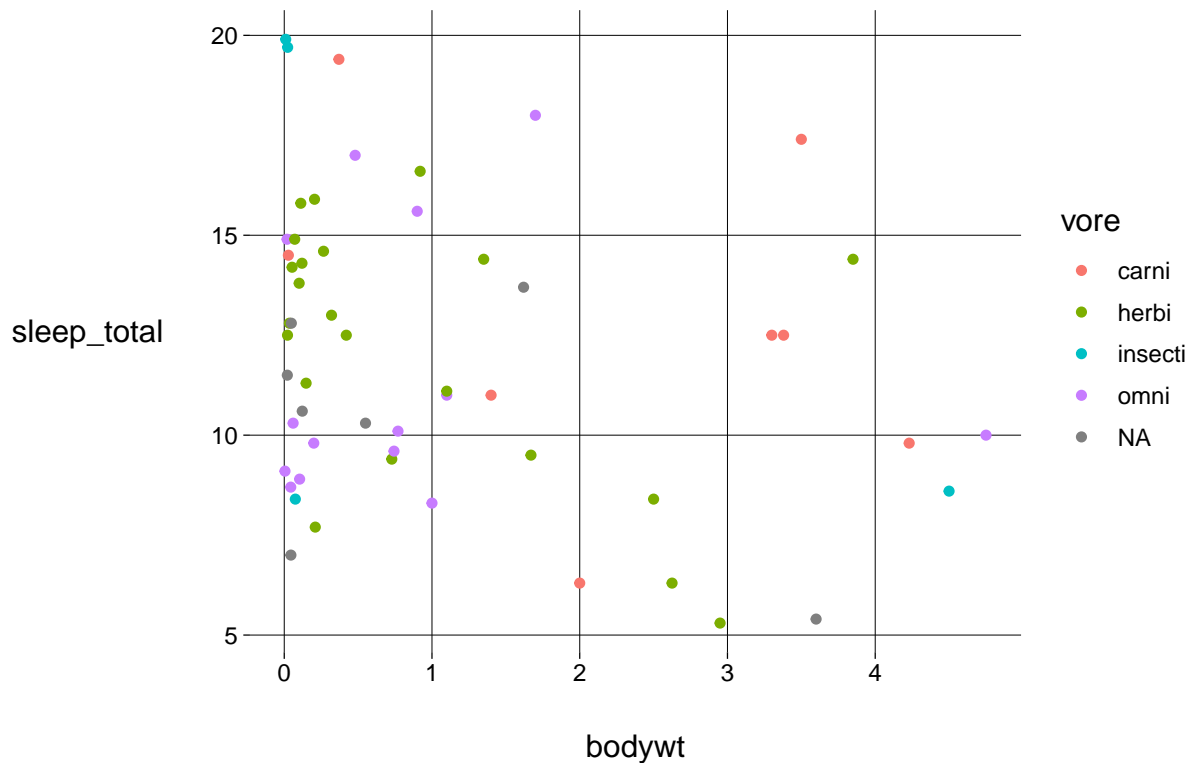


#### 4.4.4

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))
```



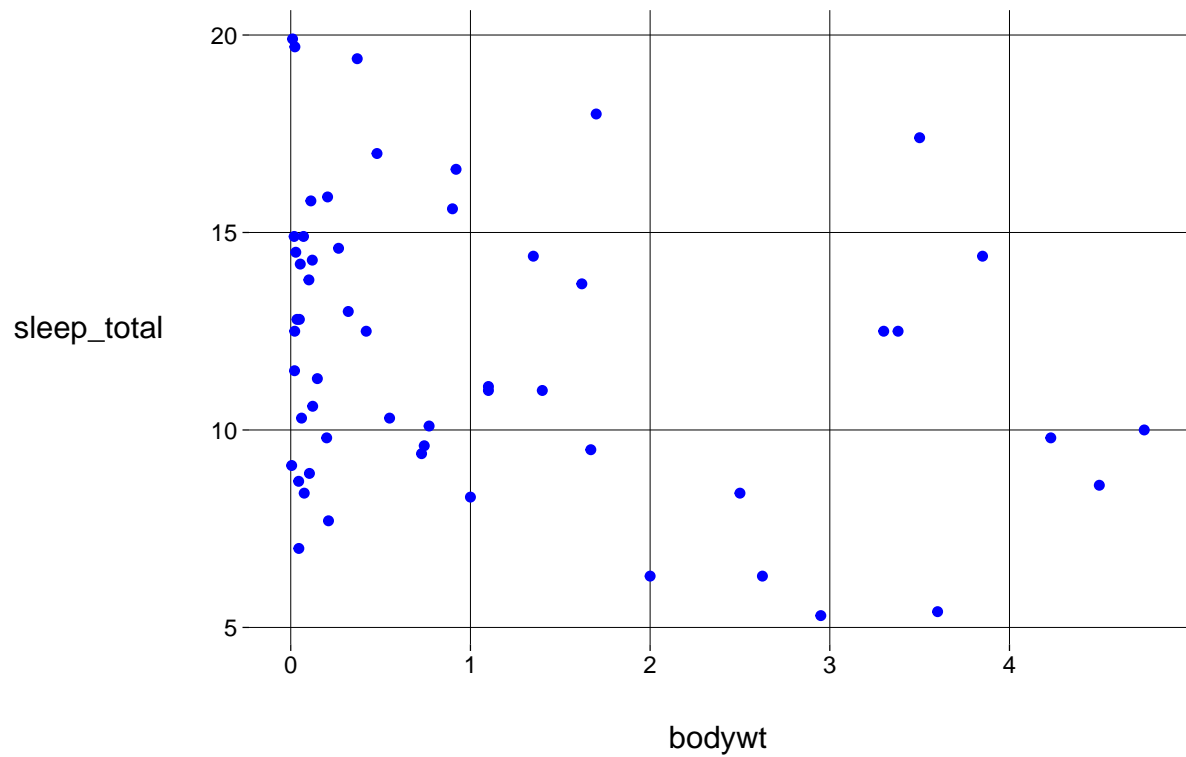
- Fazendo deste modo, a estética só vale para uma geometria. Se você acrescentar outras geometrias (linhas, por exemplo), a estética não valerá para elas.

#### 4.4.5

##### 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:

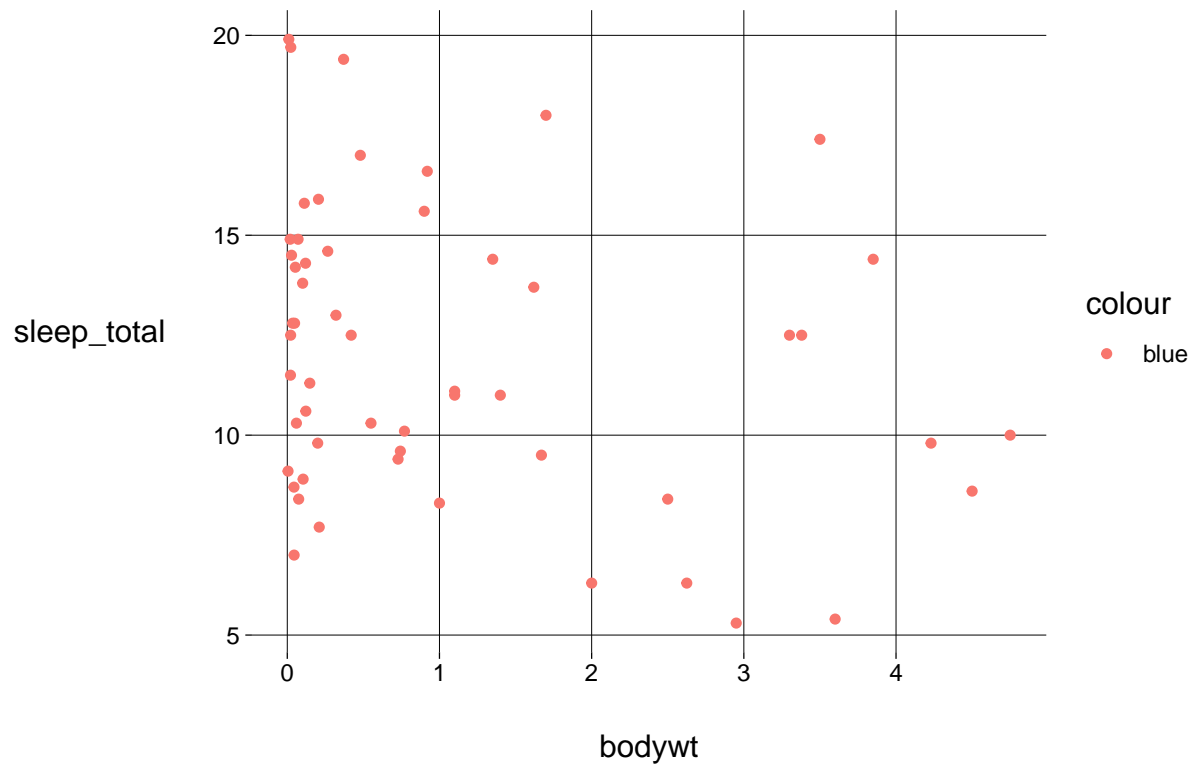
```
sono %>%
  filter(bodywt < limite) %>%
  ggplot() +
    geom_point(aes(x = bodywt, y = sleep_total), color = 'blue')
```



- 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'))
```





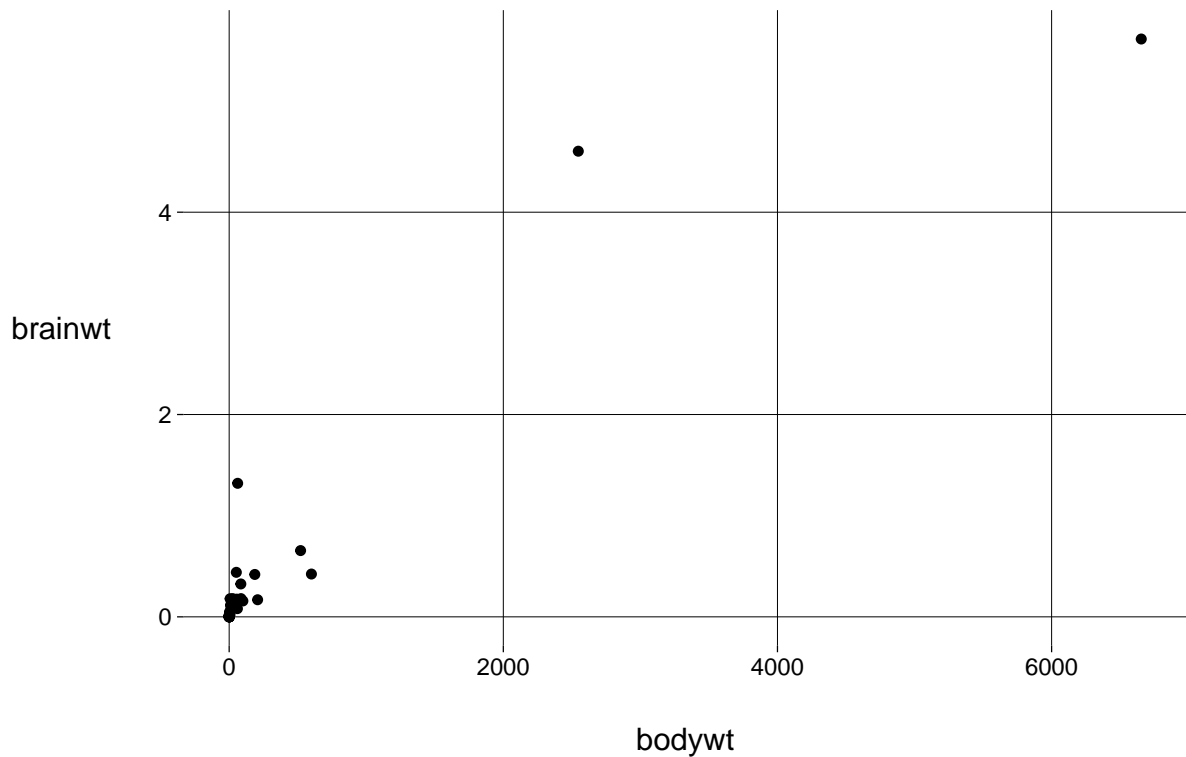
#### 4.4.6

##### 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()
```

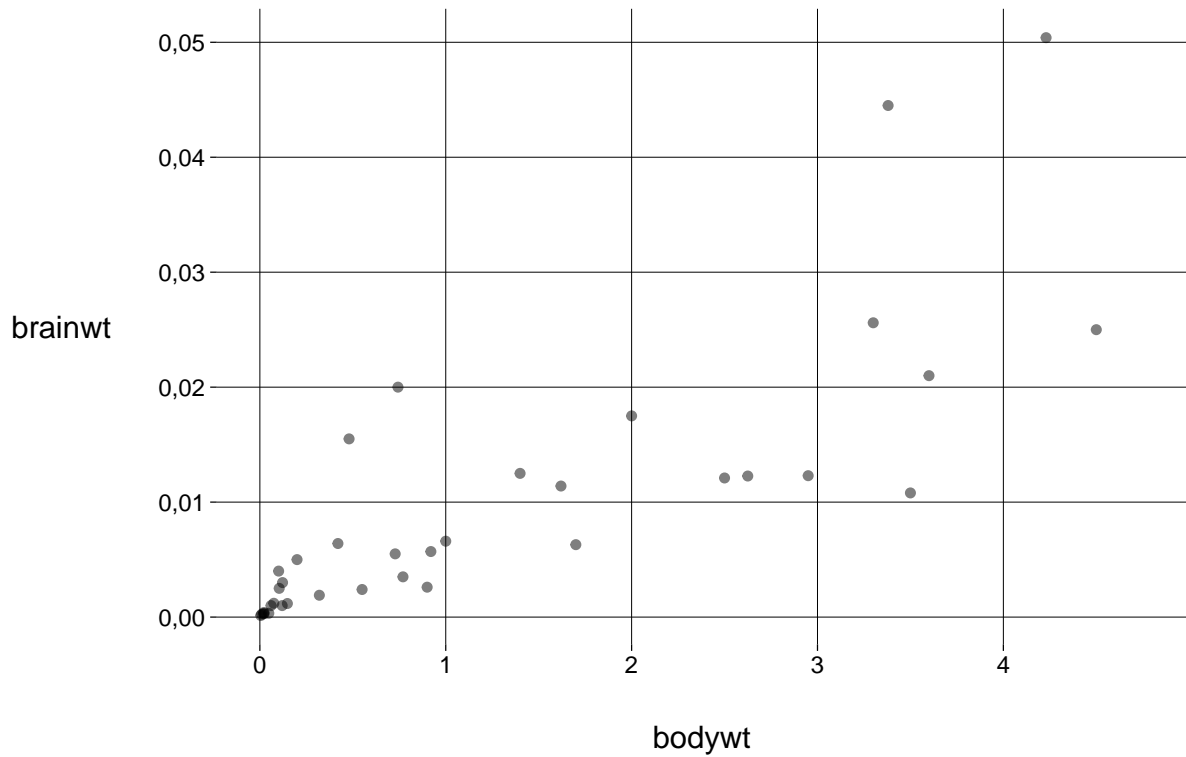
```
## # A tibble: 1 x 1
##       n
##   <int>
## 1    27
```

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

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

```
geom_point(alpha = .5)
```

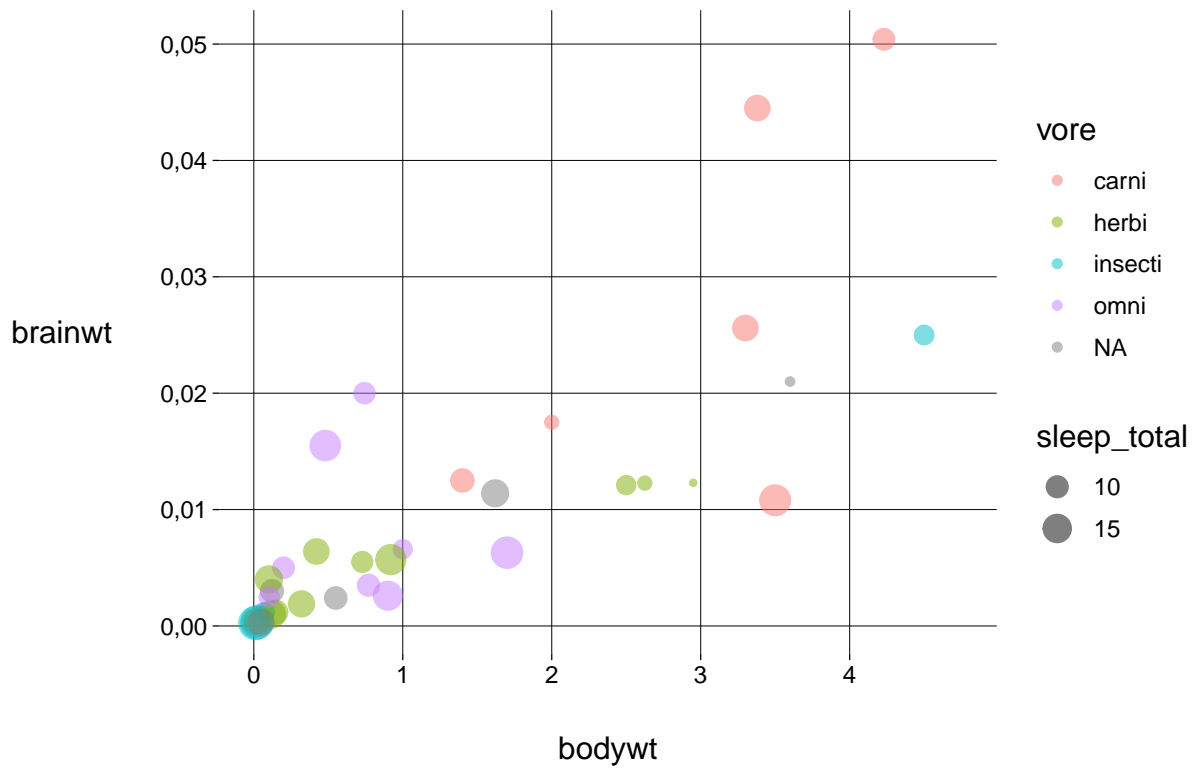
```
## Warning: Removed 18 rows containing missing values (geom_point).
```



- 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)',
  )
```

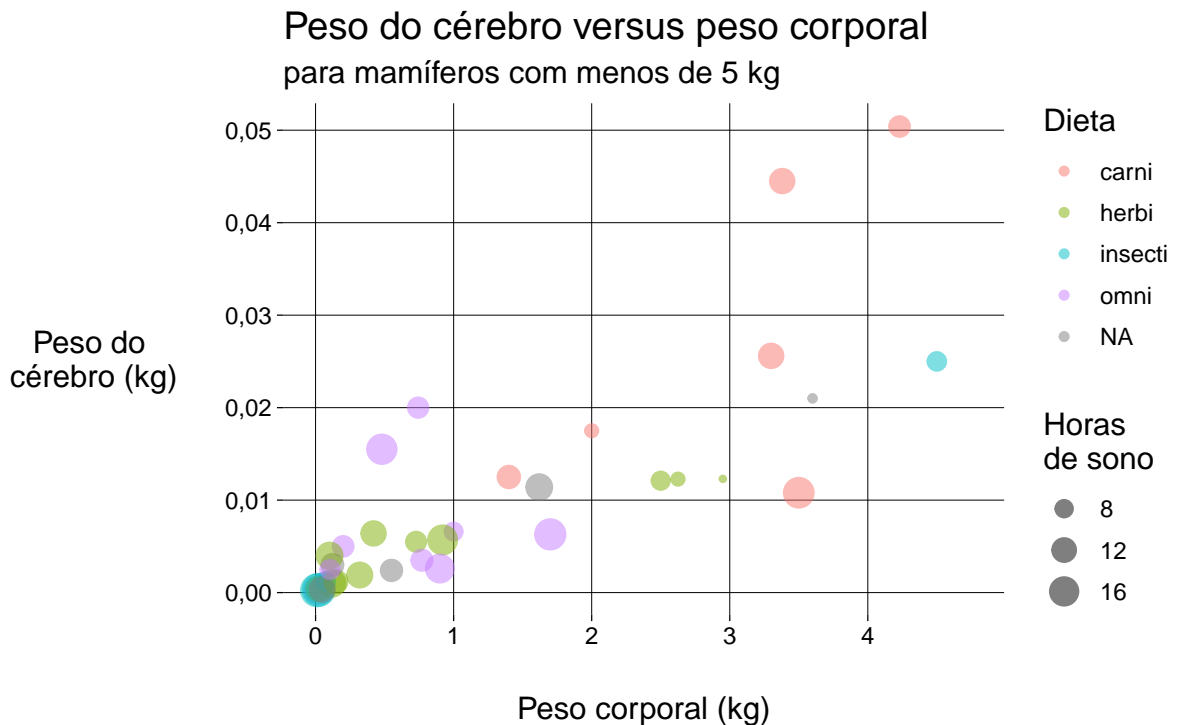
```

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.

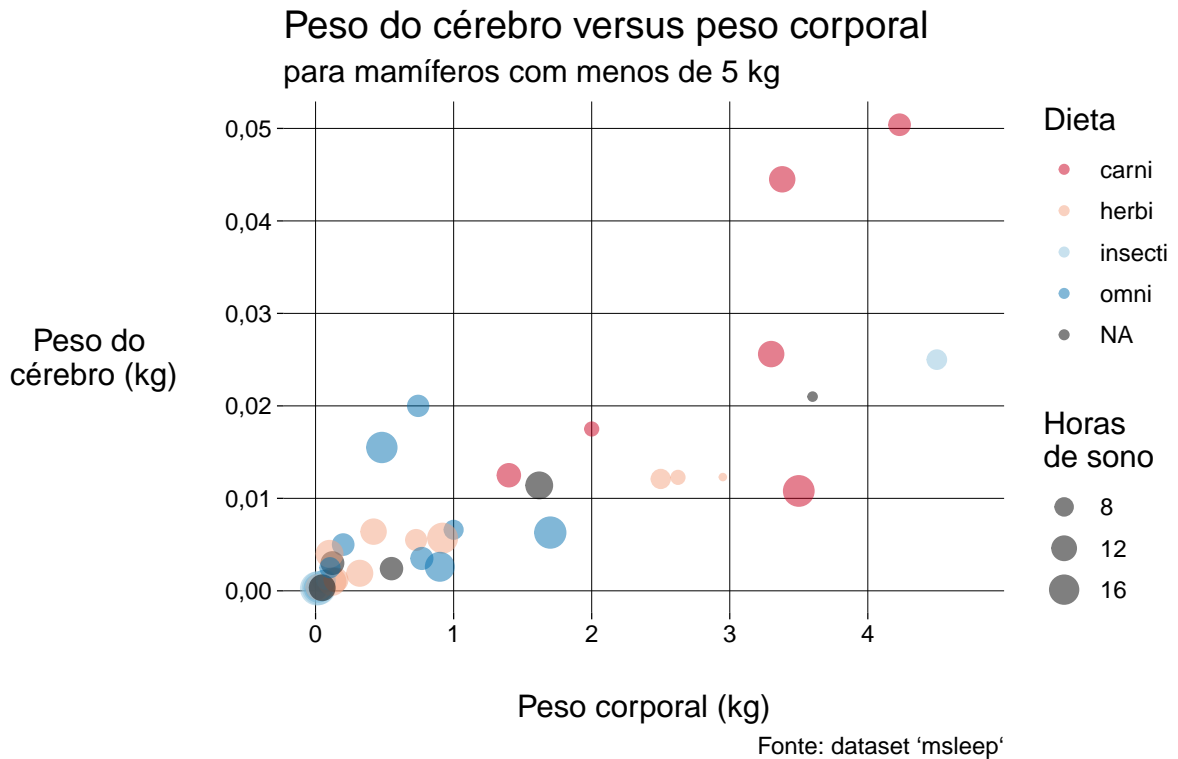
```

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

```

grafico2

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



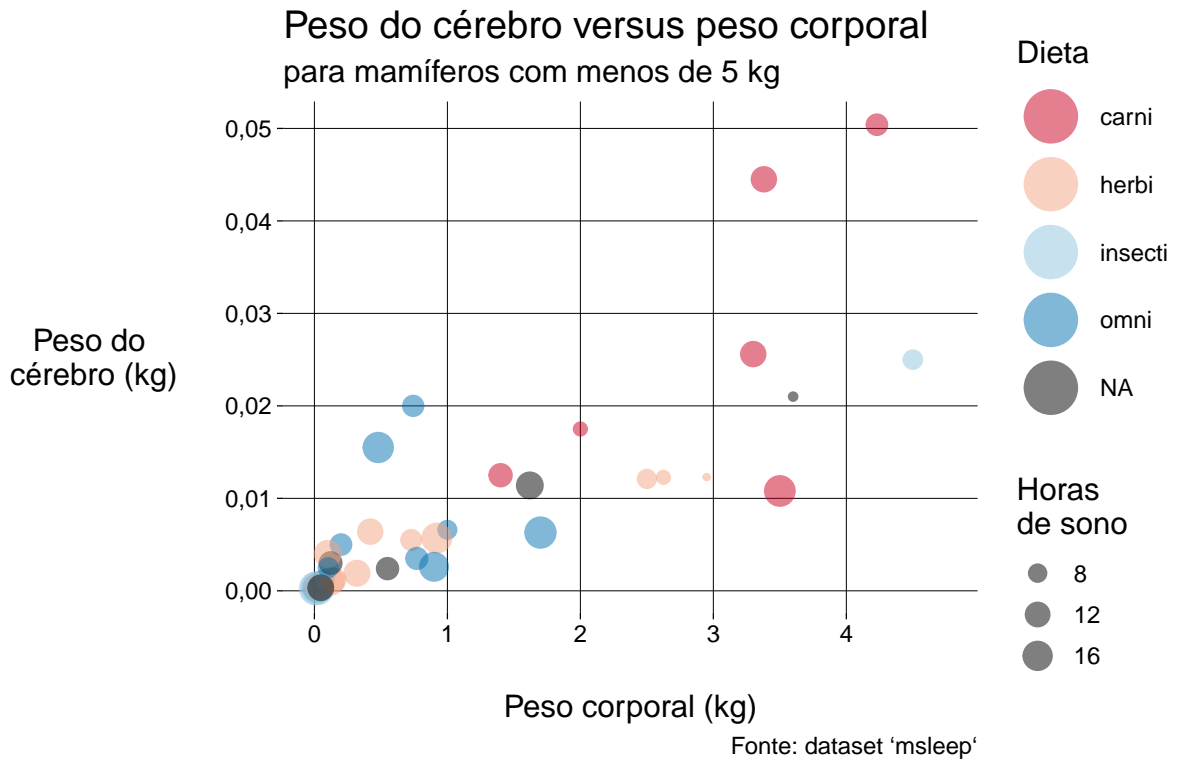
- Observe como usamos o gráfico já salvo na variável `grafico` e simplesmente acrescentamos a nova escala. Este tipo de “montagem” de gráficos `ggplot2` é bem conveniente, para evitar repetição de código.
- Um último ajuste na aparência: os pontos na legenda “Dieta” estão pequenos demais. Quase não identificamos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

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

grafico3
```

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

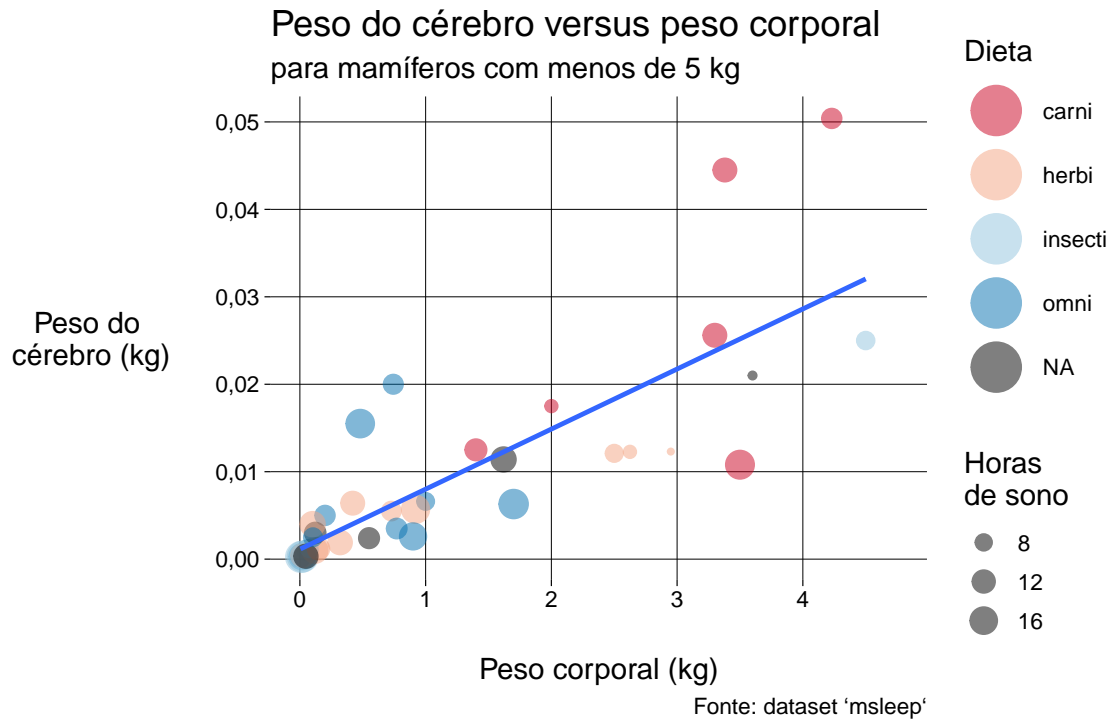


- Agora podemos finalmente comentar sobre a informação que o gráfico mostra sobre os dados:
  - De fato, existe uma correlação entre peso cerebral e peso corporal: quanto maior o peso corporal, maior o peso cerebral. Nada surpreendente.
  - Podemos fazer o `ggplot2` traçar uma reta de regressão com a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

## 4.5

---

### Vídeo 2

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

## 4.6

---

### Histogramas e cia.

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



#### 4.6.1

### 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
```

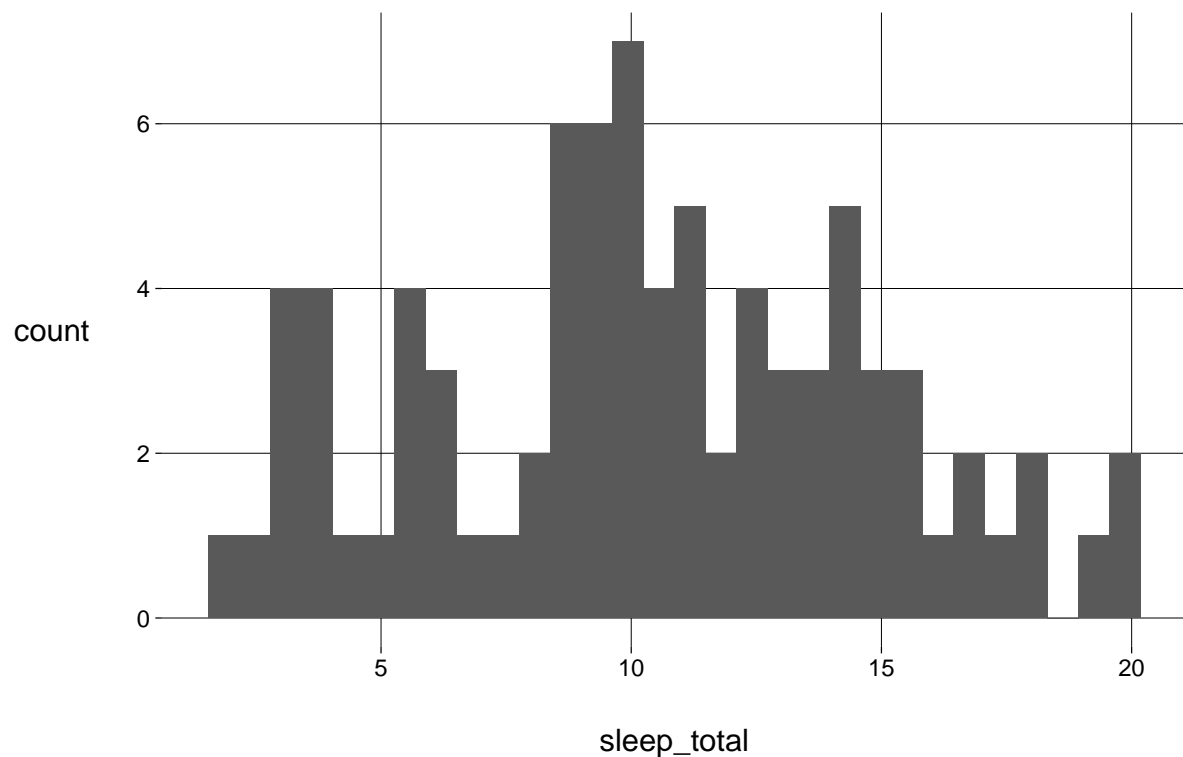
## 4.6.2

### Histograma

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

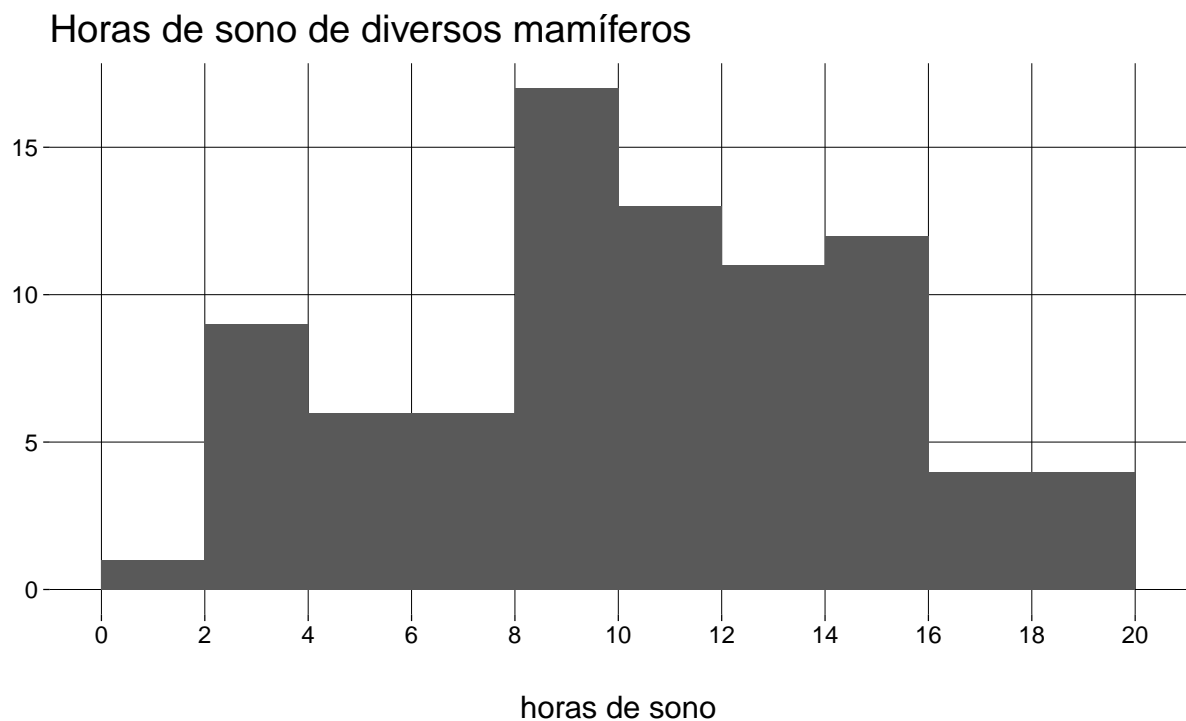
```
sono %>%
  ggplot(aes(x = sleep_total)) +
  geom_histogram()
```

## ``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)
```

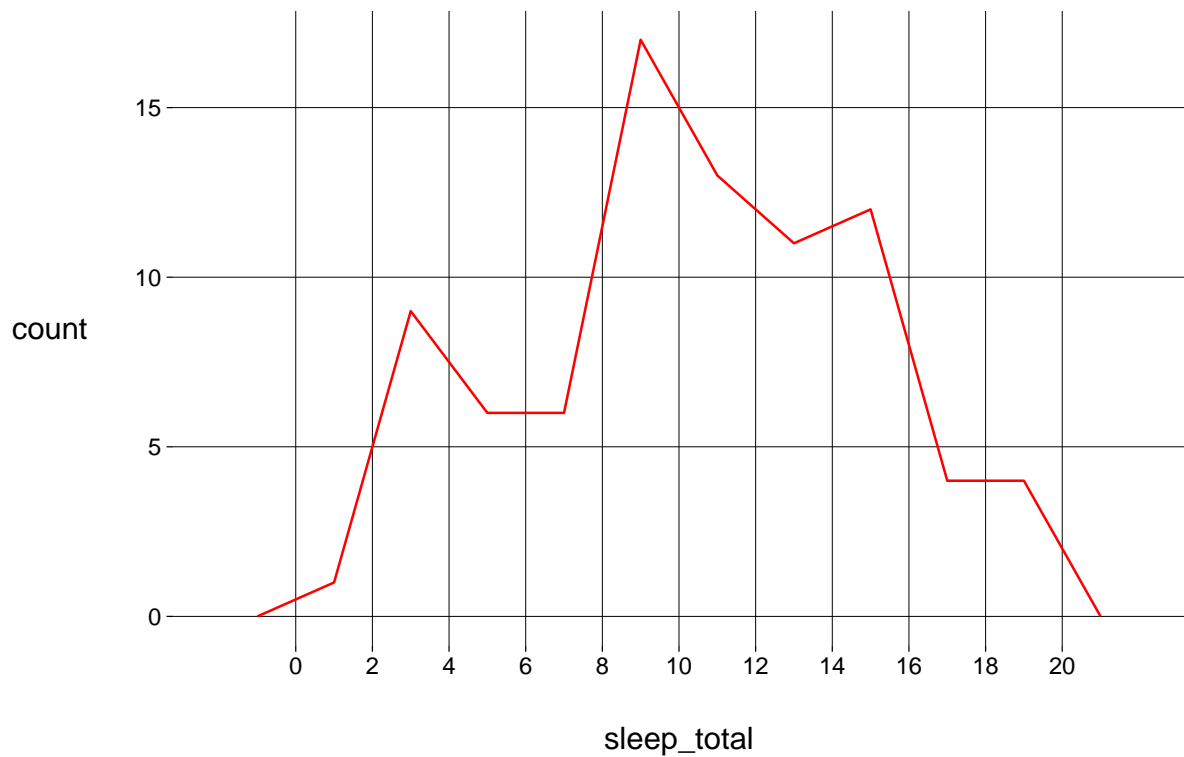
```
## [1] 10,43373
```

### 4.6.3

#### Polígono de frequência

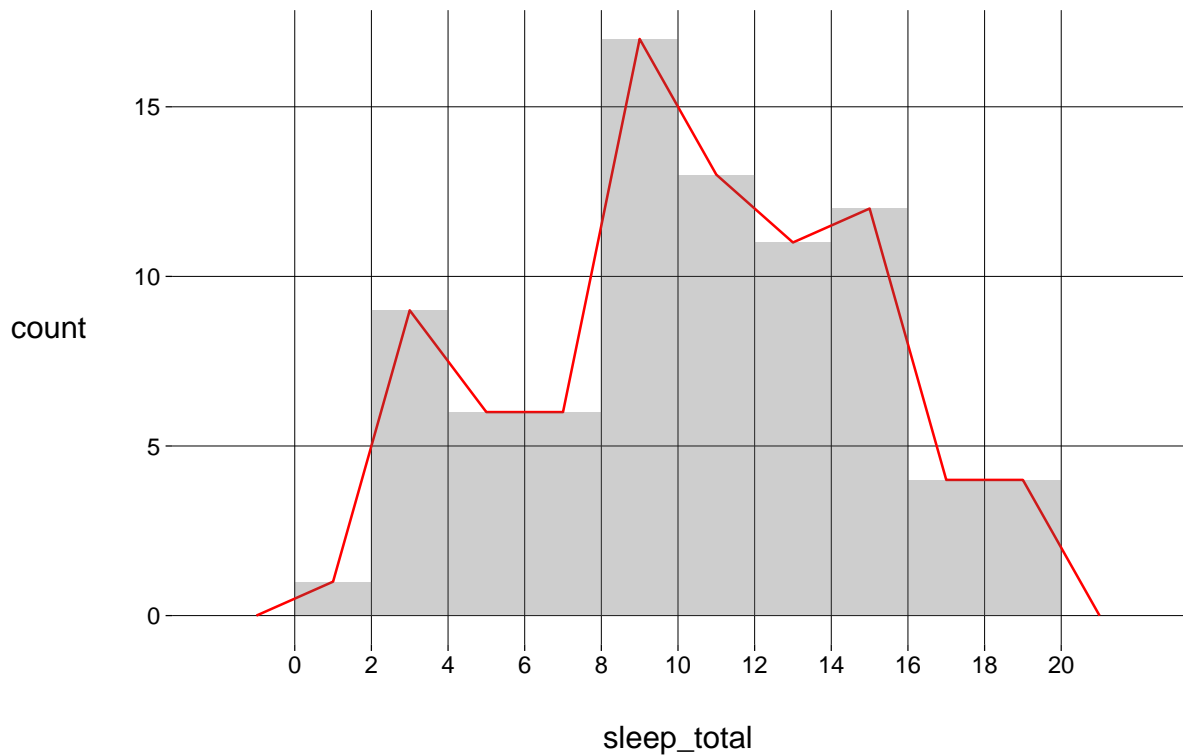
- Em vez das barras do histograma, podemos desenhar uma linha ligando seus topos.
- O resultado é um **polígono de frequência**.

```
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)
```

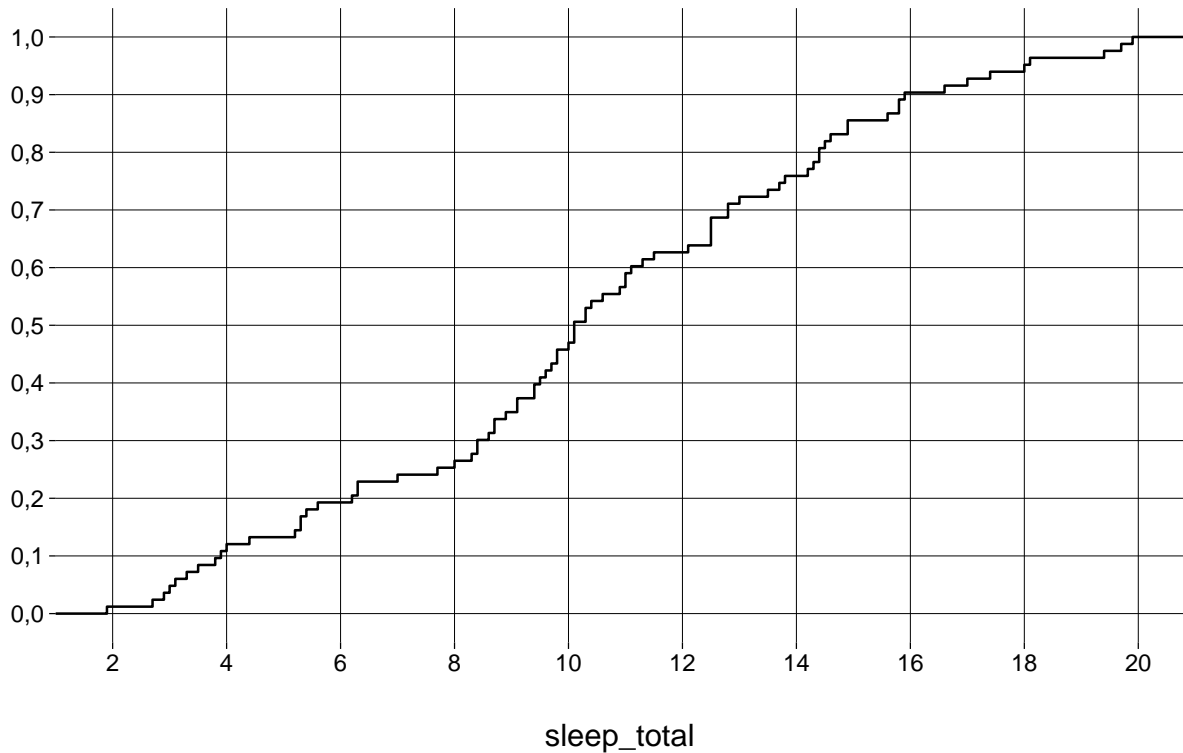


## 4.7

### 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.

## 4.8

### 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 |
##
```

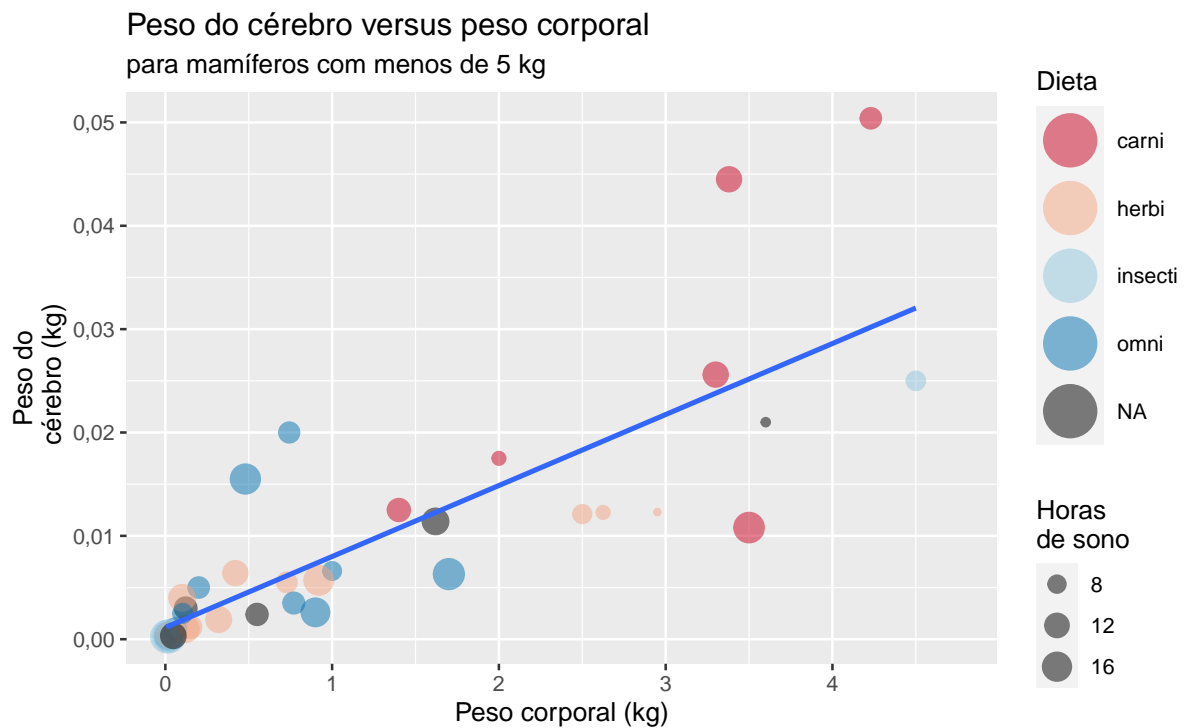
```
##    0 | 9
##    2 | 79013589
##    4 | 0423346
##    6 | 23307
##    8 | 03446779114456788
##   10 | 01113346900135
##   12 | 15555880578
##   14 | 234456996889
##   16 | 604
##   18 | 01479
```

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

## 4.9

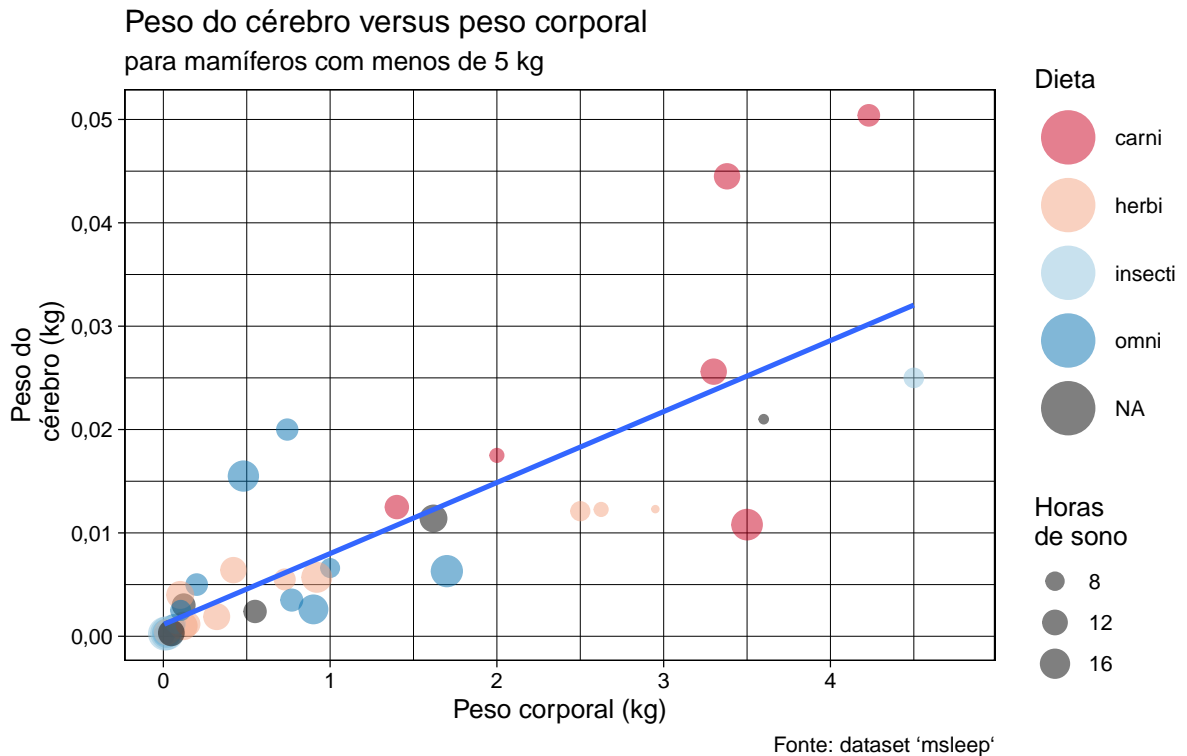
### Personalização do tema

- O `ggplot2` tem um tema *default*, chamado `theme_gray`, que gera o *scatterplot* de um exemplo anterior deste capítulo do seguinte modo:



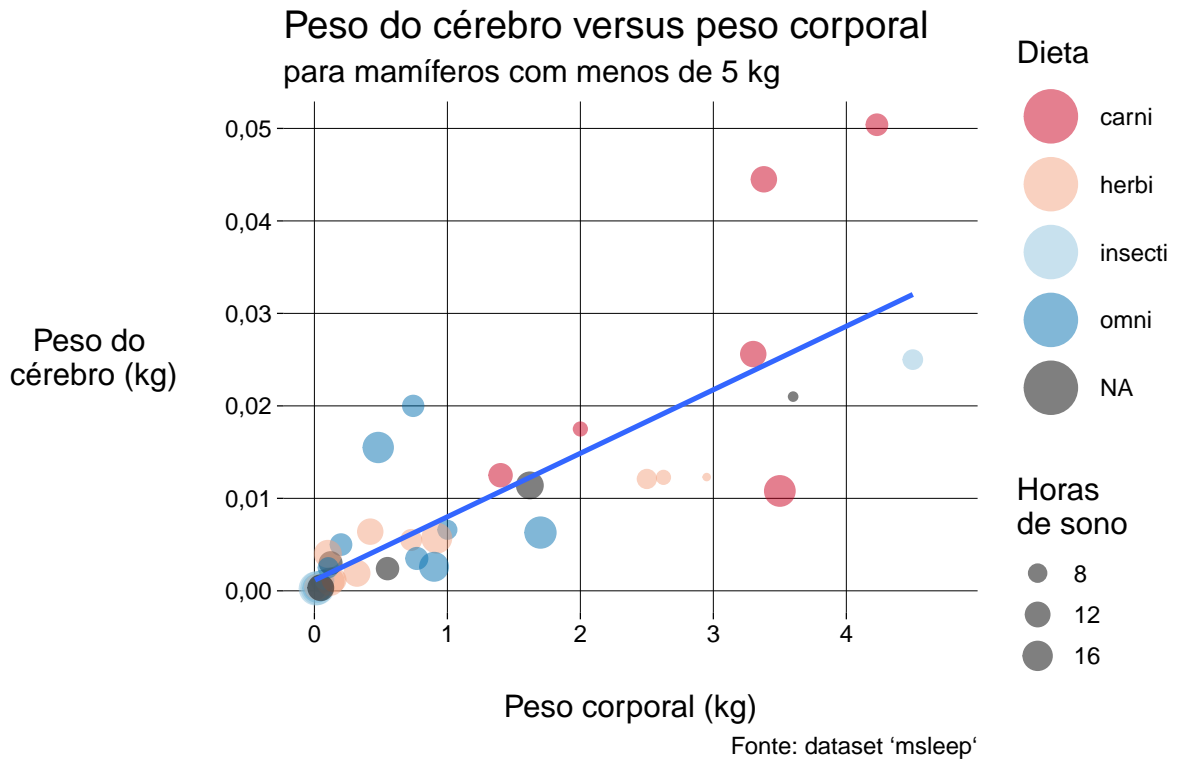
- Para este material, escolhi o tema `theme_linedraw`, que usa linhas pretas sobre fundo

branco:



- Para deixar os gráficos mais leves e facilitar a leitura, fiz as seguintes alterações no tema:
  - Mudei o tamanho do texto dos rótulos.
  - Fiz o rótulo do eixo  $y$  aparecer na horizontal; embora isto ocupe um pouco mais de espaço, evita que o leitor tenha que girar a cabeça para ler o rótulo.
  - Eliminei as linhas dos eixos, para o gráfico ficar mais leve.
  - Eliminei a moldura da área de dados, para o gráfico ficar mais leve.
  - Eliminei a grade secundária, para o gráfico ficar mais leve.
- O resultado é





- Os meus comandos para alterar o tema são

```
# Tamanho do texto depende do formato de saída (html ou pdf):
plot_text_size = ifelse(is_html_output(), 12, 13)

# Tema mais leve:
theme_set(
  theme_linedraw() +
  theme(
    # Tamanho do texto
    text = element_text(size = plot_text_size),
    # Eixo y
    axis.title.y.left = element_text(
      # Nunca girar o rótulo do eixo y
      angle = 0,
      # Separar o rótulo do eixo um pouco
      margin = margin(r = 20),
      # Posicionar verticalmente no meio
      vjust = .5
    ),
    # Eixo y secundário (à direita), quando presente
    axis.title.y.right = element_text(
      # Nunca girar o rótulo do eixo y
      angle = 0,
      # Separar o rótulo do eixo um pouco
```

```

margin = margin(l = 20),
# Posicionar verticalmente no meio
vjust = .5
),
# Não colocar marcas no eixo y secundário
axis.ticks.y.right = element_blank(),
# Separar o eixo x do rótulo um pouco mais
axis.title.x.bottom = element_text(
  margin = margin(t = 20)
),
# Eliminar linhas dos eixos
axis.line = element_blank(),
# Eliminar a moldura da área de dados
panel.border = element_blank(),
# Eliminar a grade secundária
panel.grid.minor = element_blank()
)
)

```

## 4.10

### Exercícios



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

#### 4.10.1

##### 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.

#### 4.10.2

##### 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` 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](#)), 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.

#### 4.10.3

##### Igualdade de gênero entre furacões?

[Este artigo](#) tenta achar uma relação entre o gênero do nome de um furacão e a quantidade de vítimas fatais provocadas por ele.

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](#)), 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](#)), 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](#)), 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.

---

### Visualização com ggplot2 (continuação)

---



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

#### 5.1

---

##### Vídeo 1

<https://youtu.be/TjgLDeIQHIc>

#### 5.2

---

##### *Boxplots*

##### 5.2.1

---

##### 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
```

## 5.2.2

### 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 idéia: 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
```

### 5.2.3

#### 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*.

## 5.2.4

### 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)

unnamed[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
```

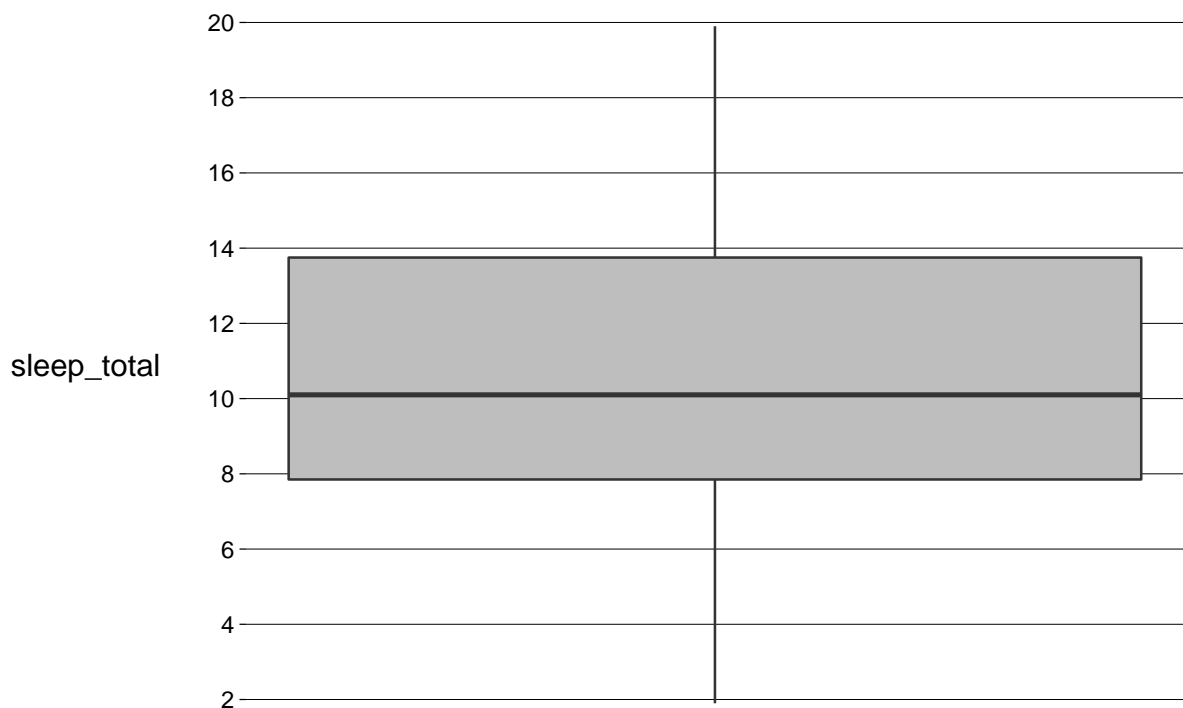
Sim! Estes animais dormem demais em comparação com os outros onívoros.

## 5.2.5

### Gerando boxplots

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

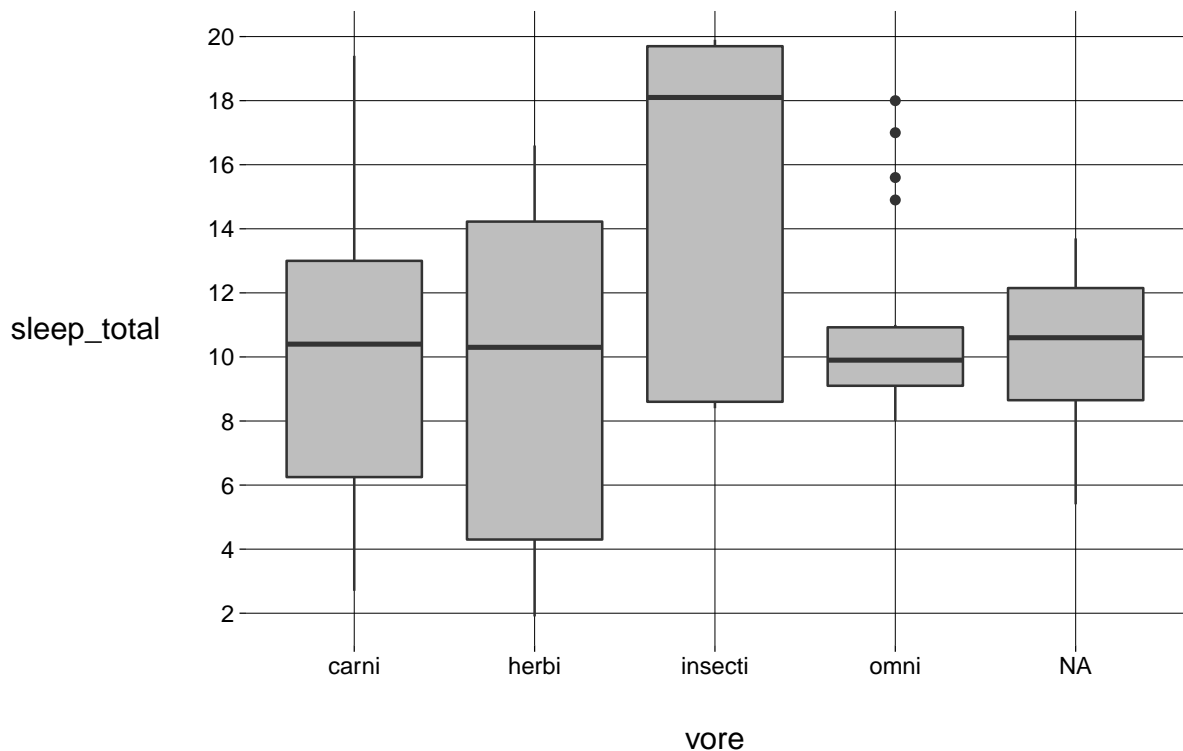
```
sono %>%  
  ggplot(aes(y = sleep_total)) +  
    geom_boxplot(fill = 'gray') +  
    scale_x_continuous(breaks = NULL) +  
    scale_y_continuous(breaks = seq(0, 20, 2))
```



- 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))
```



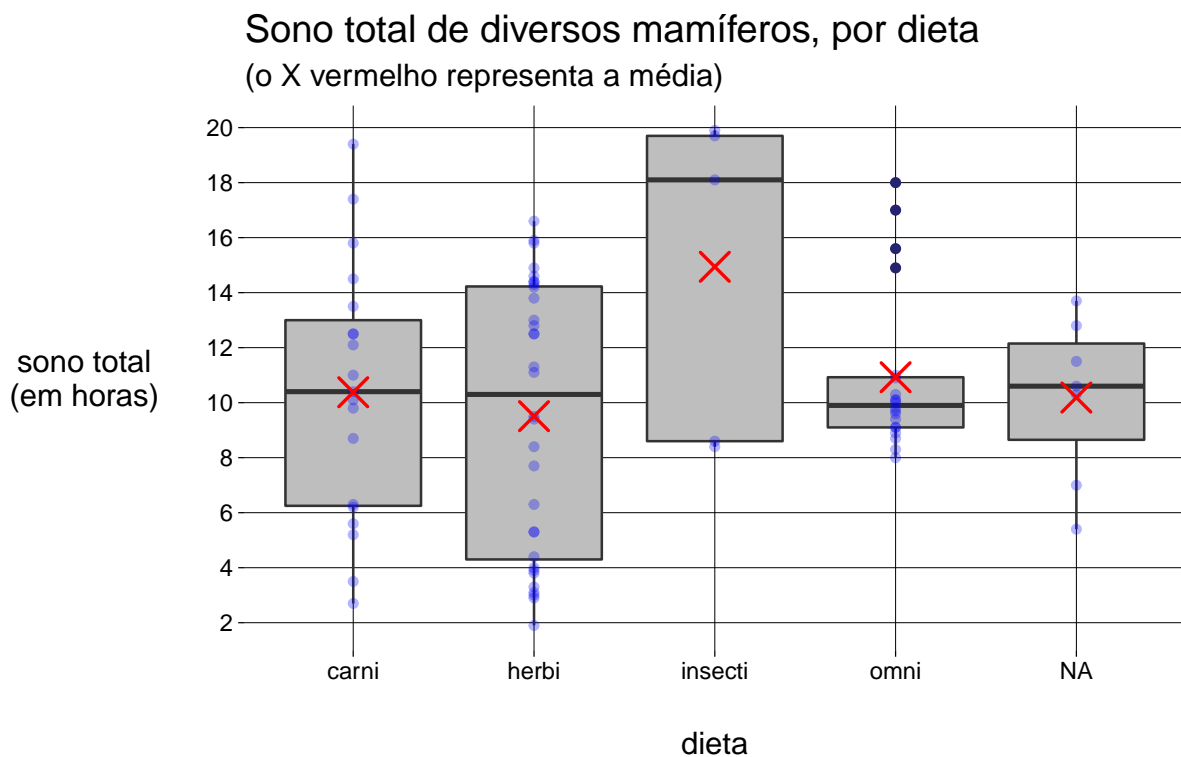
- No *boxplot* de onívoros, **os outliers aparecem como pontos isolados**, acima da caixa, além dos alcances do bigode superior (aliás, onde está bigode superior?).
- *Boxplots* lado a lado são úteis para compararmos grupos diferentes de dados.
- Veja como, com exceção dos insetívoros, as medianas dos 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.

## 5.3

---

### Vídeo 2

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

## 5.4

---

### Gráficos de barras e de colunas

#### 5.4.1

---

##### Conjunto de dados

- O R tem um *array* de 3 dimensões com dados sobre as cores dos cabelos e dos olhos de 592 alunos e alunas de uma universidade americana em 1974.
- Se pedirmos para o R exibir os dados, veremos **duas matrizes**, uma para cada sexo:

```
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
```

- Vamos transformar este *array* em um *data frame*.
- O *array* contém apenas os totais de cada classe. Vamos usar a função `uncount` para gerar uma linha para cada aluno:

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



```
df_orig
```

```
## # A tibble: 592 x 3
##   Hair Eye Sex
##   <fct> <fct> <fct>
## 1 Black Brown Male
## 2 Black Brown Male
## 3 Black Brown Male
## 4 Black Brown Male
## 5 Black Brown Male
## 6 Black Brown Male
## # ... with 586 more rows
```

- O `ggplot2` e os outros pacotes do `tidyverse` foram projetados para trabalhar com *data frames* neste formato, com uma observação (um indivíduo, um elemento) por linha. É o chamado *formato tidy*.
- Usando vetores com elementos nomeados, podemos traduzir o conteúdo do *data frame* para português:

```
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]
  )
```

- Um sumário:

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

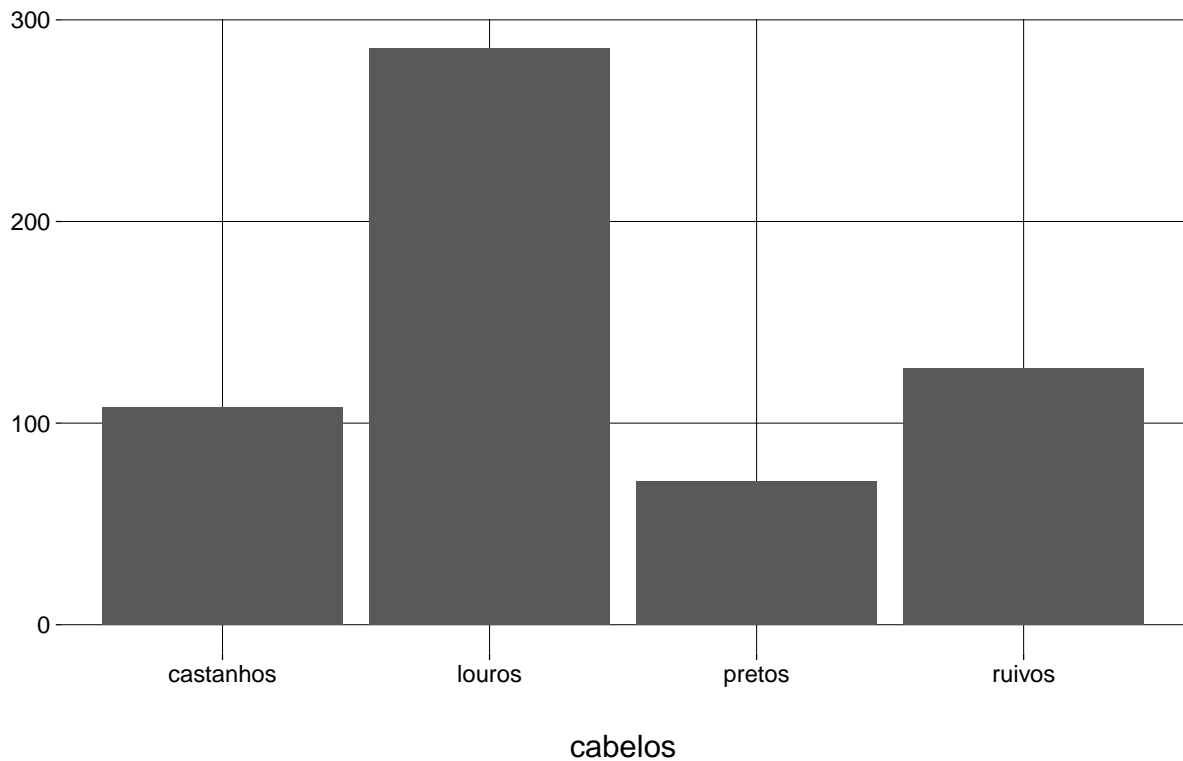
Variável	Estatísticas / Valores	Freqs (% de Válidos)	Faltante
cabelos [character]	1. castanhos	108 (18,2%)	0
	2. louros	286 (48,3%)	(0,0%)
	3. pretos	71 (12,0%)	
	4. ruivos	127 (21,5%)	
olhos [character]	1. avelã	93 (15,7%)	0
	2. azuis	215 (36,3%)	(0,0%)
	3. castanhos	220 (37,2%)	
	4. verdes	64 (10,8%)	
sexo [character]	1. homem	279 (47,1%)	0
	2. mulher	313 (52,9%)	(0,0%)

## 5.4.2

### Gerando gráficos de barras

- Um **gráfico de barras** contém uma barra para cada valor de uma **variável categórica**.
- Usamos `geom_bar` para gerar um gráfico de barras de cores de cabelo:

```
df %>%
  ggplot(aes(x = cabelos)) +
    geom_bar() +
    labs(y = NULL)
```



#### Gráfico de barras × histograma:

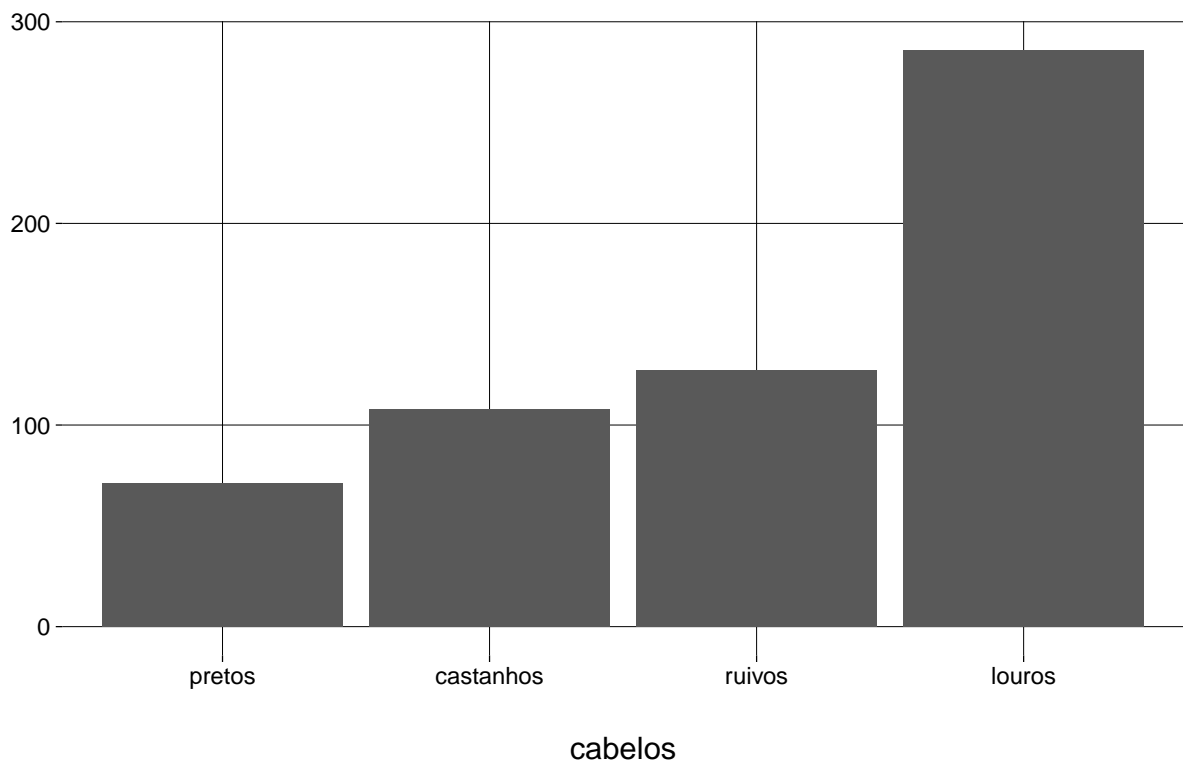
- Os dois tipos de gráficos mostram a frequência (quantidade de elementos) no eixo vertical.
- No gráfico de barras:
  - \* A variável é categórica (nominal).
  - \* Cada barra corresponde a um valor da variável.
  - \* As barras não se tocam, enfatizando o fato de que a variável é categórica.
- No histograma (veja o exemplo):
  - \* A variável é quantitativa (intervalar ou racional).
  - \* Cada barra corresponde a uma classe de valores da variável.
  - \* As barras se tocam, para enfatizar que as classes são contíguas.

- Um gráfico de barras é mais legível quando as barras são mostradas em ordem crescente ou decrescente.
- Embora os valores da variável `cabelos` sejam *strings*, podemos aplicar a eles funções que manipulam fatores.
- A função `fct_infreq`, do pacote `forcats`, ordena os valores em ordem decrescente

de frequência.

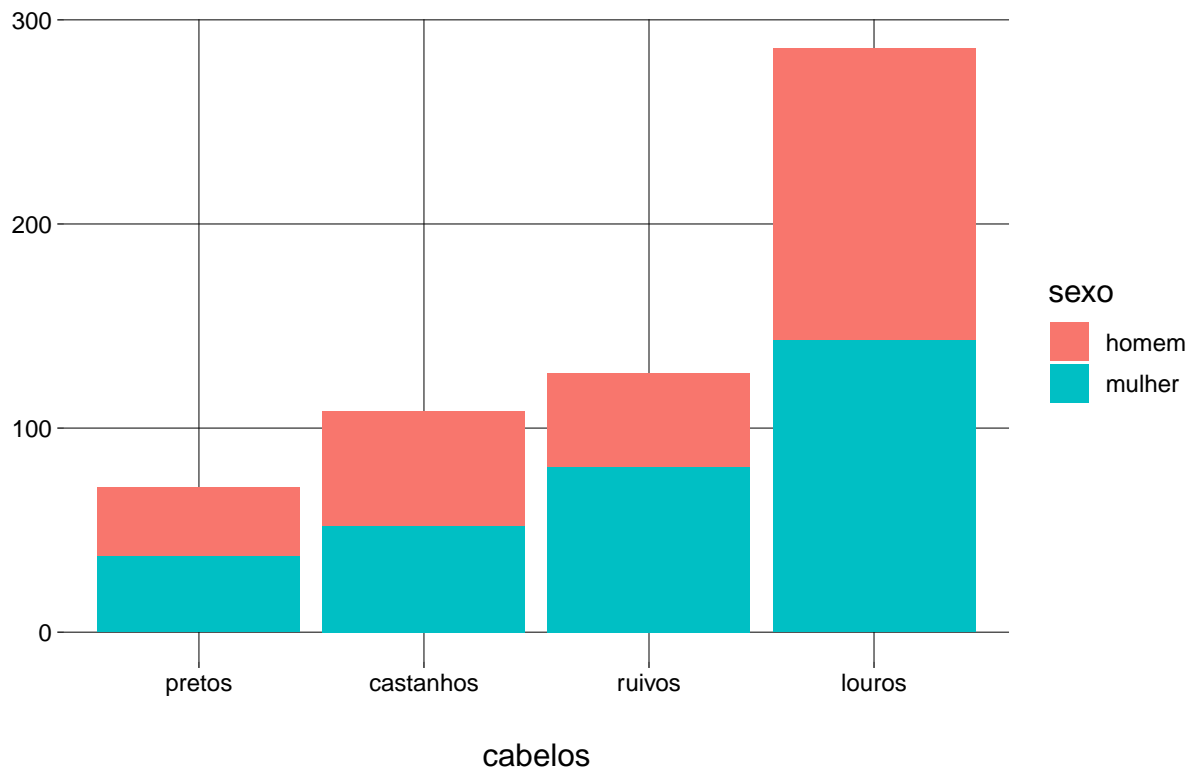
- A função `fct_rev`, também do pacote `forcats`, **inverte a ordenação**.

```
df %>%  
  ggplot(aes(x = fct_rev(fct_infreq(cabelos)))) +  
    geom_bar() +  
    labs(  
      x = 'cabelos',  
      y = NULL  
    )
```



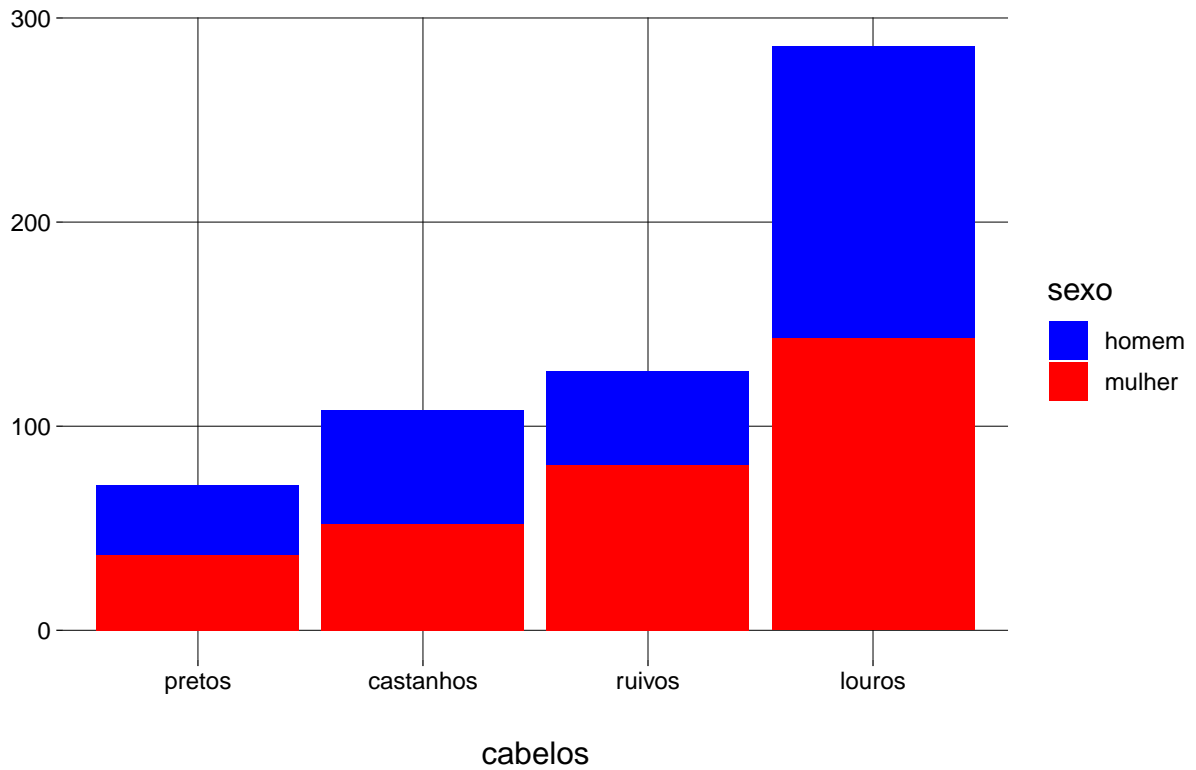
- A posição  $x$  e a altura de cada barra são estéticas: **a posição  $x$  representa a cor dos cabelos**, e **a altura representa a frequência daquela cor**.
- Vamos acrescentar mais uma estética: **a cor de preenchimento vai representar o sexo**.

```
df %>%  
  ggplot(aes(x = fct_rev(fct_infreq(cabelos)), fill = sexo)) +  
    geom_bar() +  
    labs(  
      x = 'cabelos',  
      y = NULL  
    )
```



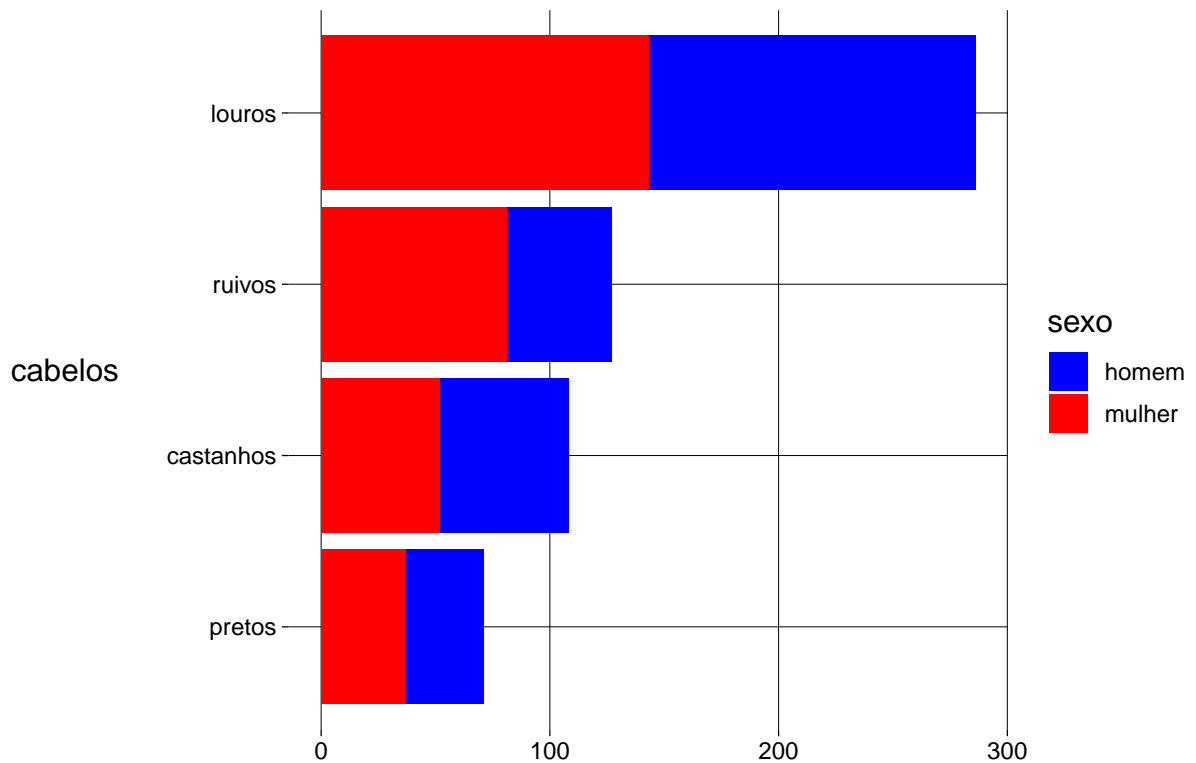
- Se a cor dos homens incomoda você, altere a escala que especifica o preenchimento (scale\_fill\_discrete):

```
df %>%  
  ggplot(aes(x = fct_rev(fct_infreq(cabelos)), fill = sexo)) +  
    geom_bar() +  
    scale_fill_discrete(type = c('blue', 'red')) +  
    labs(  
      x = 'cabelos',  
      y = NULL  
    )
```



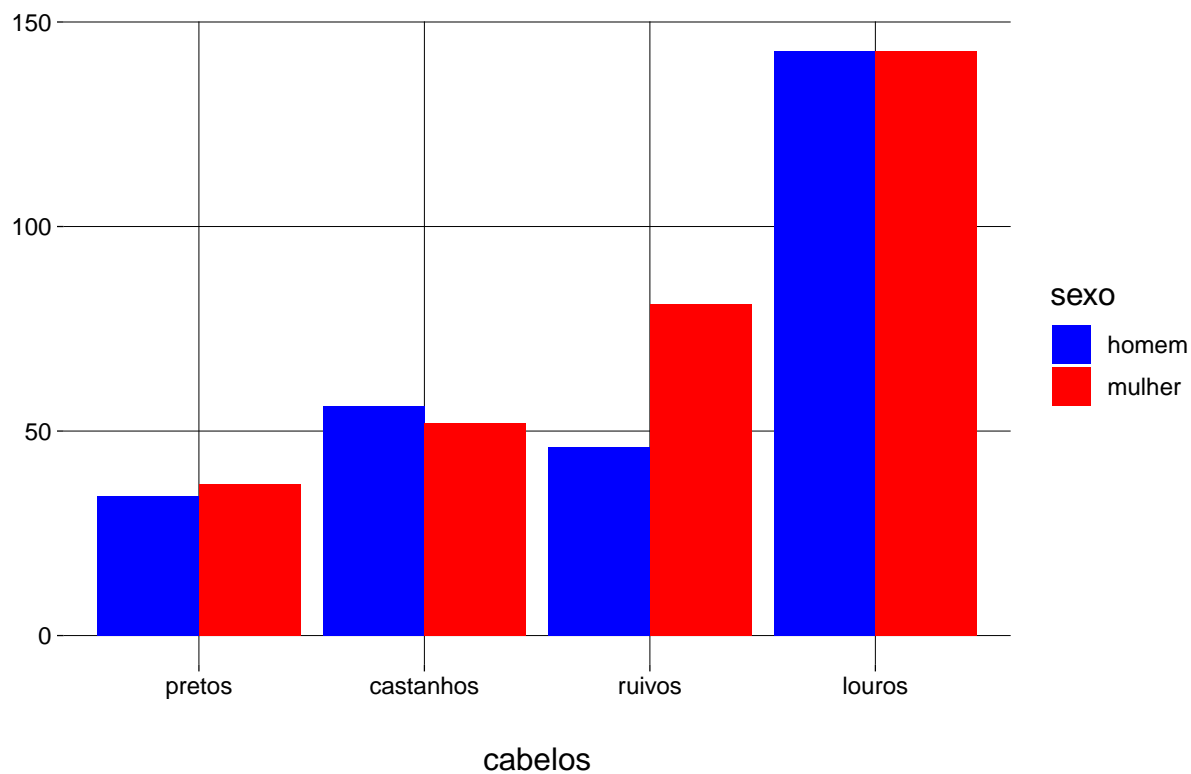
- Podemos fazer um gráfico de barras horizontais com `coord_flip`. Isto geralmente é útil quando os rótulos das barras são longos:

```
df %>%  
  ggplot(aes(x = fct_rev(fct_infreq(cabelos)), fill = sexo)) +  
    geom_bar() +  
    scale_fill_discrete(type = c('blue', 'red')) +  
    labs(  
      x = 'cabelos',  
      y = NULL  
    ) +  
    coord_flip()
```



- Você consegue dizer se há mais homens ou mulheres com cabelos pretos? E castanhos? E ruivos?
- Se posicionarmos as barras lado a lado, fica mais fácil responder.
- Usamos o argumento `position = 'dodge'` de `geom_bar`. “*Dodge*” significa “esquivar-se”, em inglês.

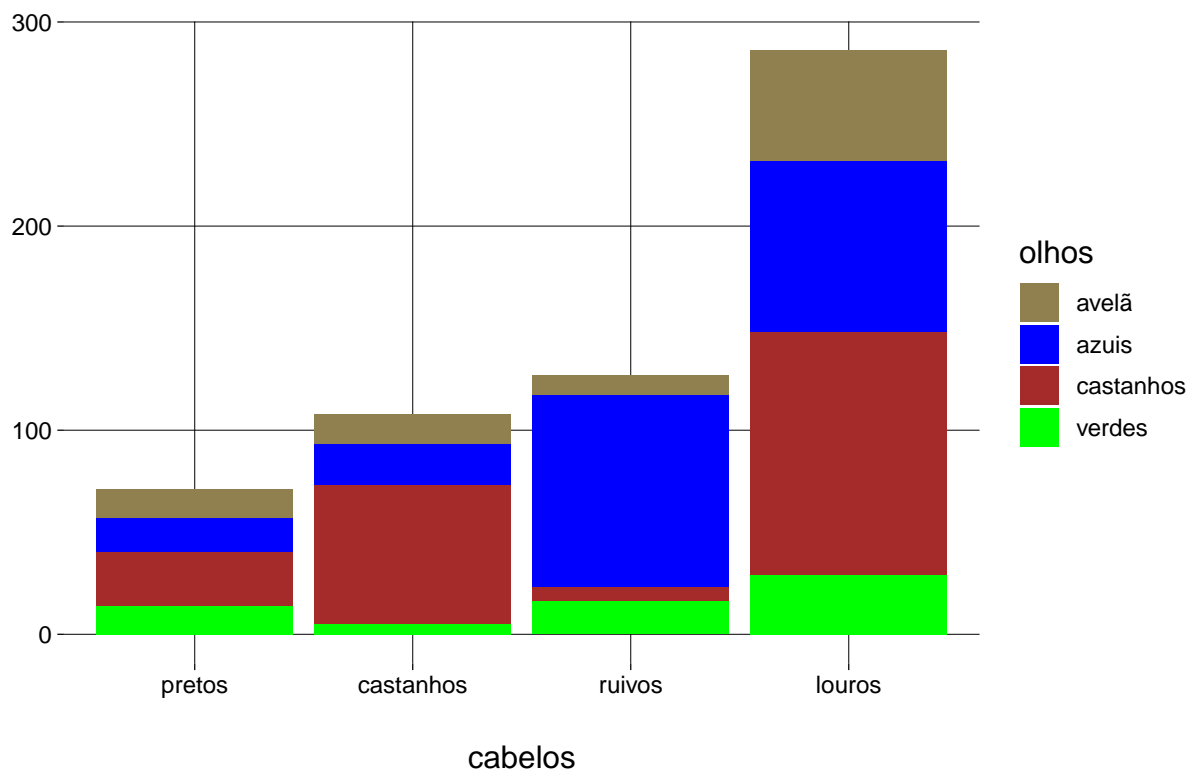
```
df %>%
  ggplot(aes(x = fct_rev(fct_infreq(cabelos)), fill = sexo)) +
  geom_bar(position = 'dodge') +
  labs(
    x = 'cabelos',
    y = NULL
  ) +
  scale_fill_discrete(type = c('blue', 'red'))
```



- Agora vamos examinar a relação entre as cores dos olhos e as cores dos cabelos:

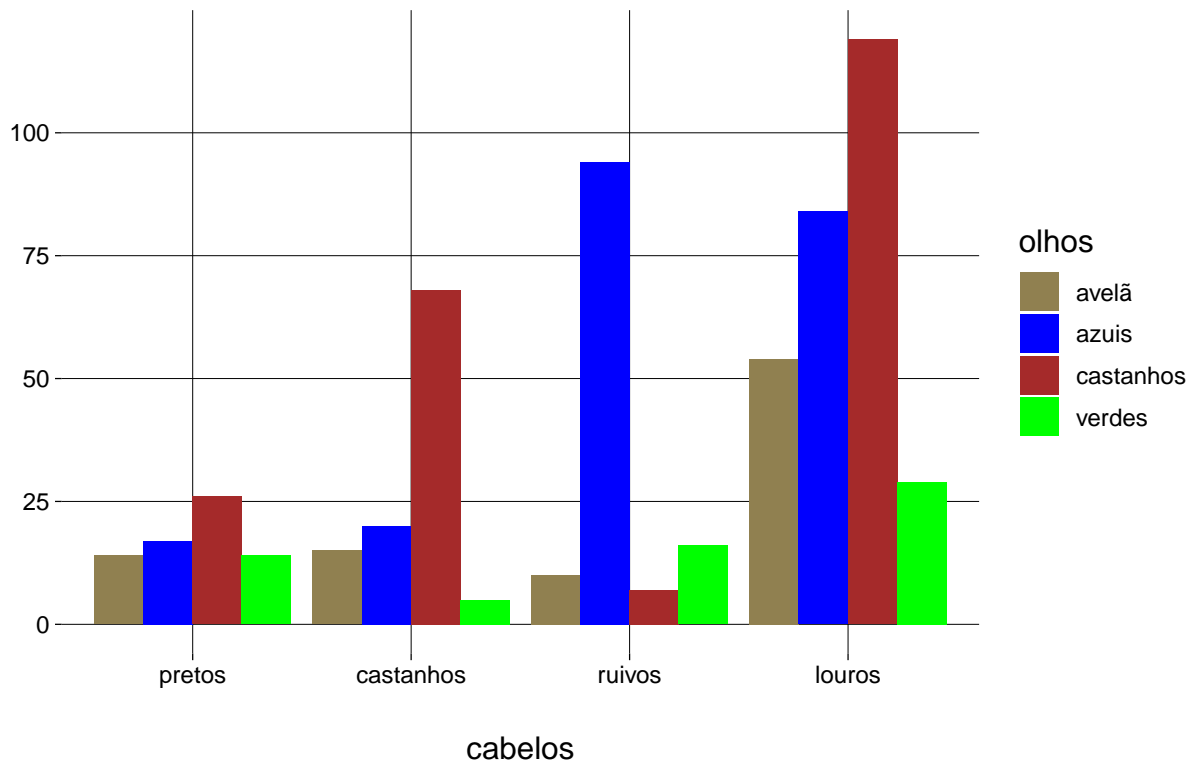
```
df %>%  
  ggplot(aes(x = fct_rev(fct_infreq(cabelos)), fill = olhos)) +  
    geom_bar() +  
    scale_fill_discrete(  
      type = c('#908050', 'blue', 'brown', 'green')  
    ) +  
    labs(  
      x = 'cabelos',  
      y = NULL  
    )
```





- Ou, com barras lado a lado:

```
df %>%
  ggplot(aes(x = fct_rev(fct_infreq(cabelos)), fill = olhos)) +
  geom_bar(position = 'dodge') +
  scale_fill_discrete(
    type = c('#908050', 'blue', 'brown', 'green')
  ) +
  labs(
    x = 'cabelos',
    y = NULL
  )
```



- Observações e perguntas:

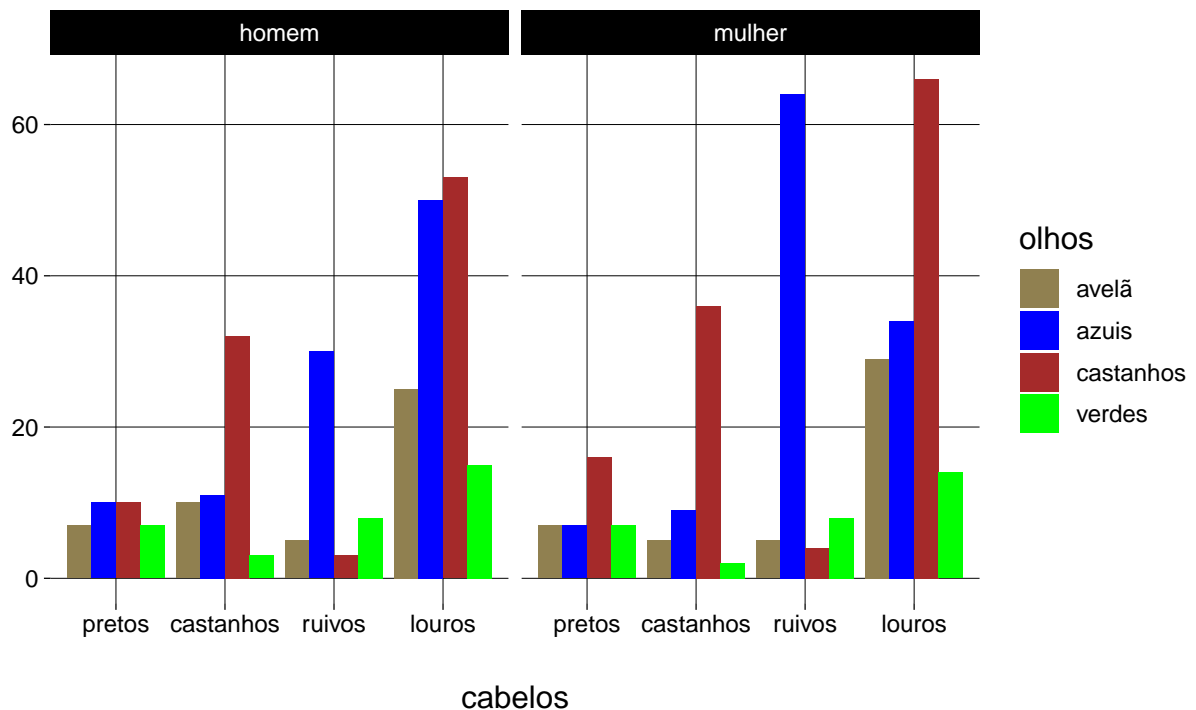
1. Há mais pessoas loiras de olhos castanhos do que loiras de olhos azuis. O esperado não seria mais pessoas loiras de olhos azuis? Pessoas loiras de olhos castanhos pintaram os cabelos?
2. Há muito mais ruivos de olhos azuis do que ruivos de olhos verdes. Não deveria ser o contrário? Também são pessoas que pintaram os cabelos de ruivo? Ou houve erro no registro das cores dos olhos?

- Para incluir o sexo, podemos **facetar** o gráfico. Usando `facet_wrap`<sup>1</sup>, geramos dois subgráficos lado a lado:

```
df %>%
  ggplot(aes(x = fct_rev(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'
  )
```

<sup>1</sup>O nome da variável segundo a qual facetar deve aparecer depois de um ~.

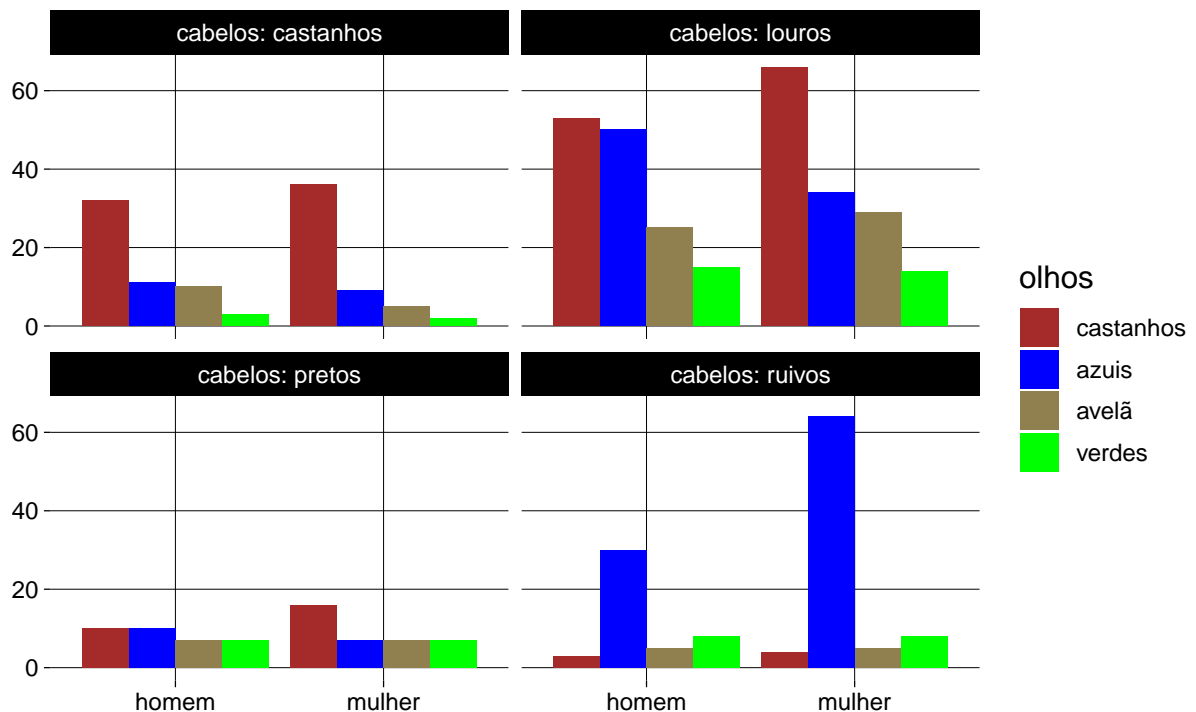
## Cores de cabelos e olhos por sexo



- Se a quantidade grande de pessoas louras de olhos castanhos (em comparação com pessoas louras de olhos azuis) for por causa da pintura de cabelos, então o gráfico acima mostra que as mulheres pintam os cabelos de loiro com mais frequência do que os homens.
- Quando facetamos por cor de cabelos, também podemos observar as mesmas diferenças entre homens e mulheres:

```
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



### 5.4.3

#### Data frame já contendo os totais

- Você percebeu que `geom_bar` analisa o *data frame* e calcula as frequências necessárias para construir o gráfico.
- Em algumas situações, o *data frame* já contém as frequências (em vez de conter uma linha por indivíduo).
- Vamos usar `count` para criar um *data frame* assim:

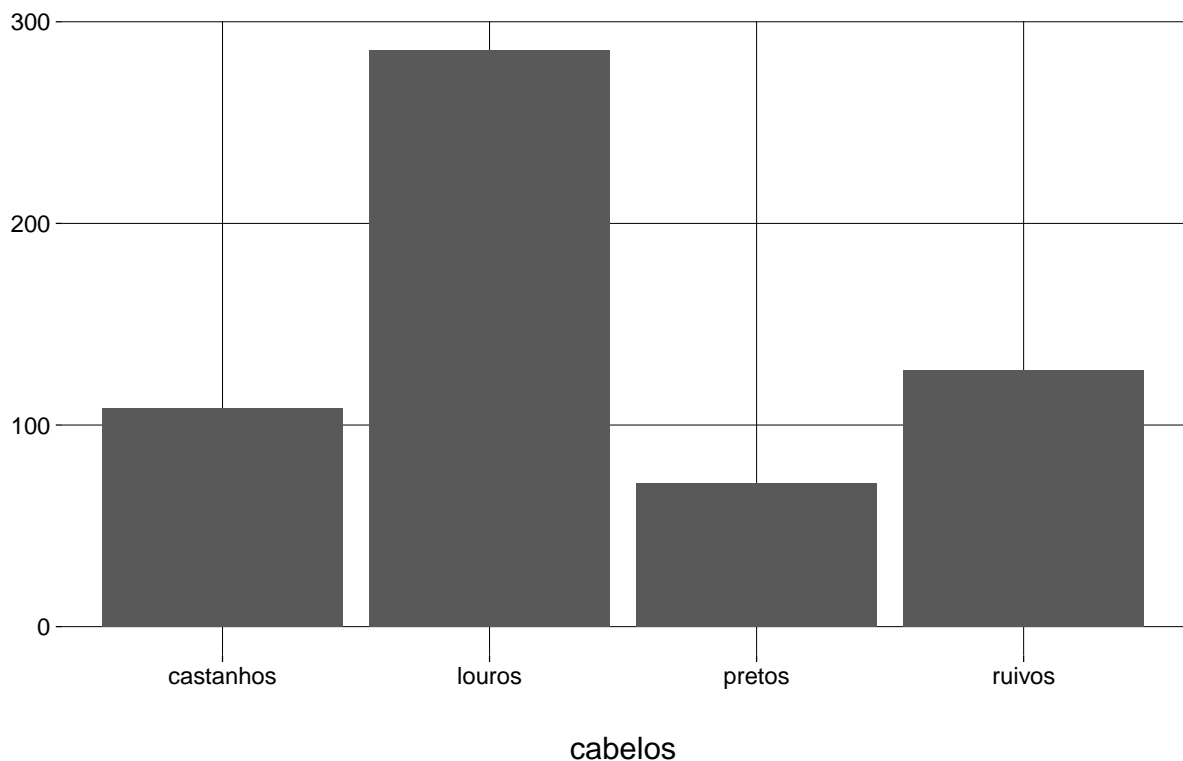
```
df_tot <- df %>%
  count(sexo, cabelos, olhos)
```

```
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
```

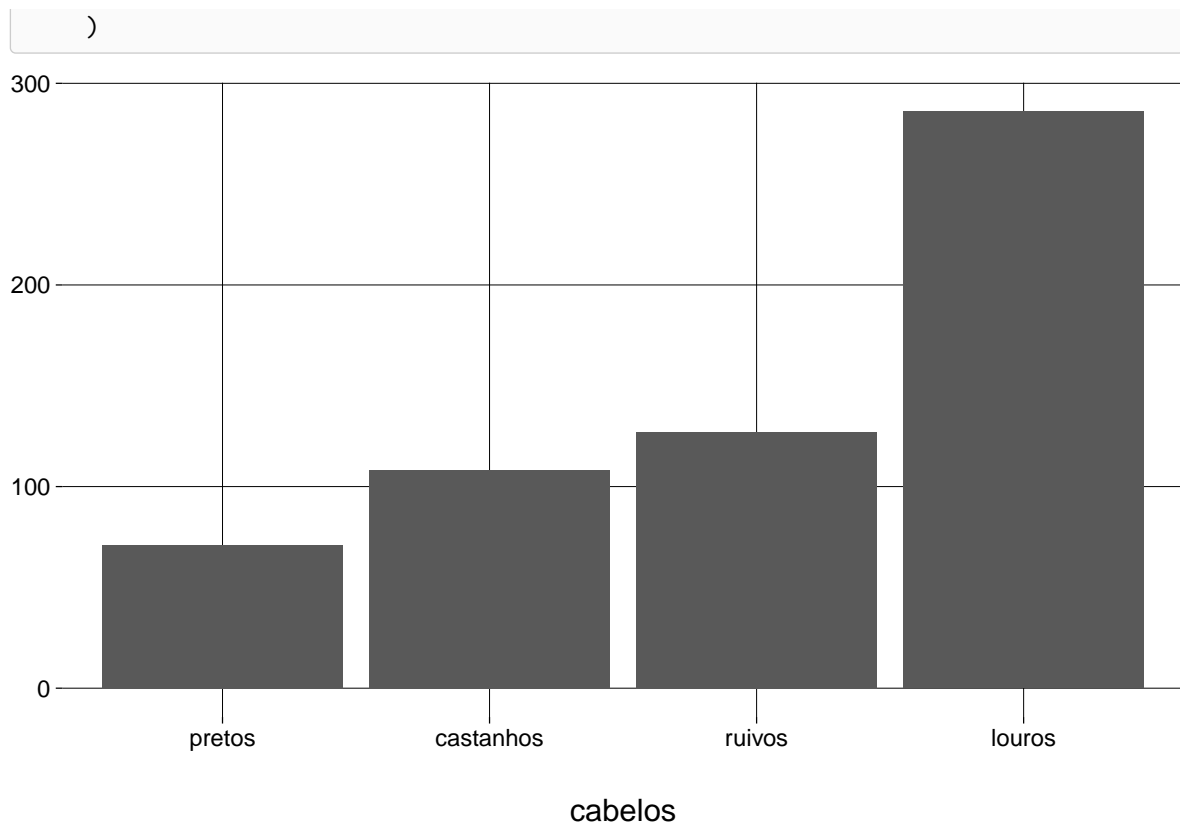
- Para 4 cores de cabelo, 4 cores de olhos, e 2 sexos, são 32 combinações possíveis.
- Com este *data frame*, podemos gerar todos os gráficos anteriores usando `geom_col` no lugar de `geom_bar`. Por exemplo:

```
df_tot %>%
  ggplot(aes(x = cabelos, y = n)) +
    geom_col() +
    labs(
      y = NULL
    )
```



- Com `geom_col`, precisamos passar a estética `y` (no nosso exemplo, a variável `n`, que contém as frequências).
- Para ordenar as barras, usamos a função `fct_reorder`, que ordena os níveis de um fator (`cabelos`) de acordo com o resultado de uma função (`sum`) aplicada sobre os valores de outra variável (`n`):

```
df_tot %>%
  ggplot(aes(x = fct_reorder(cabelos, n, sum), y = n)) +
    geom_col() +
    labs(
      x = 'cabelos',
      y = NULL
    )
```



## 5.5

### Gráficos de linha e séries temporais

#### 5.5.1

##### Conjunto de dados

- O R tem uma matriz com as quantidades de telefones em várias regiões do mundo ao longo de vários anos:

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

- Os números representam milhares.
- Os números dos anos são os nomes das linhas da matriz.

- Vamos transformar esta matriz em uma *tibble*:

```
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
```

- Esta *tibble* não está no formato *tidy*. Queremos que cada linha corresponda a uma observação, contendo
  - Ano,
  - Região,
  - Quantidade de telefones.
- Usamos a função `pivot_longer` para mudar o formato da *tibble*:

```
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
```

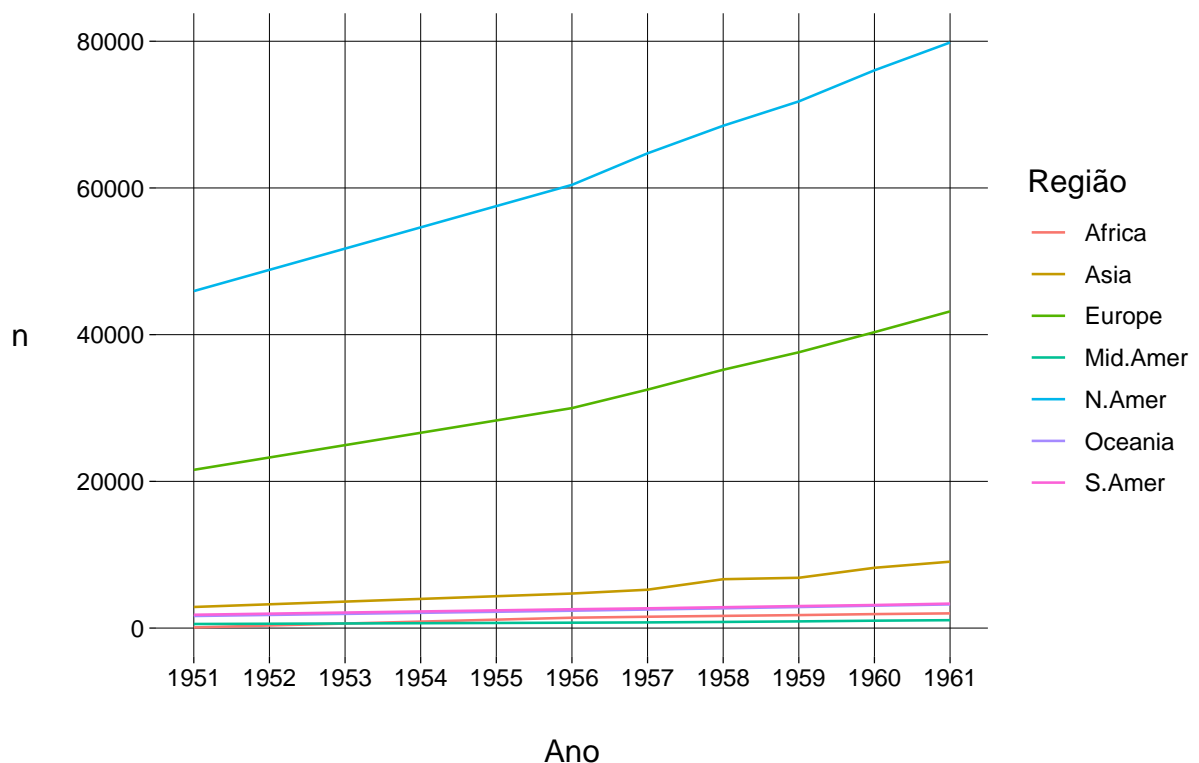
- Confira: antes, tínhamos 7 anos, com 7 quantidades por ano, uma quantidade por região. Eram 49 quantidades. Agora temos uma *tibble* de 49 linhas.

## 5.5.2

### Gerando gráficos de linha

- A geometria `geom_line` gera gráficos de linha. Perceba como geramos uma linha por região:

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



- Embora a legenda associe uma cor a cada região, a leitura seria mais fácil se a ordem das regiões na legenda coincidissem com a posição das linhas na borda direita da grade:

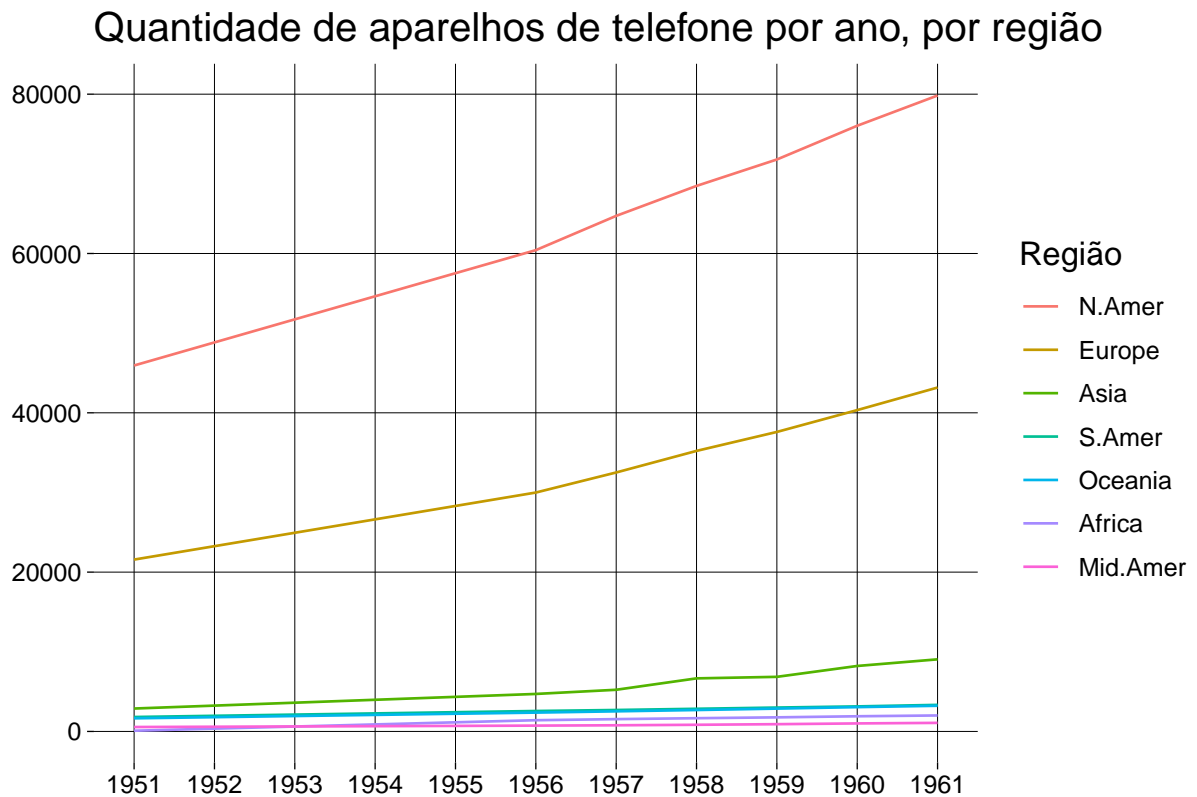
```
fones_long %>%
  ggplot(
    aes(
      x = Ano,
      y = n,
```



```

    color = fct_rev(fct_reorder(Região, n, max))
  )
) +
  geom_line() +
  scale_x_continuous(breaks = 1951:1961) +
  labs(
    color = 'Região',
    y = '',
    x = NULL,
    title = 'Quantidade de aparelhos de telefone por ano, por região'
  )

```



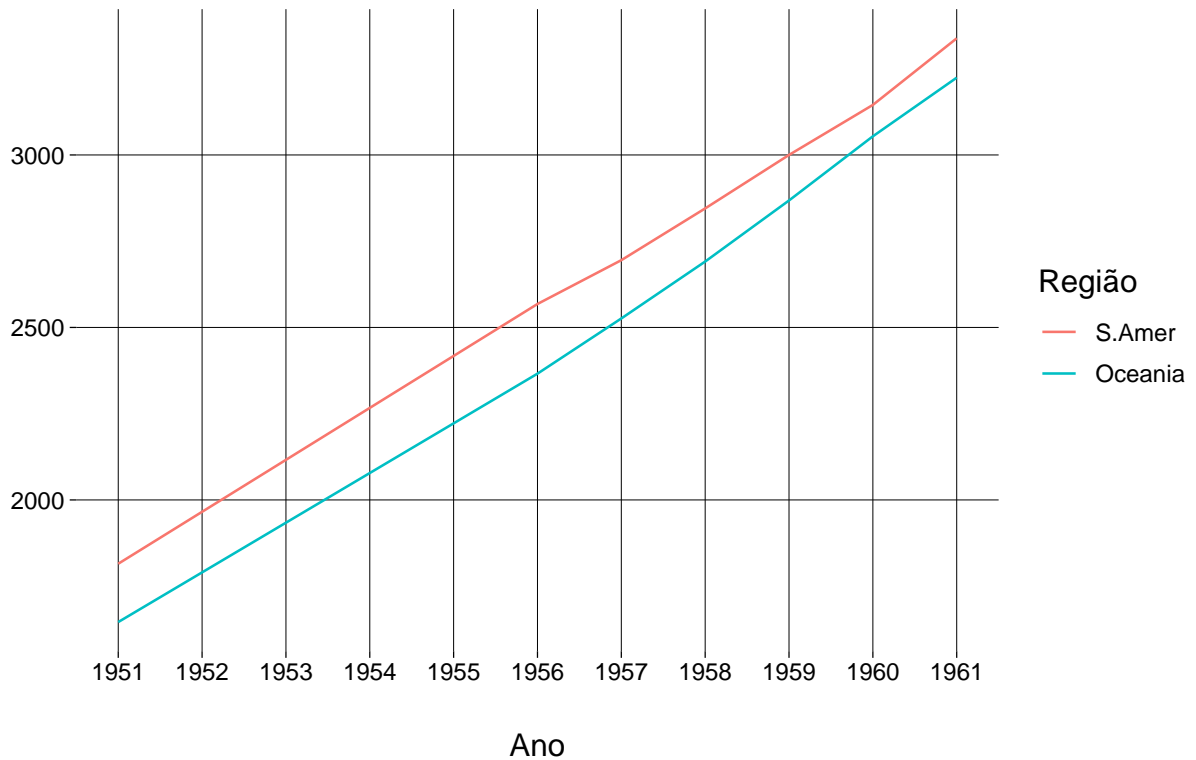
- Parece que está faltando uma linha, mas o que acontece é que as quantidades da América do Sul e da Oceania são bem parecidas:

```

fones_long %>%
  filter(Região %in% c('S.Amer', 'Oceania')) %>%
  ggplot(
    aes(
      x = Ano,
      y = n,
      color = fct_rev(fct_reorder(Região, n, max))
    )
  ) +

```

```
geom_line() +
  scale_x_continuous(breaks = 1951:1961) +
  labs(y = NULL, color = 'Região')
```



- Estamos tratando estes dados como simples números, mas, na verdade, **este conjunto de dados é uma série temporal (*time series*)**.
- R tem todo um conjunto de funções para tratar séries temporais, calcular tendências, achar padrões cíclicos, fazer estimativas, e gerar gráficos específicos, entre outras coisas.
- Mas não vamos falar mais sobre séries temporais aqui.
- O **pacote tsibble** oferece maneiras de trabalhar com séries temporais de maneira *tidy*. Você pode ler a documentação do pacote entrando

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

## 5.6

---

### Exercícios

#### 5.6.1

---

##### O bigode dos onívoros

- Examine o *data frame* `sono` para descobrir o que houve com o bigode superior do *boxplot* dos onívoros [neste gráfico](#).

#### 5.6.2

---

##### Usando `geom_col`

- Use `geom_col` para reproduzir, a partir do *data frame* `df_tot`, todos os gráficos que foram gerados com `geom_bar` na seção [Gerando gráficos de barras](#).

#### 5.6.3

---

## 5.7

---

### Referências sobre visualização e R



Busque mais informações sobre os pacotes `tidyverse` e `ggplot2` [nas referências recomendadas](#).

---

### Medidas

---

#### 6.1

---

#### Vídeo

<https://youtu.be/C96MOP4YlaY>

#### 6.2

---

### Medidas de centralidade

#### 6.2.1

---

#### Média

- A **média de uma população** é escrita como  $\mu$ , e é definida como

$$\mu = \frac{\sum_{i=1}^N x_i}{N}$$

- $\sum_{i=1}^N x_i$  é a soma de todos os dados da população.
- $N$  é a quantidade de elementos na população.

- A **média de uma amostra** é escrita como  $\bar{x}$ , e é definida como:

$$\bar{x} = \frac{\sum_{i=1}^n x_i}{n}$$

- $\sum_{i=1}^n x_i$  é a soma de todos os dados da amostra.
- $n$  é a quantidade de elementos na amostra.
- O cálculo é essencialmente o mesmo. Só mudam os símbolos:  $N$  versus  $n$ , e  $\mu$  versus  $\bar{x}$ .

### 6.2.1.1

#### Exemplo

- Idades dos alunos de uma turma:

```
idades <- c(
  20, 20, 20, 20, 20, 20, 21, 21, 21, 21,
  22, 22, 22, 23, 23, 23, 23, 24, 24,
  65
)
```

- Média **com** o velhinho de 65 anos:

```
mean(idades)
```

```
## [1] 23,75
```

- Média **sem** o velhinho:

```
mean(idades[-length(idades)])
```

```
## [1] 21,57895
```

### 6.2.2

#### Mediana

- Já aprendemos sobre a mediana na [seção sobre boxplots](#).
- A ideia é que, depois de ordenar os dados, 50% dos dados estarão à esquerda da mediana, e 50% à direita.
- A mediana não é tão sensível a *outliers* quanto à média.

### 6.2.2.1

---

#### Exemplo

- Mediana **com** o velhinho:

```
median(idades)
```

```
## [1] 21,5
```

- Mediana **sem** o velhinho:

```
median(idades[-length(idades)])
```

```
## [1] 21
```

### 6.2.3

---

#### Moda

- A **moda** é o **valor mais frequente** do conjunto de dados.
- Pode haver mais de uma moda.
- **Não existe uma função para a moda em R base.** Por quê?
- Por incrível que pareça, **é complicado definir a moda de forma a conseguir resultados interessantes.**
- **Por exemplo**, vamos definir um conjunto de 1000 valores numéricos distribuídos normalmente<sup>1</sup>, com média igual a 5 e desvio-padrão<sup>2</sup> igual a 2:

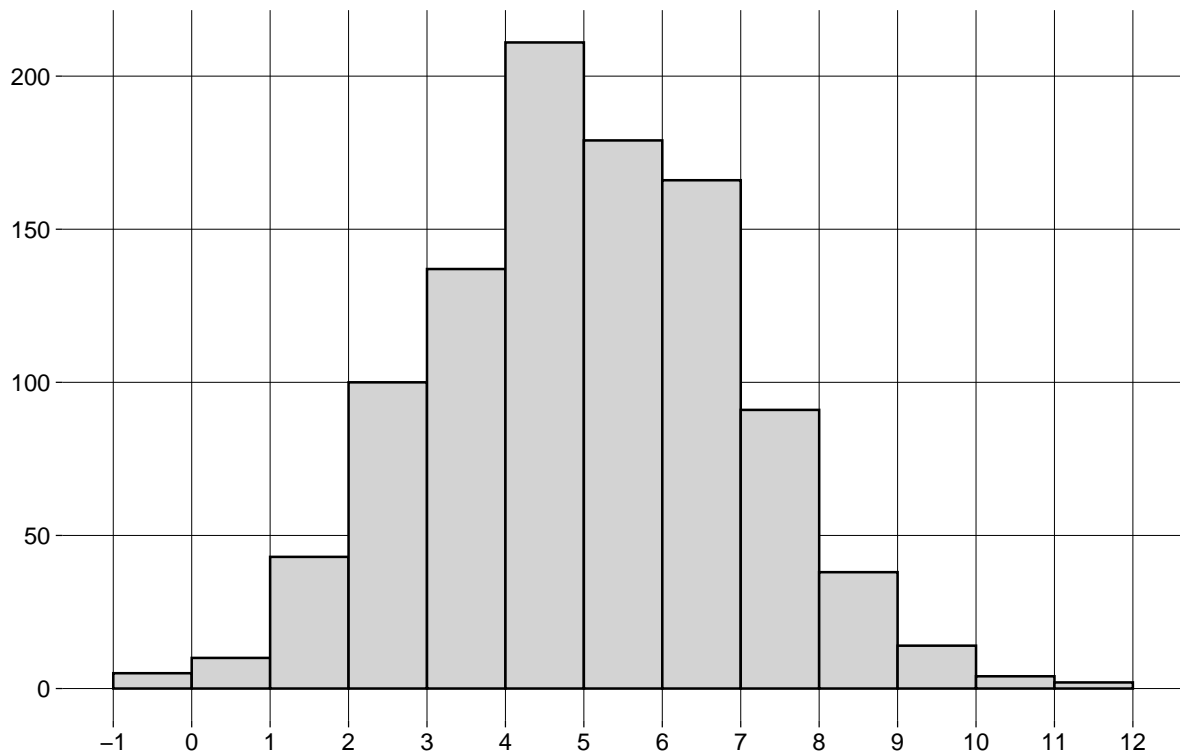
```
normal <- rnorm(1000, mean = 5, sd = 2)
```

- O histograma dos nossos dados é

---

<sup>1</sup>Mais sobre a distribuição normal no capítulo ???.

<sup>2</sup>Mais sobre o desvio-padrão daqui a pouco.



- Vamos calcular a moda com a função `mfv` (*most frequent value*), do pacote `modeest`:

```
# Pacote com funções para calcular modas
library(modeest)
```

```
## Registered S3 method overwritten by 'rmutil':
##   method      from
##   print.response httr
```

```
# Por causa de um bug na função mfv,
# precisamos de números com ponto decimal
# (em vez de vírgula):
options(OutDec = '.')
mfv(normal)
```

```
##      [1] -0.55298295 -0.49879679 -0.38565757 -0.28251884 -0.04510405
##      [6]  0.02831001  0.28307531  0.28962328  0.30535535  0.37548108
##     [11]  0.48832809  0.58475232  0.63463085  0.81457263  0.90675406
##     [16]  1.04119634  1.09909979  1.15767635  1.19952655  1.20882480
##     [21]  1.29357411  1.30653613  1.37552267  1.38683882  1.44625887
##     [26]  1.44939762  1.45830584  1.47233603  1.50062392  1.51705986
##     [31]  1.53151865  1.54930057  1.57741695  1.62520318  1.64460300
##     [36]  1.68445739  1.69481381  1.76193559  1.82081059  1.83176197
##     [41]  1.83888157  1.85700960  1.86151643  1.86903902  1.87782065
##     [46]  1.88974147  1.90018763  1.90027248  1.91034177  1.91256702
##     [51]  1.92182823  1.92611864  1.93077173  1.93590589  1.94149690
```

##	[56]	1.95739918	1.95975225	1.98822697	2.00797867	2.00827251
##	[61]	2.03090944	2.03575976	2.04019901	2.04024904	2.04802186
##	[66]	2.05646645	2.07449795	2.07926124	2.13601562	2.17599116
##	[71]	2.18740562	2.20643651	2.21222505	2.21701330	2.23359350
##	[76]	2.25992683	2.27324721	2.29178308	2.30821314	2.31407224
##	[81]	2.32166561	2.32807181	2.32863138	2.34727308	2.35642099
##	[86]	2.35715160	2.35741029	2.36598991	2.37068222	2.37867935
##	[91]	2.37929394	2.40236281	2.40672962	2.42904864	2.43585918
##	[96]	2.44467044	2.45365890	2.48543109	2.48765811	2.49136419
##	[101]	2.49387704	2.49514738	2.49680630	2.50902470	2.51122509
##	[106]	2.52669435	2.55650809	2.55816004	2.61440244	2.63627293
##	[111]	2.64656611	2.65784358	2.66164217	2.67084104	2.67462881
##	[116]	2.67992270	2.68058719	2.69949309	2.69981226	2.70626096
##	[121]	2.72651779	2.72735345	2.72952943	2.73064274	2.73655107
##	[126]	2.73735457	2.75264397	2.75700313	2.76307423	2.76476342
##	[131]	2.78307482	2.79815166	2.84206494	2.84464167	2.84876293
##	[136]	2.85535067	2.85760130	2.87592706	2.87676990	2.88809310
##	[141]	2.89559424	2.90707270	2.90828208	2.90985508	2.92364922
##	[146]	2.92507393	2.92744569	2.93275880	2.93474571	2.93654348
##	[151]	2.94308224	2.94990237	2.95307648	2.95315509	2.95957251
##	[156]	2.97294203	2.98074651	2.98831012	3.01161807	3.04655947
##	[161]	3.06581454	3.06697074	3.07176378	3.07640552	3.07866700
##	[166]	3.09645316	3.12260566	3.12348707	3.12901305	3.13608295
##	[171]	3.14312376	3.15139616	3.15621470	3.16542288	3.16641412
##	[176]	3.16671445	3.16679641	3.18773157	3.18812586	3.19815400
##	[181]	3.21115563	3.24452530	3.24512854	3.25168775	3.25293357
##	[186]	3.26522675	3.26624534	3.27127793	3.28841478	3.29740979
##	[191]	3.30009576	3.32315279	3.32403405	3.32580517	3.33759587
##	[196]	3.33964567	3.34553398	3.35958051	3.36013876	3.36677745
##	[201]	3.36765012	3.37883086	3.37884220	3.38842582	3.39185144
##	[206]	3.40769815	3.43289517	3.44283617	3.47359409	3.47952602
##	[211]	3.48528413	3.49374667	3.49628301	3.50381883	3.50755418
##	[216]	3.51225591	3.51641167	3.53075944	3.53622665	3.54217437
##	[221]	3.54282179	3.54903673	3.55617683	3.55779837	3.56390169
##	[226]	3.57205173	3.57878237	3.57965778	3.58068177	3.59283167
##	[231]	3.60174417	3.60183685	3.60445172	3.64619769	3.64680119
##	[236]	3.66921337	3.69239117	3.69445380	3.69775084	3.69936237
##	[241]	3.71391516	3.71557606	3.71754932	3.71936933	3.73164209
##	[246]	3.73523087	3.74595914	3.76886339	3.77288150	3.77560549
##	[251]	3.77881741	3.78323824	3.78574444	3.79628891	3.80398692
##	[256]	3.80488853	3.81021175	3.81102173	3.82400100	3.82818833
##	[261]	3.82945493	3.83208231	3.83914185	3.84602360	3.85927332
##	[266]	3.86654717	3.86745599	3.87257744	3.87763251	3.88059924
##	[271]	3.88145121	3.88609031	3.88732383	3.89130100	3.89801146
##	[276]	3.90363895	3.91558265	3.91621565	3.91959061	3.92124877
##	[281]	3.92289530	3.93498932	3.93536316	3.93588923	3.93667460
##	[286]	3.94254979	3.95316257	3.97053331	3.97785447	3.98231287
##	[291]	3.98264833	3.98363309	3.98499650	3.98580738	3.98811757



##	[296]	4.00356282	4.00810209	4.01058327	4.01284349	4.01337990
##	[301]	4.02438428	4.02699334	4.03325417	4.03469725	4.03764576
##	[306]	4.04030860	4.04962732	4.05934682	4.06451979	4.06935748
##	[311]	4.07351665	4.07615261	4.07921910	4.07986833	4.08013891
##	[316]	4.08436169	4.08575326	4.08631961	4.10096657	4.11168927
##	[321]	4.11698633	4.12011892	4.12030145	4.12606052	4.13802174
##	[326]	4.14492441	4.14500656	4.15166682	4.15606168	4.15659248
##	[331]	4.15783519	4.15844876	4.16021689	4.16253298	4.18248717
##	[336]	4.18372571	4.18909802	4.20255923	4.20287881	4.20337597
##	[341]	4.20905744	4.21027604	4.22454598	4.22455396	4.22956248
##	[346]	4.23231015	4.23405103	4.23583784	4.23819240	4.24227622
##	[351]	4.24659044	4.25385768	4.25777043	4.26036465	4.26445989
##	[356]	4.26580175	4.27722985	4.28061770	4.28218355	4.30268057
##	[361]	4.30274626	4.32204272	4.32583042	4.34837280	4.34935072
##	[366]	4.35734070	4.35836029	4.36613056	4.36905871	4.37325871
##	[371]	4.37830955	4.37963220	4.38337637	4.38434860	4.39301146
##	[376]	4.39691753	4.40366445	4.41284836	4.41613792	4.41681460
##	[381]	4.42218496	4.42914724	4.43872900	4.44167825	4.44230429
##	[386]	4.44623367	4.45513316	4.45728677	4.45777263	4.45977459
##	[391]	4.46042816	4.47315222	4.47341891	4.48405738	4.48619896
##	[396]	4.48842625	4.49614995	4.50288023	4.50401936	4.50965422
##	[401]	4.51306778	4.51570666	4.52620371	4.53145419	4.53193546
##	[406]	4.53421678	4.54403904	4.54544148	4.55732130	4.56028714
##	[411]	4.56138625	4.56235018	4.56252108	4.56419603	4.56620699
##	[416]	4.57903711	4.58827658	4.58931447	4.59815190	4.60106470
##	[421]	4.60369899	4.61090488	4.61998370	4.62516356	4.62836331
##	[426]	4.63041684	4.63131522	4.63524489	4.63544745	4.64594528
##	[431]	4.64941629	4.64941996	4.65279650	4.65672504	4.66437109
##	[436]	4.66620689	4.66691680	4.67063417	4.68468651	4.69980501
##	[441]	4.70187781	4.70895886	4.71239349	4.71398302	4.71635116
##	[446]	4.71962200	4.72098194	4.72285675	4.72683177	4.72809556
##	[451]	4.73322000	4.73444722	4.74423149	4.75003686	4.75552732
##	[456]	4.75625550	4.76386213	4.76415192	4.76435277	4.77468068
##	[461]	4.78086551	4.78649718	4.78677733	4.78682435	4.79205784
##	[466]	4.79496292	4.79996314	4.80131839	4.80345578	4.80664397
##	[471]	4.81314971	4.81335887	4.81433156	4.82153985	4.82586248
##	[476]	4.83717037	4.84244778	4.84318606	4.84471517	4.84636633
##	[481]	4.85155148	4.85690565	4.85833171	4.86060051	4.86695864
##	[486]	4.87589466	4.87696698	4.87764565	4.88108345	4.92577982
##	[491]	4.92720478	4.92931722	4.92938699	4.93323893	4.93337122
##	[496]	4.93415475	4.93625551	4.94117859	4.94173232	4.95036805
##	[501]	4.95352560	4.96165619	4.96621465	4.96859981	4.97089387
##	[506]	4.97398882	5.00121261	5.00538212	5.01146263	5.01384018
##	[511]	5.01395529	5.01516776	5.02066096	5.02481496	5.02635428
##	[516]	5.02725849	5.03100674	5.03545869	5.03782800	5.05358326
##	[521]	5.05597167	5.05997768	5.06538845	5.07031513	5.07674708
##	[526]	5.07704334	5.08546073	5.09363800	5.10237835	5.10698528
##	[531]	5.11113070	5.11803101	5.12688984	5.12706517	5.13140434

##	[536]	5.14604222	5.16246998	5.17980310	5.18235546	5.19011452
##	[541]	5.19185767	5.19449644	5.19851675	5.19913544	5.20673147
##	[546]	5.21160511	5.21465557	5.21493773	5.21765280	5.21973529
##	[551]	5.22539722	5.23339211	5.24567053	5.25005415	5.25889400
##	[556]	5.26179521	5.27520061	5.28097540	5.28464377	5.29406106
##	[561]	5.29805759	5.30503315	5.30794819	5.30892420	5.31014338
##	[566]	5.31571883	5.31617153	5.32403404	5.33967122	5.33973769
##	[571]	5.34221980	5.34946305	5.35740523	5.35820711	5.37808704
##	[576]	5.37844495	5.39212685	5.39283535	5.40033200	5.40728223
##	[581]	5.41213556	5.43333091	5.44001557	5.45819481	5.46047153
##	[586]	5.46335796	5.46902799	5.47116483	5.47395051	5.47411296
##	[591]	5.48652153	5.49713766	5.50295075	5.50358912	5.50787320
##	[596]	5.50897007	5.51161434	5.51184048	5.51621159	5.51818103
##	[601]	5.52012727	5.52910074	5.53387542	5.53966888	5.54238008
##	[606]	5.55089042	5.55416897	5.55693478	5.56000829	5.56002912
##	[611]	5.56077419	5.56113076	5.56182896	5.56829934	5.57281591
##	[616]	5.57697967	5.58196384	5.59765919	5.59851818	5.59929037
##	[621]	5.59941991	5.60026739	5.60579673	5.62745023	5.63177594
##	[626]	5.64033135	5.64416530	5.65264876	5.65320907	5.65796778
##	[631]	5.66800033	5.66952922	5.68014681	5.68278780	5.68672541
##	[636]	5.69390639	5.69894634	5.70179403	5.74265968	5.75183381
##	[641]	5.75877653	5.76615862	5.76844027	5.78483515	5.78483886
##	[646]	5.78797843	5.78925585	5.79023620	5.79273433	5.79429256
##	[651]	5.80051311	5.81876355	5.82013255	5.82017015	5.82227102
##	[656]	5.82811013	5.83605848	5.83975732	5.84108253	5.84203485
##	[661]	5.84517784	5.85262021	5.86136581	5.86190318	5.86190392
##	[666]	5.86689081	5.87610061	5.88008567	5.88552488	5.89390049
##	[671]	5.91667122	5.91996452	5.93117299	5.94084428	5.94716996
##	[676]	5.94729836	5.95485897	5.96432537	5.96788813	5.97141379
##	[681]	5.97806985	5.98820622	5.99410469	5.99692182	5.99774524
##	[686]	6.00679810	6.01909539	6.03644274	6.04264914	6.04303696
##	[691]	6.04498570	6.04710531	6.05978805	6.06207311	6.06442442
##	[696]	6.06686707	6.06852720	6.06887783	6.08321732	6.08737333
##	[701]	6.10215328	6.10770121	6.11059399	6.11529596	6.11683134
##	[706]	6.11769756	6.11844230	6.11952047	6.12006514	6.13042852
##	[711]	6.13212514	6.13294152	6.14482411	6.15275767	6.15382478
##	[716]	6.15865925	6.16161979	6.16572529	6.17620677	6.17737072
##	[721]	6.18180982	6.18863240	6.18918016	6.19009820	6.19240524
##	[726]	6.19587198	6.20094989	6.20940869	6.21037568	6.21227914
##	[731]	6.21240687	6.22112211	6.22657865	6.23090937	6.23864156
##	[736]	6.24090629	6.24256280	6.24833852	6.24947761	6.25293645
##	[741]	6.25710568	6.26498851	6.27983004	6.28583547	6.30162303
##	[746]	6.30497820	6.30659979	6.31989443	6.33176409	6.33263611
##	[751]	6.33740066	6.34906022	6.35820937	6.35832455	6.37118414
##	[756]	6.38300812	6.38686578	6.38950560	6.38987288	6.39078192
##	[761]	6.39179530	6.39345397	6.39688486	6.39834340	6.40807268
##	[766]	6.41118991	6.42066613	6.42992441	6.43407757	6.44883427
##	[771]	6.45141413	6.45231701	6.45488372	6.46148981	6.46232873

##	[776]	6.46671796	6.47156367	6.49600468	6.50641430	6.50977008
##	[781]	6.51421666	6.52445707	6.52567548	6.53008852	6.54301857
##	[786]	6.55721224	6.56302983	6.56439293	6.57178089	6.57584589
##	[791]	6.58479210	6.58583325	6.58640144	6.59333655	6.60245757
##	[796]	6.60618102	6.61081118	6.61475479	6.62186239	6.62339424
##	[801]	6.62560196	6.62790874	6.65071417	6.66243547	6.67251352
##	[806]	6.68061368	6.70920733	6.71220359	6.72486707	6.72556584
##	[811]	6.72837023	6.73059887	6.74131671	6.74235209	6.76385137
##	[816]	6.76766581	6.77628496	6.78545593	6.78600332	6.79133921
##	[821]	6.79179445	6.79306048	6.80660799	6.82473470	6.82888724
##	[826]	6.83731947	6.84644935	6.85270937	6.85587657	6.85624555
##	[831]	6.86915798	6.88023095	6.88161480	6.88330162	6.88334851
##	[836]	6.89067202	6.89121579	6.89673562	6.90065010	6.90907115
##	[841]	6.90918201	6.90998987	6.96228176	6.96728381	6.97337717
##	[846]	6.98004759	6.98007378	6.98577699	6.98898522	6.99310468
##	[851]	6.99676928	7.00716853	7.01443664	7.01841628	7.02343207
##	[856]	7.02983168	7.04027324	7.04512647	7.04586758	7.04747222
##	[861]	7.05222187	7.05455098	7.06857980	7.07323254	7.07782742
##	[866]	7.09018637	7.09095061	7.09103765	7.10307306	7.13761014
##	[871]	7.14800766	7.14822727	7.19084074	7.22837723	7.22999816
##	[876]	7.23281976	7.23743165	7.24496039	7.26020676	7.26395016
##	[881]	7.28674751	7.29005847	7.29303653	7.30065346	7.30973170
##	[886]	7.32664438	7.33356097	7.34435322	7.34608054	7.34731652
##	[891]	7.37432766	7.38516607	7.39430019	7.40031033	7.40885686
##	[896]	7.41829991	7.41900368	7.46817921	7.52175194	7.52230139
##	[901]	7.56221099	7.57112789	7.57930352	7.58092216	7.59841350
##	[906]	7.61684835	7.63035783	7.63827571	7.64997338	7.66075388
##	[911]	7.67295083	7.68100482	7.69345737	7.72707155	7.73035912
##	[916]	7.73947046	7.74797377	7.74858349	7.75629338	7.77436512
##	[921]	7.77934240	7.79243959	7.79402067	7.80584685	7.82193890
##	[926]	7.82622540	7.83857988	7.84077738	7.85349050	7.86012248
##	[931]	7.87693739	7.89112381	7.89908856	7.90294501	7.90351537
##	[936]	7.90480929	7.90505722	7.90902382	7.91694906	7.92664911
##	[941]	7.92667783	7.94273781	8.02781838	8.03694329	8.03942951
##	[946]	8.05704670	8.05888797	8.08909533	8.09747670	8.12904701
##	[951]	8.14232012	8.14956661	8.15994965	8.17336926	8.18687845
##	[956]	8.20222723	8.21494264	8.22216088	8.24874470	8.26018368
##	[961]	8.29672475	8.30476792	8.38405279	8.39437847	8.47694831
##	[966]	8.48948571	8.49338578	8.51173690	8.51277220	8.54423838
##	[971]	8.55339253	8.56277531	8.57374779	8.69861379	8.78270887
##	[976]	8.82712365	8.85339348	8.85885463	8.86646136	8.87807179
##	[981]	9.05289332	9.19120279	9.32207425	9.36313668	9.44871276
##	[986]	9.62273152	9.64067818	9.64794046	9.69179875	9.69880084
##	[991]	9.72349588	9.92682852	9.94263946	9.96619218	10.01598475
##	[996]	10.08362377	10.24087712	10.42457197	11.26067939	11.54295365

```
# Voltamos para a vírgula como separador decimal:
options(OutDec = ',')
```

- O que houve?!
- O problema é que não há valores repetidos no conjunto de dados! Por isso, todos os 1000 valores são modais.
- Uma maneira de evitar isto é definir a moda como o **centro do intervalo mais curto que contém metade dos dados**. Usamos a função `mlv` (*most likely value*):

```
moda <- mlv(normal, method = 'venter')
moda
```

```
## [1] 5,037732
```

- Esta moda estimada pode nem estar no conjunto de dados:

```
moda %in% normal
```

```
## [1] FALSE
```

- Mas o resultado de `mlv()` é útil, pois nos diz que, embora não haja valores repetidos, valores próximos de 5 são mais frequentes, como mostra o histograma.

#### 6.2.3.1

#### Exercícios

- Arredonde os valores no vetor `normal` para 2 casas decimais e ache a(s) moda(s)
  1. com a função `mfv`, e
  2. com a função `mlv`, usando o método `venter`.

Considerando o histograma, qual das respostas você prefere? Por quê?

## 6.3

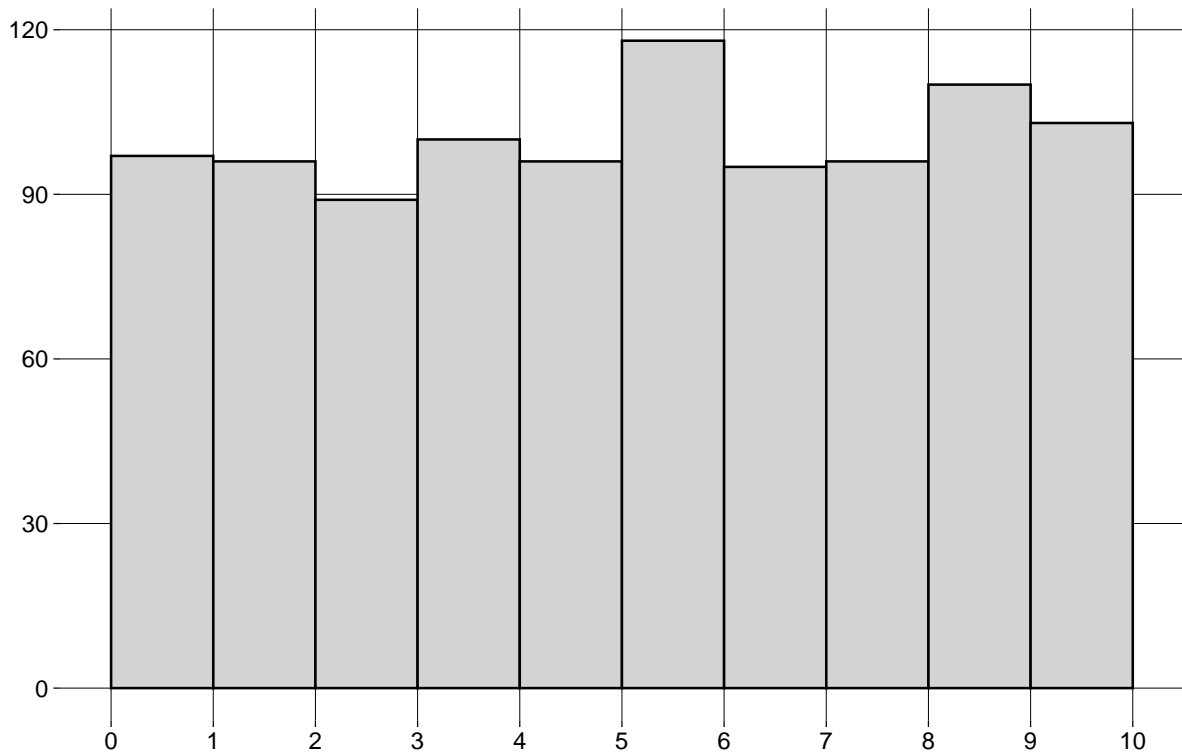
### Formas de uma distribuição

- A forma do histograma mostra aspectos importantes da distribuição dos dados.

### 6.3.1

#### Distribuição Uniforme

- Se o histograma tem todas as barras aproximadamente da mesma altura, dizemos que a distribuição é **uniforme**:



- A **distribuição uniforme não tem moda**, já que todos os valores têm aproximadamente a mesma frequência.

### 6.3.2

#### Simetria

- Se o histograma for simétrico (i.e., os lados esquerdo e direito são “espelhados”), dizemos que a distribuição é **simétrica**.
- A distribuição normal **do exemplo acima** é simétrica.
- A distribuição uniforme também é simétrica.
- Para distribuições simétricas, a média, a mediana e a moda **(quando existe e é única)** são bem próximas.
  - Para a distribuição normal do exemplo:

```
mean(normal)
```

```
## [1] 5,01997
```

```
median(normal)
```

```
## [1] 4,951947
```

```
mlv(normal, method = 'venter')
```

```
## [1] 5,037732
```

- Para a distribuição uniforme do exemplo:

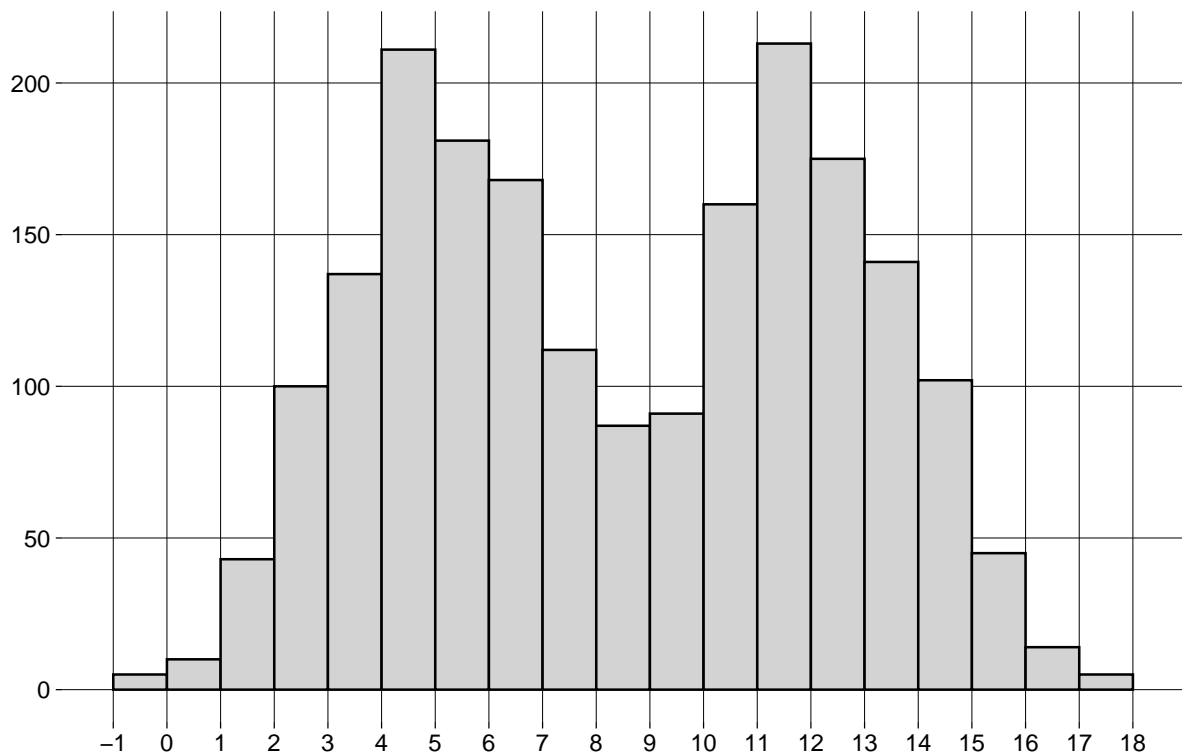
```
mean(uniforme)
```

```
## [1] 5,098863
```

```
median(uniforme)
```

```
## [1] 5,19571
```

- Uma distribuição pode ser **simétrica**, mas ter **duas (ou mais) modas diferentes**:



- Algumas distribuições não são simétricas, mas têm uma **cauda longa** à esquerda ou à direita.
- Dependendo da cauda, as distribuições são chamadas de **assimétricas à esquerda** ou **assimétricas à direita**.
- Um exemplo: receitas anuais (em milhões de dólares) de CEOs de grandes empresas:

```
df <- read_csv(
  './data/CEO_Salary_2012.csv',
  show_col_types = FALSE
)
glimpse(df)
```

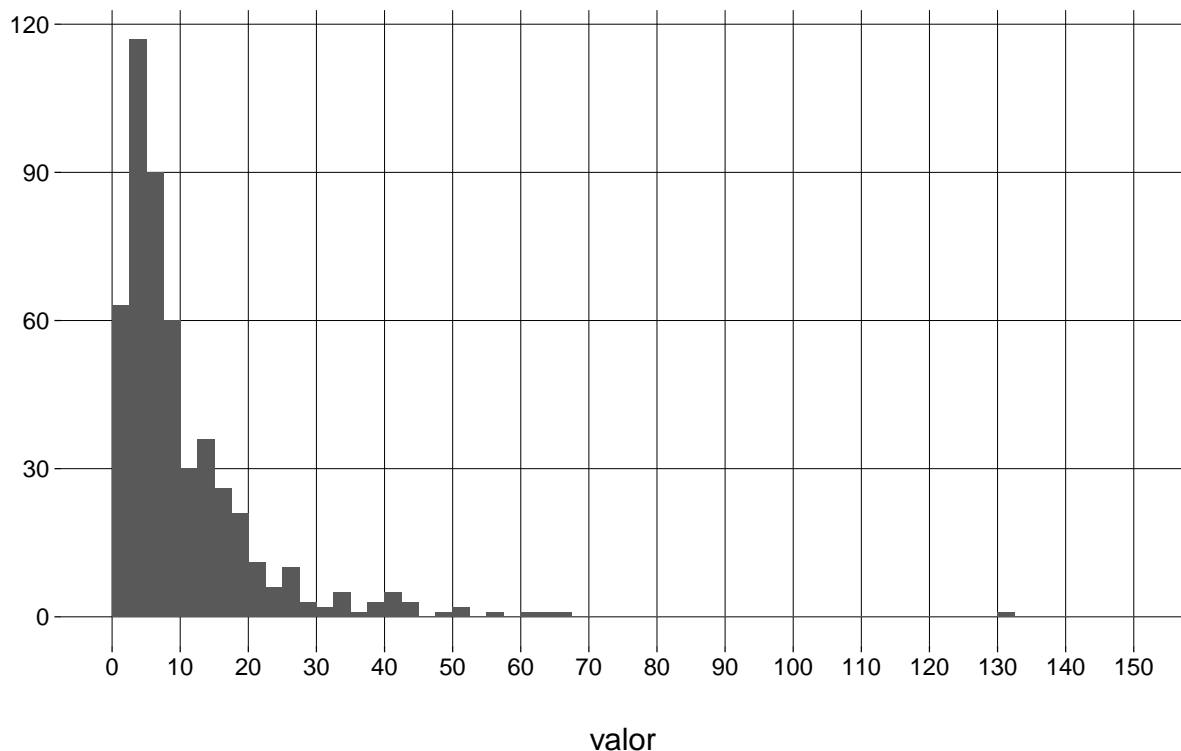
```
## Rows: 500
## Columns: 9
## $ Rank          <dbl> 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, ~
## $ Name          <chr> "John H Hammergren", "Ralph Lauren", "M~
## $ Company       <chr> "McKesson", "Ralph Lauren", "Vornado Re~
## $ `1-Year Pay ($mil)` <dbl> 131,190, 66,650, 64,405, 60,940, 55,790~
## $ `5 Year Pay ($mil)` <dbl> 285,020, 204,060, NA, 60,940, 96,110, 1~
## $ `Shares Owned ($mil)` <dbl> 51,9, 5010,4, 171,7, 8582,3, 21,5, 47,3~
## $ Age           <dbl> 53, 72, 55, 67, 59, 57, 55, 59, 61, 60,~
## $ Efficiency    <dbl> 121, 84, NA, NA, 138, 36, 12, NA, 91, 1~
## $ `Log Pay`     <dbl> 8,117901, 7,823800, 7,808920, 7,784902,~
```

- Vamos usar apenas os nomes e os valores anuais:

```
salarios <- df %>%
  select(Name, valor = `1-Year Pay ($mil)`)
```

- Um histograma:

```
salarios %>%
  ggplot(aes(x = valor)) +
    geom_histogram(breaks = seq(0, 150, 2.5)) +
    scale_x_continuous(breaks = seq(0, 150, 10)) +
    labs(y = NULL)
```



- É uma distribuição **assimétrica à direita**: a maior parte dos CEOs têm receitas anuais “baixas”, de menos de 10 milhões. À medida que examinamos valores maiores, a quantidade de CEOs vai diminuindo lentamente.
- Observe que **a longa cauda à direita “puxa” a média para um valor mais alto do que a mediana**.
- A moda, que corresponde à barra mais alta do histograma, é menor que a mediana (e que a média):

```
sumario <- salarios %>%
  summarise(
    moda = mlv(valor, method = 'venter'),
    mediana = median(valor),
    media = mean(valor)
  )

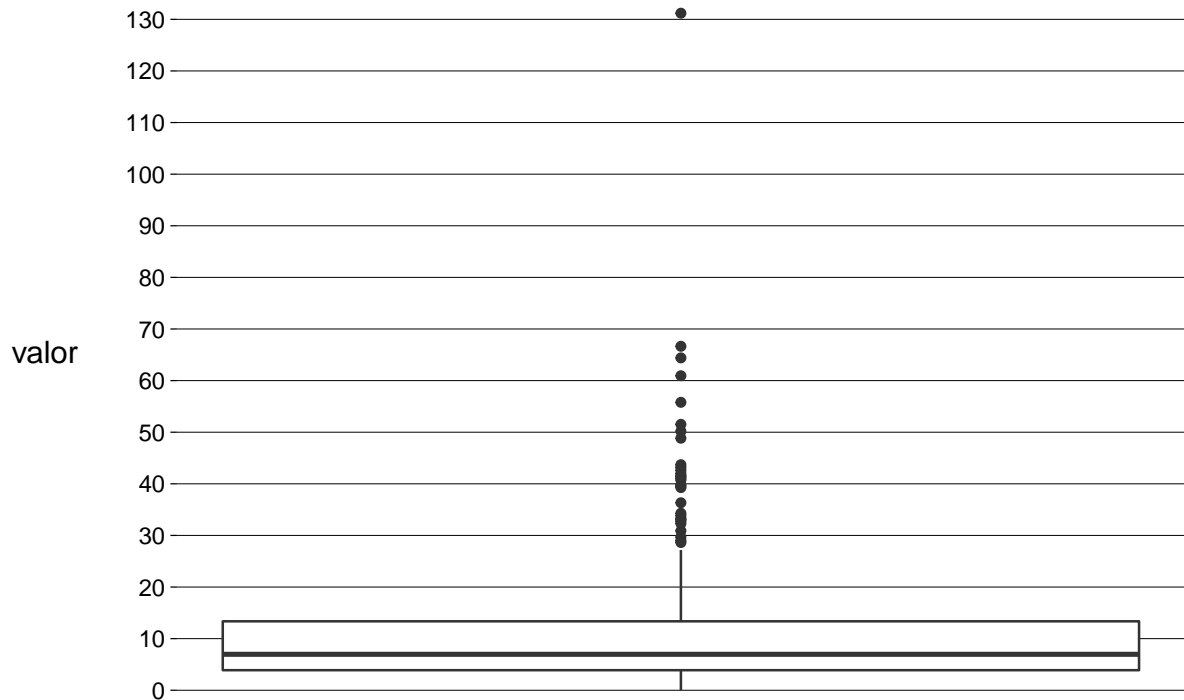
sumario
```

```
## # A tibble: 1 x 3
##   moda mediana media
##   <dbl>   <dbl> <dbl>
## 1  4.60    6.97  10.5
```

- Em um *boxplot*, também é possível detectar a assimetria pela grande quantidade de *outliers* em um extremo:



```
salarios %>%
  ggplot(aes(y = valor)) +
    geom_boxplot() +
    scale_x_continuous(breaks = NULL) +
    scale_y_continuous(breaks = seq(0, 150, 10))
```



- Com distribuições **assimétricas à esquerda**, a situação se inverte: **a média é menor que a mediana, que é menor que a moda.**

### 6.3.3

#### Exercícios

- Ache um conjunto de dados com uma distribuição assimétrica à esquerda.
- Faça um histograma.
- Calcule a média, a mediana, e a moda dos dados.

## 6.4

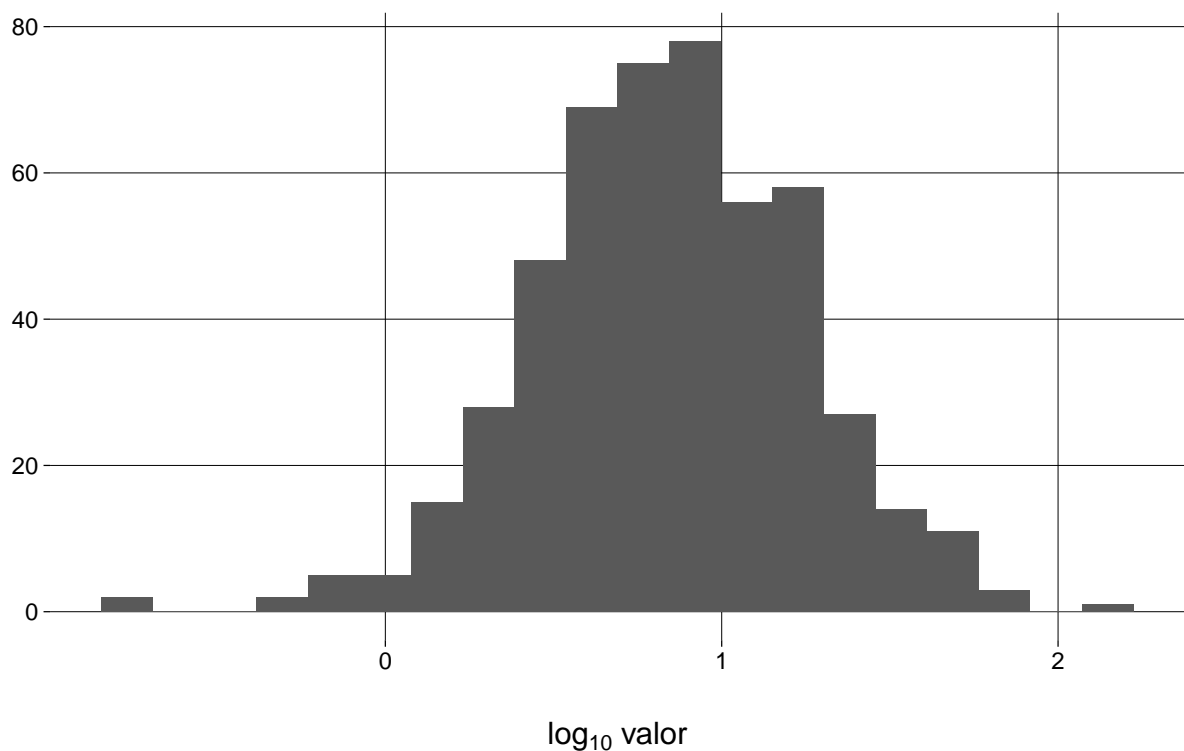
### Re-expressão

- Muitas vezes, é recomendável **transformar a escala dos dados** para que uma distribuição assimétrica se torne simétrica.
- No exemplo das receitas dos CEOs, podemos tomar os **logaritmos** dos valores, em vez dos valores:

```
salarios_log <- salarios %>%  
  mutate(log_valor = log10(valor))
```

```
salarios_log %>%  
  ggplot(aes(x = log_valor)) +  
    geom_histogram(bins = 20) +  
    labs(  
      x = TeX('$\\log_{10}$ valor'),  
      y = NULL  
    )
```

## Warning: Removed 3 rows containing non-finite values (stat\_bin).



- O logaritmo de um número na base 10 é, essencialmente, a quantidade de dígitos do número, vista como uma grandeza contínua.
- Logaritmos negativos vêm de valores entre 0 e 1.

- Logaritmo zero vem do valor 1.
- Valores iguais ou menores que zero não têm logaritmo definido.
- Por isso a mensagem de aviso sobre 3 valores removidos. São valores iguais a zero:

```
salarios_log %>%
  filter(valor == 0)
```

```
## # A tibble: 3 x 3
##   Name          valor log_valor
##   <chr>         <dbl>    <dbl>
## 1 Malon Wilkus      0      -Inf
## 2 Matthew J Lambiase 0      -Inf
## 3 Larry Page       0      -Inf
```

- Uma vantagem desta escala logarítmica é que podemos entender melhor o histograma. Os dados não estão amontoados de um lado só.

#### 6.4.1

##### Exercícios

- Quais são os registros com  $\log_{10} \text{valor} < 0$ ?
- Faça um *boxplot* dos logaritmos das receitas.

## 6.5

### Medidas de posição

#### 6.5.1

##### Quantis

- Na *seção sobre boxplots*, falamos sobre **quantis**, que são medidas de posição.
- Em R, a função `quantile` calcula quantis de um vetor:

```
salarios %>%
  pull(valor) %>%
  quantile()
```

```
##           0%          25%          50%          75%          100%
## 0,00000  3,88500  6,96750 13,36125 131,19000
```

- Você pode passar frações entre 0 e 1 para `quantile`. Por exemplo, para calcular o primeiro, o quinto, e o décimo percentis<sup>3</sup> das receitas dos CEOs:

<sup>3</sup>Um percentil é um quantil da forma  $k/100$ , para  $k$  natural,  $0 \leq k \leq 100$ .

```
salarios %>%  
  pull(valor) %>%  
  quantile(c(.01, .05, .1))
```

```
##      1%      5%     10%  
## 0,48695 1,48405 2,19400
```

## 6.6

---

### Medidas de dispersão

- Tão importantes quanto as medidas de centralidade são as medidas de dispersão (ou **espalhamento**).
- Elas informam **o quanto os dados variam**.

#### 6.6.1

---

##### Amplitude

- Uma medida simples é a **diferença entre o valor máximo e o valor mínimo**.
- Lembrando do nosso exemplo das idades dos alunos:

```
idades
```

```
## [1] 20 20 20 20 20 20 20 21 21 21 21 22 22 22 23 23 23 23 24 24 65
```

- A função `range` retorna o mínimo e o máximo:

```
range(idades)
```

```
## [1] 20 65
```

- A amplitude destes dados é, então

```
range(idades)[2] - range(idades)[1]
```

```
## [1] 45
```

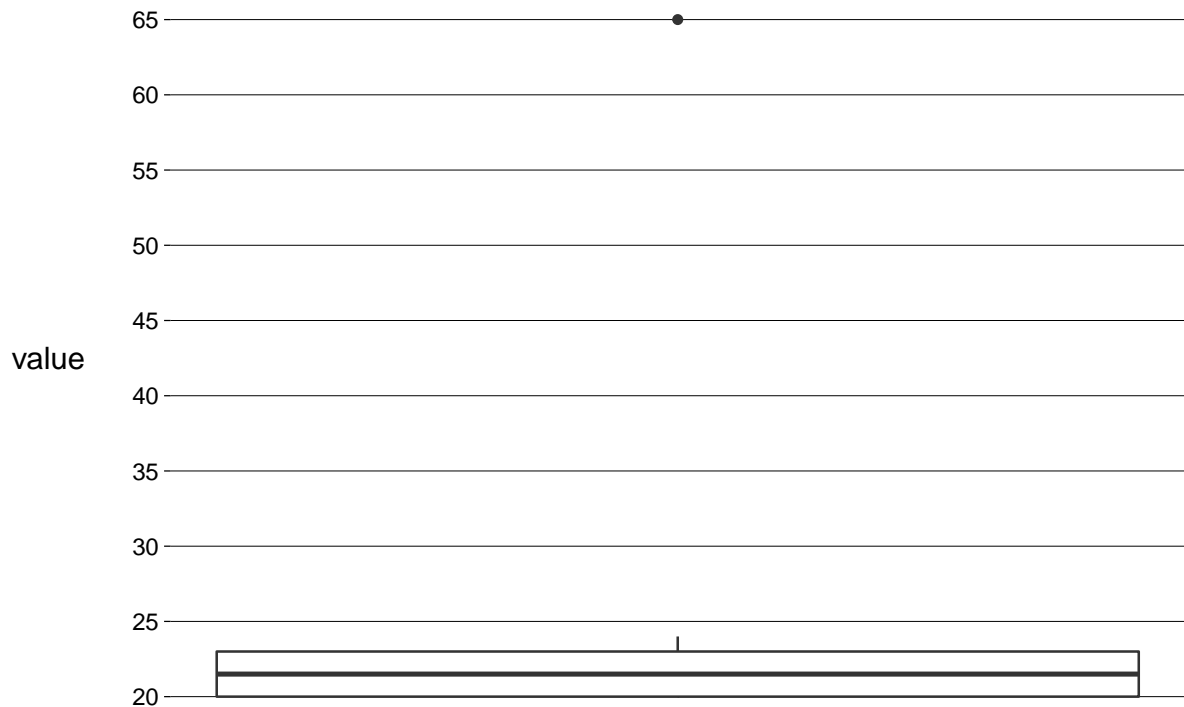
- A diferença de idade entre o aluno mais novo e o mais velho é de 45 anos, um valor alto, por causa do velhinho.

## 6.6.2

### IQR

- Na *seção sobre boxplots*, também falamos sobre o **intervalo interquartil** (IQR).
- No *boxplot*, é a **altura da caixa**. Para as idades dos alunos:

```
idades %>%  
  as_tibble() %>%  
  ggplot(aes(y = value)) +  
    geom_boxplot() +  
    scale_x_continuous(breaks = NULL) +  
    scale_y_continuous(breaks = seq(20, 70, 5))
```



- O IQR é a diferença entre o primeiro e o terceiro quartis:

```
summary(idades)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.   
##  20,00  20,00   21,50   23,75  23,00   65,00
```

```
unnamed[5] - unnamed[2]
```

```
## [1] 3
```

```
IQR(idades)
```

```
## [1] 3
```

- Ou seja, os 50% centrais dos alunos têm idade entre 20 e 23 anos, um IQR de 3.
- É uma variação pequena, porém mais fiel à realidade do que a amplitude, que é alta por causa do velhinho.
- Quanto maior o IQR, mais espalhados estão os dados.

### 6.6.3

#### Variância

- Agora, vamos trabalhar com os pesos (kg) e alturas (m) de um time de basquete:

```
medidas <- tibble(  
  altura = .025 *  
    c(72, 74, 68, 76, 74, 69, 72, 79, 70, 69, 77, 73),  
  peso = 0.45 *  
    c(180, 168, 225, 201, 189, 192, 197, 162, 174, 171, 185, 210)  
)  
  
medidas
```

```
## # A tibble: 12 x 2  
##   altura peso  
##   <dbl> <dbl>  
## 1    1.8   81  
## 2    1.85  75.6  
## 3    1.7  101.  
## 4    1.9   90.4  
## 5    1.85  85.0  
## 6    1.72  86.4  
## # ... with 6 more rows
```

```
summary(medidas$altura)
```

```
##   Min. 1st Qu.  Median    Mean 3rd Qu.    Max.  
##  1,700  1,744   1,812   1,819  1,863   1,975
```

```
summary(medidas$peso)
```

```
##   Min. 1st Qu.  Median    Mean 3rd Qu.    Max.  
##   72,90  77,96   84,15   84,53  89,10  101,25
```

- A variância é a maneira mais usada de medir o espalhamento em torno da média.

- Para calcular a variância das alturas e a variância dos pesos, precisamos calcular valores intermediários.
- O **desvio** de um valor é a **diferença entre o valor e a média**. O desvio pode ser positivo ou negativo.

```
d_medidas <- medidas %>%
  mutate(
    d_altura = altura - mean(altura),
    d_peso = peso - mean(peso)
  )

d_medidas
```

```
## # A tibble: 12 x 4
##   altura peso d_altura d_peso
##   <dbl> <dbl>   <dbl>  <dbl>
## 1    1.8   81    -0.0188 -3.53
## 2    1.85  75.6    0.0312 -8.92
## 3    1.7  101.   -0.119  16.7
## 4    1.9   90.4    0.0813  5.92
## 5    1.85  85.0    0.0312  0.525
## 6    1.72  86.4   -0.0938  1.88
## # ... with 6 more rows
```

- Vamos calcular o desvio médio das alturas e o desvio médio dos pesos:

```
d_medidas %>%
  summarize(
    d_medio_altura = mean(d_altura),
    d_medio_peso = mean(d_peso)
  )
```

```
## # A tibble: 1 x 2
##   d_medio_altura d_medio_peso
##           <dbl>         <dbl>
## 1              0    -3.55e-15
```

- Não foi uma boa idéia. **O desvio médio sempre é igual a zero.**<sup>4</sup> (O R pode mostrar algum valor diferente de zero por causa da precisão limitada dos números de ponto flutuante.)
- Como resolver isto? Elevando os desvios ao quadrado:

```
dq_medidas <- d_medidas %>%
  mutate(
    dq_altura = d_altura^2,
```

---

<sup>4</sup>Você vai provar isto em um exercício.

```
    dq_peso = d_peso^2
  )
```

```
dq_medidas
```

```
## # A tibble: 12 x 6
##   altura peso d_altura d_peso dq_altura dq_peso
##   <dbl> <dbl>   <dbl>  <dbl>   <dbl>   <dbl>
## 1   1.8   81    -0.0188 -3.53   0.000352  12.4
## 2   1.85  75.6    0.0312 -8.92   0.000977  79.7
## 3   1.7  101.   -0.119  16.7    0.0141   280.
## 4   1.9   90.4    0.0813  5.92    0.00660   35.1
## 5   1.85  85.0    0.0312  0.525   0.000977   0.276
## 6   1.72  86.4   -0.0938  1.88    0.00879   3.52
## # ... with 6 more rows
```

- Agora temos os **desvios quadrados**, que são todos **positivos**.
- O **desvio quadrado médio** vai ser a **variância**:

```
dq_medidas %>%
  summarize(
    var_altura = mean(dq_altura),
    var_peso = mean(dq_peso)
  )
```

```
## # A tibble: 1 x 2
##   var_altura var_peso
##   <dbl>    <dbl>
## 1   0.00678    63.3
```

- Uma vantagem da variância é que *outliers* (que têm desvios quadrados maiores) contribuem mais do que elementos próximos à média (que têm desvios quadrados menores).
- Uma desvantagem da variância é que a **sua unidade é o quadrado da unidade dos valores**.
- Neste exemplo, as unidades são  $m^2$  e  $kg^2$ !

#### 6.6.4

#### Desvio-padrão

- É melhor trabalhar com **a raiz quadrada da variância**, que chamamos de **desvio-padrão**.
- As unidades são as mesmas que as unidades dos dados.



```
dq_medidas %>%
  summarize(
    dp_altura = sqrt(mean(dq_altura)),
    dp_peso = sqrt((mean(dq_peso)))
  )
```

```
## # A tibble: 1 x 2
##   dp_altura dp_peso
##   <dbl>    <dbl>
## 1    0.0824    7.96
```

- Claro que o R tem funções para calcular isso: `var` e `sd` (*standard deviation*):

```
medidas %>%
  summarize(
    altura_var = var(altura),
    altura_dp = sd(altura),
    peso_var = var(peso),
    peso_dp = sd(peso)
  )
```

```
## # A tibble: 1 x 4
##   altura_var altura_dp peso_var peso_dp
##   <dbl>    <dbl>    <dbl>    <dbl>
## 1    0.00740    0.0860    69.1    8.31
```

- Mas os valores são diferentes dos que calculamos. Por quê?

## 6.6.5

### Definições

- Para uma **população** com  $N$  elementos e média  $\mu$ , a **variância** é

$$\sigma^2 = \frac{\sum_{i=1}^N (x_i - \mu)^2}{N}$$

- e o **desvio-padrão** é

$$\sigma = \sqrt{\frac{\sum_{i=1}^N (x_i - \mu)^2}{N}}$$

- Para uma **amostra** com  $n$  elementos e média  $\bar{x}$ , a **variância** é

$$s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n - 1}$$

e o desvio-padrão é

$$s = \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n - 1}}$$

- Nós calculamos a versão populacional destas medidas.
- R calcula a versão amostral destas medidas.
- Reveja os cálculos e entenda a diferença.
- Note, também, que as medidas populacionais são representadas por letras gregas —  $\mu, \sigma^2, \sigma$  —, enquanto as medidas amostrais são representadas por letras latinas —  $\bar{x}, s^2, s$ .



Mais adiante no curso, você vai entender por que o denominador da variância amostral é  $n - 1$ , em vez de  $n$ .

Nada é por acaso, nem mesmo em Estatística.

#### 6.6.6

#### Exercícios

- Quando a variância e o desvio-padrão de um conjunto de dados são iguais a zero?
- Mostre que o desvio médio de qualquer conjunto de valores é igual a zero.

Ou seja, considere o conjunto

$$\{x_1, x_2, \dots, x_n\}$$

e prove que

$$\sum_{i=1}^n (x_i - \bar{x}) = 0$$

Manipule apenas as variáveis  $x_i$ . Não use exemplos, pois eles não provam o enunciado geral.

Dica: lembre que  $\bar{x} = \frac{\sum_{i=1}^n x_i}{n}$ .

## Coeficiente de variação

- Em um conjunto de dados, o desvio-padrão é uma medida importante da variação dos dados.
- Mas a unidade do desvio-padrão muda de um conjunto de dados para outro: alturas em metros, pesos em quilos etc.
- Podemos eliminar as unidades expressando o desvio-padrão em termos da média.
- O resultado é a fração  $\frac{\sigma}{\mu}$  (na população) ou  $\frac{s}{\bar{x}}$  na amostra.
- Esta fração é o coeficiente de variação (CV).
- O CV não tem unidades.
- Para as alturas do exemplo dos jogadores de basquete:  
A média das alturas é 1,82 metros.  
O desvio-padrão das alturas é 0,09 metros.  
O CV é aproximadamente 0,0473.

```
statip::cv(medidas$altura)
```

```
## [1] 0,04729982
```

Em outras palavras, para as alturas, um desvio-padrão corresponde a 4,73% da média.

- Para os pesos:  
A média dos pesos é 84,53 quilos.  
O desvio-padrão dos pesos é 8,31 quilos.  
O CV é aproximadamente 0,0983.

```
statip::cv(medidas$peso)
```

```
## [1] 0,09834649
```

Em outras palavras, para os pesos, um desvio-padrão corresponde a 9,83% da média.

- Segundo estes valores, a variação dos pesos é cerca de 2 vezes maior do que a variação das alturas.



O coeficiente de variação sempre faz sentido para dados do **nível racional** (veja a definição) — i.e., dados onde o zero é absoluto.

Para dados apenas intervalares, o uso do CV pode levar a conclusões absurdas, como você terá chance de ver no exercício.

## 6.7.1

### Exercícios

- Considere o seguinte conjunto de temperaturas (em graus Celsius):

```
celsius <- c(0, 10, 20, 30, 40)
```

- E as **mesmas temperaturas** (em graus Fahrenheit):

```
fahrenheit <- 9 * celsius / 5 + 32
```

- Calcule para cada um dos dois vetores acima:

1. A média,
2. O desvio-padrão,
3. O coeficiente de variação.

- As temperaturas são as mesmas (apenas em unidades diferentes), mas os CVs são diferentes. Por quê?

- Agora, convertamos as **mesmas temperaturas** para a **escala Kelvin**:

```
kelvin <- celsius + 273.15
```

- E para a **escala Rankine**:

```
rankine <- fahrenheit + 459.67
```

- Calcule para cada um dos dois vetores acima:

1. A média,
2. O desvio-padrão,
3. O coeficiente de variação.

- Compare:

1. As médias de `celsius` e `kelvin`,
2. As médias de `fahrenheit` e `rankine`,
3. Os desvios-padrão de `celsius` e `kelvin`,

4. Os desvios-padrão de fahrenheit e rankine,
  5. Os coeficientes de variação de kelvin e rankine.
- Explique o que houve.

## 6.8

---

### Escores-padrão

- Para qualquer conjunto de dados, a unidade usada é uma escolha arbitrária.
- Para alturas, por exemplo, podemos usar metros, centímetros, pés, polegadas etc.
- A escolha de unidades é tão arbitrária que podemos escolher uma unidade (que dificilmente vai ter nome) que faça com que a média do conjunto de dados seja zero e que o desvio-padrão seja igual a 1.
- Isto equivale a tomar, como unidade, o *desvio-padrão acima da média*.
- Os valores, nesta nova unidade, são chamados de *escores-padrão*.
- Dizemos que os valores foram *padronizados*.
- Vamos usar as alturas dos jogadores de basquete.
- Para fazer a altura média virar zero, basta subtrair, de cada altura, a altura média:

```
alturas <- medidas$altura
mean(alturas)
```

```
## [1] 1,81875
```

```
alturas_deslocadas <- alturas - mean(alturas)
mean(alturas_deslocadas)
```

```
## [1] 0
```

- Para fazer o desvio-padrão ser igual a 1, basta dividir estes valores pelo desvio-padrão dos dados originais:

```
sd(alturas)
```

```
## [1] 0,08602656
```

```
alturas_padronizadas <- alturas_deslocadas / sd(alturas)
sd(alturas_padronizadas)
```

```
## [1] 1
```

- Tome, por exemplo, o seguinte jogador:

```
altura <- alturas[1]
altura
```

```
## [1] 1,8
```

```
altura_padronizada <- (alturas[1] - mean(alturas)) / sd(alturas)
altura_padronizada
```

```
## [1] -0,217956
```

Faça as contas: o valor da altura padronizada deste jogador significa que a altura dele está 0,217956 **desvios-padrão abaixo da altura média**.

- **No geral:**

- Se a média for  $\bar{x}$ , e
- Se o desvio-padrão for  $s$ ,
- Os escores-padrão  $z_i$  vão ser

$$z_i = \frac{x_i - \bar{x}}{s}$$

- Em R, a função `scale` faz isso:

```
medidas <- medidas %>%
  mutate(altura_padronizada = scale(altura)[,1])

medidas %>%
  select(altura, altura_padronizada)
```

```
## # A tibble: 12 x 2
##   altura altura_padronizada
##   <dbl>         <dbl>
## 1   1.8         -0.218
## 2   1.85         0.363
## 3   1.7        -1.38
## 4   1.9         0.944
## 5   1.85         0.363
## 6   1.72        -1.09
## # ... with 6 more rows
```

```
mean(medidas$altura_padronizada)
```

```
## [1] -0,000000000000000004610683
```

```
sd(medidas$altura_padronizada)
```

```
## [1] 1
```

- A função `scale` foi feita para receber e retornar **matrizes**. Como estamos trabalhando com **vetores**, usamos `scale(altura)[,1]` para tomar apenas a primeira (e única) coluna do resultado.

### 6.8.1

---

#### Exercícios

- Por que, **quando calculamos as alturas deslocadas divididas pelo desvio-padrão das alturas**, temos certeza de que a média dos valores resultantes não mudou?
- Padronize os pesos dos jogadores de basquete.
- Confira a média e o desvio-padrão dos pesos padronizados.
- Crie um *scatterplot* de peso por altura.
- Crie um *scatterplot* de peso padronizado por altura padronizada.
- Compare os dois *scatterplots*. O que muda de um para outro?

### 6.9

---

#### Teorema de Tchebychev

- Grosso modo, **quanto mais alto o desvio-padrão, maior é a distância da média até os valores**.
- Ou seja, **quanto menor o desvio-padrão, maior é a proporção de valores que estão próximos à média**.
- O teorema de Tchebychev quantifica esta idéia:  
Em **qualquer** distribuição, a proporção de valores **dentro de  $\pm k$  desvios-padrão ( $k > 1$ ) da média** é de, **no mínimo**

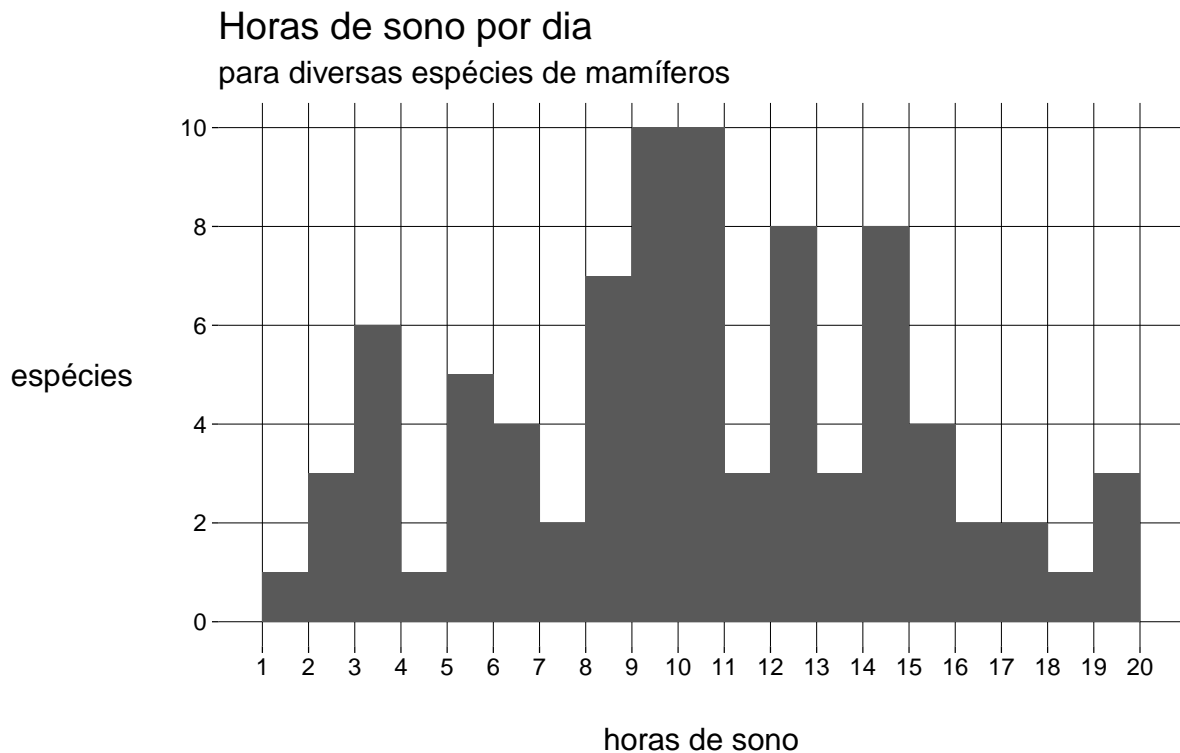
$$1 - \frac{1}{k^2}$$

### 6.9.1

---

#### Exemplo

- Lembre-se do **conjunto de dados sobre os totais de horas de sono de diversos mamíferos**:



- Média e desvio-padrão:

```
media <- mean(df$value)
media
```

```
## [1] 10,43373
```

```
dp <- sd(df$value)
dp
```

```
## [1] 4,450357
```

- Qual a proporção de espécies que estão a 1,3 ou menos desvios-padrão de distância da média?

```
k <- 1.3
inicio <- media - k * dp
inicio
```

```
## [1] 4,648271
```

```
fim <- media + k * dp
fim
```

```
## [1] 16,2192
```

- O teorema diz que no mínimo a seguinte proporção das espécies está dentro deste

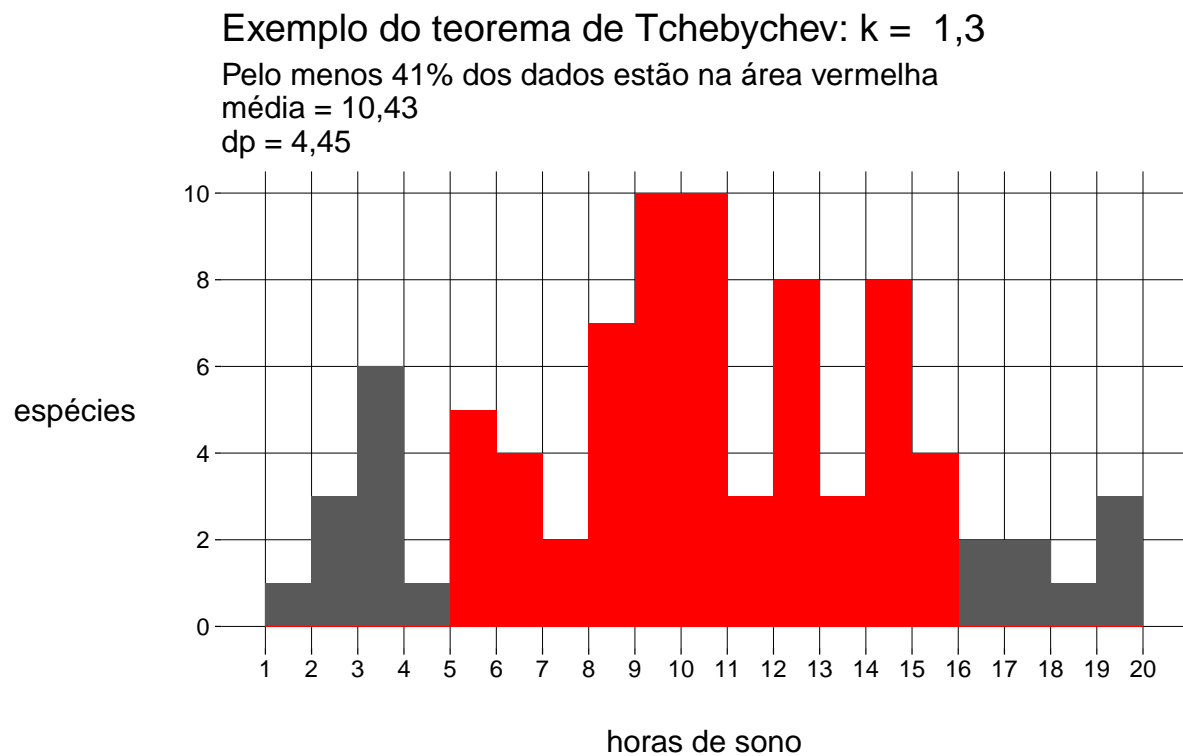


intervalo:

```
proporcao_teorema <- 1 - 1 / k^2  
proporcao_teorema
```

```
## [1] 0,408284
```

- Graficamente:



- Vamos conferir:

```
total_especies <- df %>% nrow()  
total_especies
```

```
## [1] 83
```

```
especies_intervalo <- df %>%  
  filter(value >= inicio & value <= fim) %>%  
  nrow()  
especies_intervalo
```

```
## [1] 64
```

```
proporcao_real <- especies_intervalo / total_especies  
proporcao_real
```

## [1] 0,7710843

- Como o teorema usa apenas a média e o desvio-padrão, e **mais nenhuma informação sobre a distribuição dos valores** — forma, simetria etc. — o que ele garante é, muitas vezes, mais fraco do que a realidade.
- Neste exemplo, o teorema garantia **no mínimo** 40,83% das espécies a 1,30 ou menos desvios-padrão de distância da média.
- A proporção verdadeira é 77,11% das espécies.
- O teorema está certo (claro), mas, sem mais informações sobre a distribuição dos dados, o teorema não pode ser mais preciso.

### 6.9.2

---

#### Exercício

- Uma loja recebe uma média de 20 clientes por dia, com desvio-padrão de 2 clientes.
- Em um ano (365 dias), em quantos dias, no mínimo, o número de clientes ficou entre 16 e 24, inclusive?