

Daten rekodieren

Unit 5

Ziele für heute

1. Datentypen und -klassen in einem Datensatz identifizieren und deren Bedeutung erklären
2. Spalten basierend auf Bedingungen umkodieren
3. Kategorien in Daten mit `{forcats}` umkodieren und sortieren

Datentypen

Warum sollten wir uns für Datentypen interessieren?

Beispiel

```
1 cat_lovers <- read_csv("data/cat-lovers.csv")
```

name	number_of_cats	handedness
Bernice Warren	0	left
Woodrow Stone	0	left
Willie Bass	1	left
Tyrone Estrada	3	left
Alex Daniels	3	left
Jane Bates	2	left
Latoya Simpson	1	left
Darin Woods	1	left
Annes Cobb	0	left

Durschnittliche Anzahl

```
1 cat_lovers |>
2   summarise(mean_cats = mean(number_of_cats))
```

```
1 # A tibble: 1 × 1
2   mean_cats
3     <dbl>
4   1       NA
```

Warum funktioniert es nicht?!

```
1 cat_lovers |>
2   summarise(mean_cats = mean(number_of_cats, na.rm = TRUE))
```

```
1 # A tibble: 1 × 1
2   mean_cats
3     <dbl>
4   1       NA
```

Warum funktioniert es immer noch nicht?!

Einatmen... und sich die Daten anschauen

Welchen Typ hat die Variable `number_of_cats`?

```
1 cat_lovers |>  
2 glimpse()
```

```
1 Rows: 60  
2 Columns: 3  
3 $ name           <chr> "Bernice Warren", "Woodrow Stone", "Willie Bass", "Tyro...  
4 $ number_of_cats <chr> "0", "0", "1", "3", "3", "2", "1", "1", "0", "0", "...  
5 $ handedness     <chr> "left", "left", "left", "left", "left", "left", "left", ...
```

Noch Einmal einen Blick darauf Werfen

```
1 cat_lovers |> count(number_of_cats)

1 # A tibble: 7 × 2
2   number_of_cats      n
3   <chr>                <int>
4   1 0                  32
5   2 1                  15
6   3 1.5 - honestly I think one of my cats is half human  1
7   4 2                  4
8   5 3                  6
9   6 4                  1
10  7 three               1
```

Manchmal musst Du auf deine Befragten aufpassen

```

1 cat_lovers |>
2   mutate(number_of_cats = case_when(
3     number_of_cats == "1.5 - honestly I think one of my cats is half human" ~ 2,
4     number_of_cats == "three" ~ 3,
5     .default = as.numeric(number_of_cats)
6   )) |>
7   summarise(mean_cats = mean(number_of_cats))

```

Warning: There was 1 warning in `mutate()`.
 i In argument: `number_of_cats = case_when(...)`.
 Caused by warning in `vec_case_when()`:
 ! NAs introduced by coercion

```

1 # A tibble: 1 × 1
2   mean_cats
3   <dbl>
4 1     0.833

```

Immer Datentypen Respektieren!

```

1 cat_lovers |>
2   mutate(
3     number_of_cats = case_when(
4       number_of_cats == "1.5 - honestly I think one of my cats is half human" ~ "2",
5       number_of_cats == "three" ~ "3",
6       .default = number_of_cats
7     ),
8     number_of_cats = as.numeric(number_of_cats)
9   ) |>
10   summarise(mean_cats = mean(number_of_cats))

```

```

1 # A tibble: 1 × 1
2   mean_cats
3   <dbl>
4   1     0.833

```

Jetzt, wo wir wissen, was wir tun...

```
1 cat_lovers <- cat_lovers |>
2   mutate(
3     number_of_cats = case_when(
4       number_of_cats == "1.5 - honestly I think one of my cats is half human" ~ "2",
5       number_of_cats == "three" ~ "3",
6       .default = number_of_cats
7     ),
8     number_of_cats = as.numeric(number_of_cats)
9   )
```

Moral der Geschichte

- Wenn sich Deine Daten nicht so verhalten, wie du es erwartest, könnte ein *type coercion* beim Einlesen der Daten die Ursache sein.
- Gehe hinein, untersuche deine Daten, wende den Fix an, speichere deine Daten und lebe glücklich bis ans Ende deiner Tage.

Datentypen in R

Atomic vectors

logical: TRUE,
FALSE

character:
“Hallo”, “a”,
“TRUE”

integer: 2L, 34L,
0L

double: 1, 2.4,
pi

Datentypen in R

`typeof()` → wie R das Objekt speichert

logical (TRUE/FALSE)

```
1 typeof(TRUE)
1 [1] "logical"
1 typeof(c(TRUE, FALSE))
1 [1] "logical"
```

double (*floating point*)

```
1 typeof(3.56)
1 [1] "double"
1 typeof(c(4, 3))
1 [1] "double"
```

character (Text)

```
1 typeof("Hallo")
1 [1] "character"
1 typeof(c("a", "aa"))
1 [1] "character"
```

integer (Ganzzahl)

```
1 typeof(4L)
1 [1] "integer"
1 typeof(1:4)
1 [1] "integer"
```

Expliziter vs. Impliziter Typenzwang

- ***Explicit coercion*** `as.logical()`, `as.numeric()`,
`as.integer()`, `as.double()`, `as.character()`.

```
1 x <- c(TRUE, FALSE)
2 as.character(x)

1 [1] "TRUE"   "FALSE"
```

- ***Implicit coercion*** z. B. R konvertiert Variablen
 gemischter Typen in einen einzelnen Typ.

```
1 c(15, "Danke")
1 [1] "15"      "Danke"

1 c(3L, pi)
1 [1] 3.000000 3.141593
```

... und das ist nicht immer eine gute Sache!

Praktikum: *Type Coercion*

prak-05a-type-coercion.qmd

Welcher Typ sind die angegebenen Vektoren?

Daten-Rekodierung

`if_else()`, `case_when()`

TRUE/FALSE: `if_else()`

Schnabellänge kategorisieren: “überdurchschnittlich”,
“unterdurchschnittlich”

```
1 penguins |>
2   summarise(median_bill_length = median(bill_length_mm, na.rm = TRUE))

1 # A tibble: 1 × 1
2   median_bill_length
3             <dbl>
4   1                 44.4
```

TRUE/FALSE: `if_else()`

`if_else(stimmt_das, das_passiert, sonst_das_passiert)`

```
1 penguins |>
2   mutate(
3     bl_cat = if_else(bill_length_mm < 44.45, "unterdurchschnittlich", "überdurchschnittlich"),
4   ) |>
5   count(bl_cat)
```

```
1 # A tibble: 3 × 2
2   bl_cat           n
3   <chr>          <int>
4   1 unterdurchschnittlich    171
5   2 überdurchschnittlich    171
6   3 <NA>                  2
```

TRUE/FALSE: `if_else()`

`if_else(stimmt_das, das_passiert, sonst_das_passiert,
NA_so_behandeln)`

```
1 penguins |>
2   mutate(
3     bl_cat =
4       if_else(
5         bill_length_mm < 44.45, "unterdurchschnittlich", "überdurchschnittlich",
6         )
7     ) |>
8     count(bl_cat)
```

```
1 # A tibble: 3 × 2
2   bl_cat           n
3   <chr>          <int>
4   1 unbekannt      2
5   2 unterdurchschnittlich 171
6   3 überdurchschnittlich 171
```

Mehrere Bedingungen: `case_when()`

Schnabellänge kategorisieren: *short, medium, long*

```
1 penguins |>  
2   select(bill_length_mm) |>  
3   summary()
```

```
1 bill_length_mm  
2 Min.    :32.10  
3 1st Qu.:39.23  
4 Median   :44.45  
5 Mean     :43.92  
6 3rd Qu.:48.50  
7 Max.     :59.60  
8 NA's      :2
```

Mehrere Bedingungen: case_when()

```

1 penguins |>
2   mutate(
3     bl_cat = case_when(
4       is.na(bill_length_mm) ~ NA,
5       bill_length_mm < 39.2 ~ "short",
6       between(bill_length_mm, 39.2, 48.5) ~ "medium",
7       .default = "long"
8     )
9   ) |>
10   count(bl_cat)

```

```

1 # A tibble: 4 × 2
2   bl_cat     n
3   <chr>   <int>
4   1 long      84
5   2 medium    175
6   3 short     83
7   4 <NA>       2

```

Praktikum: Daten rekodieren

prak-05b-cond-mutate.qmd

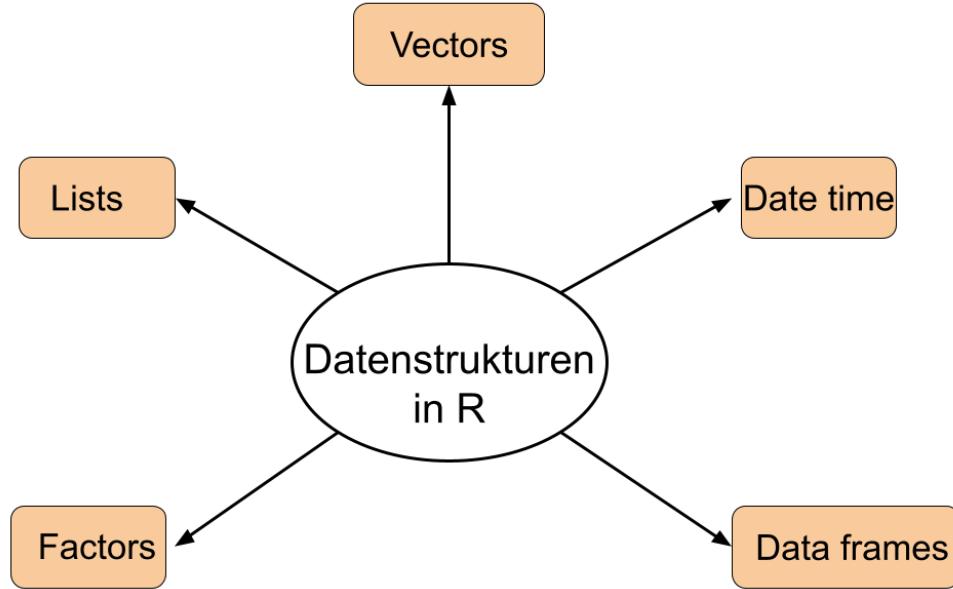
Break



10:00

Datenstrukturen

Datenstrukturen



`class()` → wie sich das
Objekt verhält

Factors

→ Kategoriale Variablen: Character + Ganzzahl

```
1 (x <- c("BS", "MS", "PhD", "MS"))
1 [1] "BS"   "MS"   "PhD"  "MS"
1 typeof(x)
1 [1] "character"
1 class(x)
1 [1] "character"
```

```
1 (y <- factor(x))
1 [1] BS   MS   PhD  MS
2 Levels: BS MS PhD
1 typeof(y)
1 [1] "integer"
1 class(y)
1 [1] "factor"
1 as.integer(y)
1 [1] 1 2 3 2
```

Dates

Ganzzahl = Anzahl Tage seit Ursprung

```
1 (y <- as.Date("1990-01-01"))
1 [1] "1990-01-01"
1 typeof(y)
1 [1] "double"
1 class(y)
1 [1] "Date"
1 as.integer(y)
1 [1] 7305
1 as.integer(y) / 365
1 [1] 20.0137
```



~ 20 Jahre nach dem 1970-01-01

Lists

Generische Vektorcontainers: Vektoren jeglicher Typ und Länge

```
1 l <- list(  
2   x = 1:4,  
3   y = c("Hallo", "hello", "salut"),  
4   z = c(TRUE, FALSE)  
5 )  
6 l  
  
1 $x  
2 [1] 1 2 3 4  
3  
4 $y  
5 [1] "Hallo" "hello" "salut"  
6  
7 $z  
8 [1] TRUE FALSE
```

Data Frames

Spezielle Liste mit Vektoren gleicher Länge

```
1 (df <- data.frame(x = 1:2, y = 3:4))
```

```
1 x y
2 1 1 3
3 2 2 4
```

```
1 class(df)
```

```
1 [1] "data.frame"
```

```
1 df |>
2 pull(y)
```

```
1 [1] 3 4
```

```
1 df$y
```

```
1 [1] 3 4
```

```
1 (df <- tibble(x = 1:2, y = 3:4))
```

```
1 # A tibble: 2 × 2
2   x     y
3   <int> <int>
4 1     1     3
5 2     2     4
```

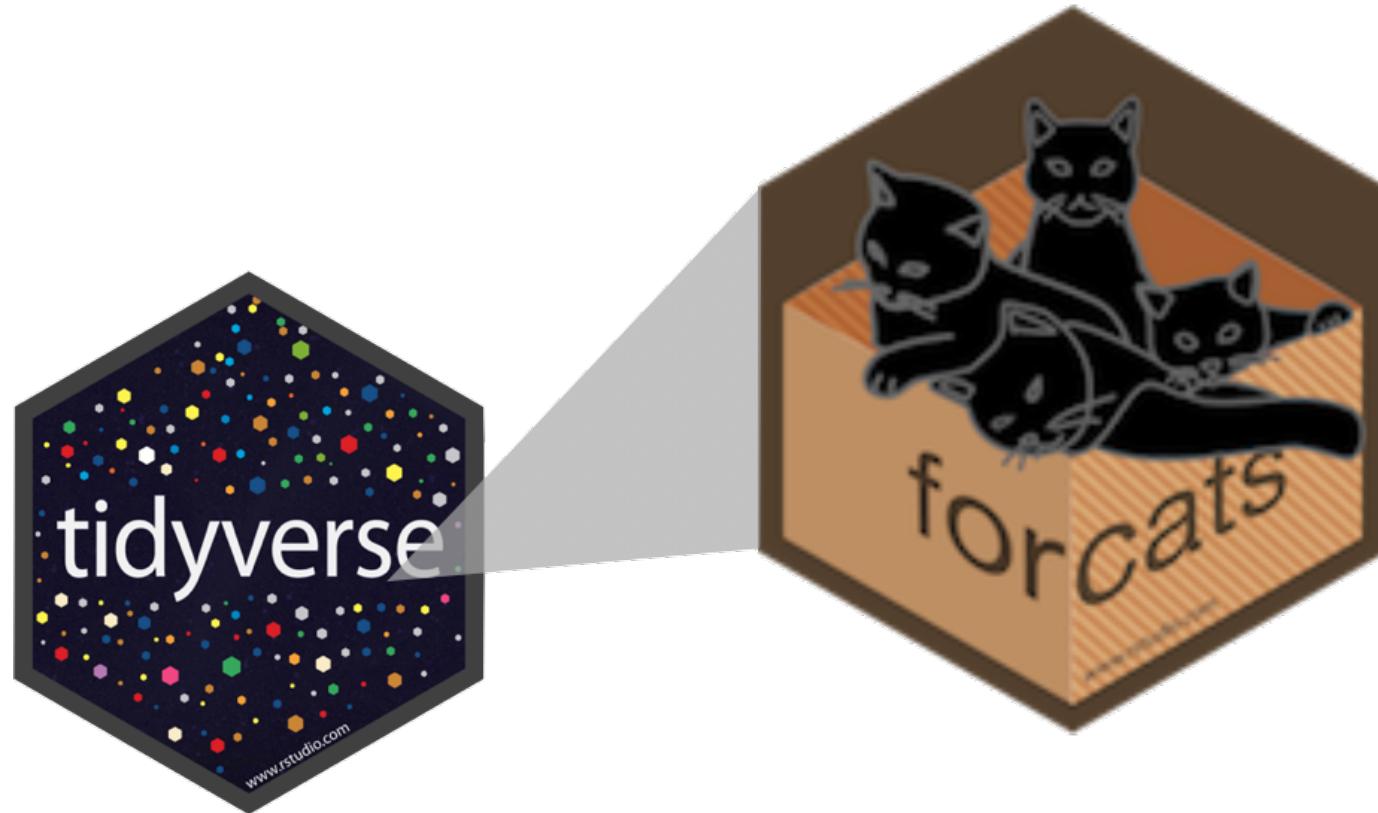
```
1 class(df)
```

```
1 [1] "tbl_df"      "tbl"
```

"data.frame"

rstatsBL - Data Science mit R

Mit Factors Arbeiten: {forcats}



Daten

```

1 penguins |>
2 glimpse()

1 Rows: 344
2 Columns: 8
3 $ species           <fct> Adelie, Adelie, Adelie, Adelie, Adelie, Adelie, Adel...
4 $ island             <fct> Torgersen, Torgersen, Torgersen, Torgersen, Torgers...
5 $ bill_length_mm    <dbl> 39.1, 39.5, 40.3, NA, 36.7, 39.3, 38.9, 39.2, 34.1, ...
6 $ bill_depth_mm     <dbl> 18.7, 17.4, 18.0, NA, 19.3, 20.6, 17.8, 19.6, 18.1, ...
7 $ flipper_length_mm <int> 181, 186, 195, NA, 193, 190, 181, 195, 193, 190, 186...
8 $ body_mass_g        <int> 3750, 3800, 3250, NA, 3450, 3650, 3625, 4675, 3475, ...
9 $ sex                <fct> male, female, female, NA, female, male, female, male...
10 $ year              <int> 2007, 2007, 2007, 2007, 2007, 2007, 2007, 2007, 2007...

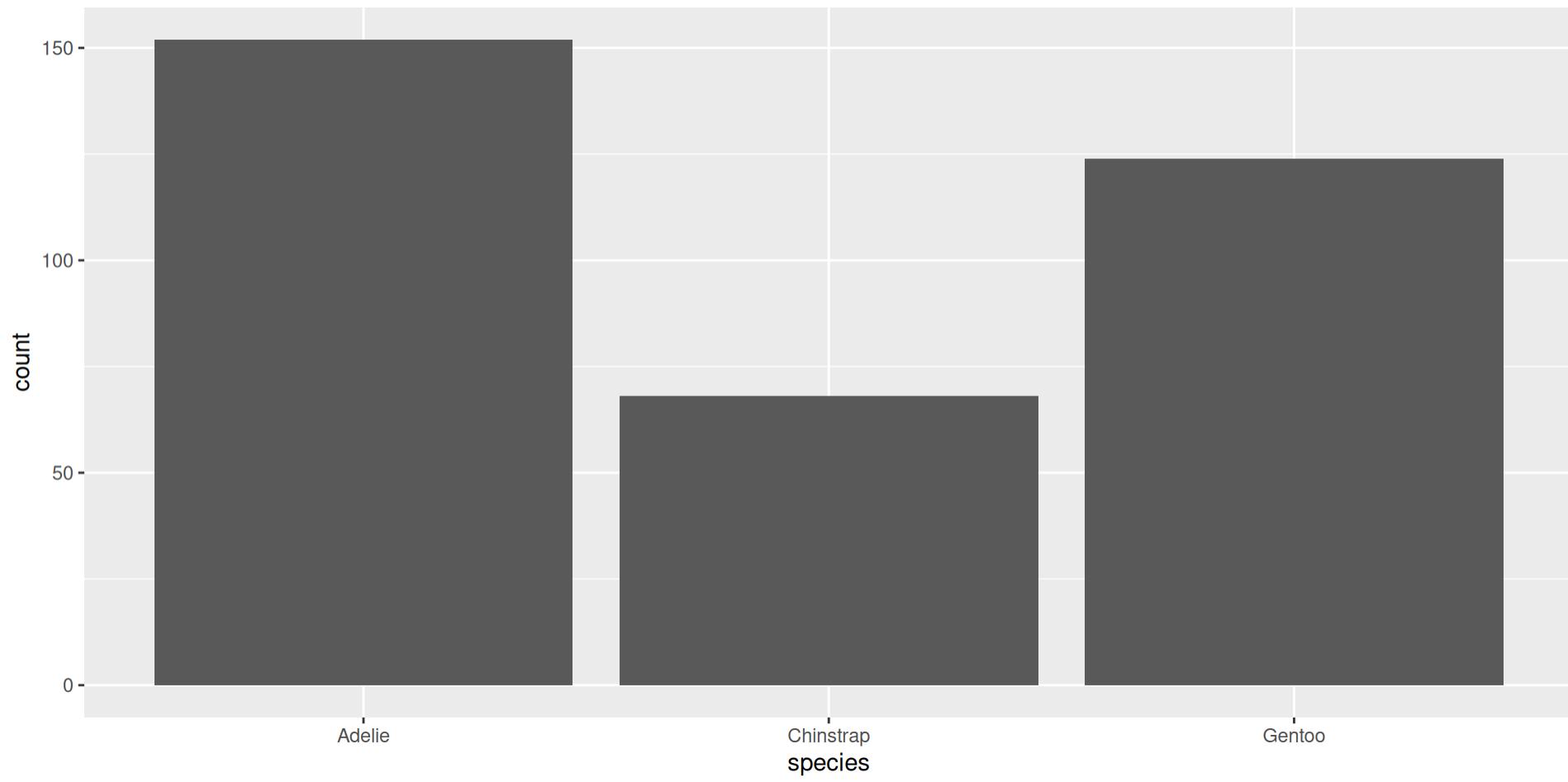
```

▼ Code

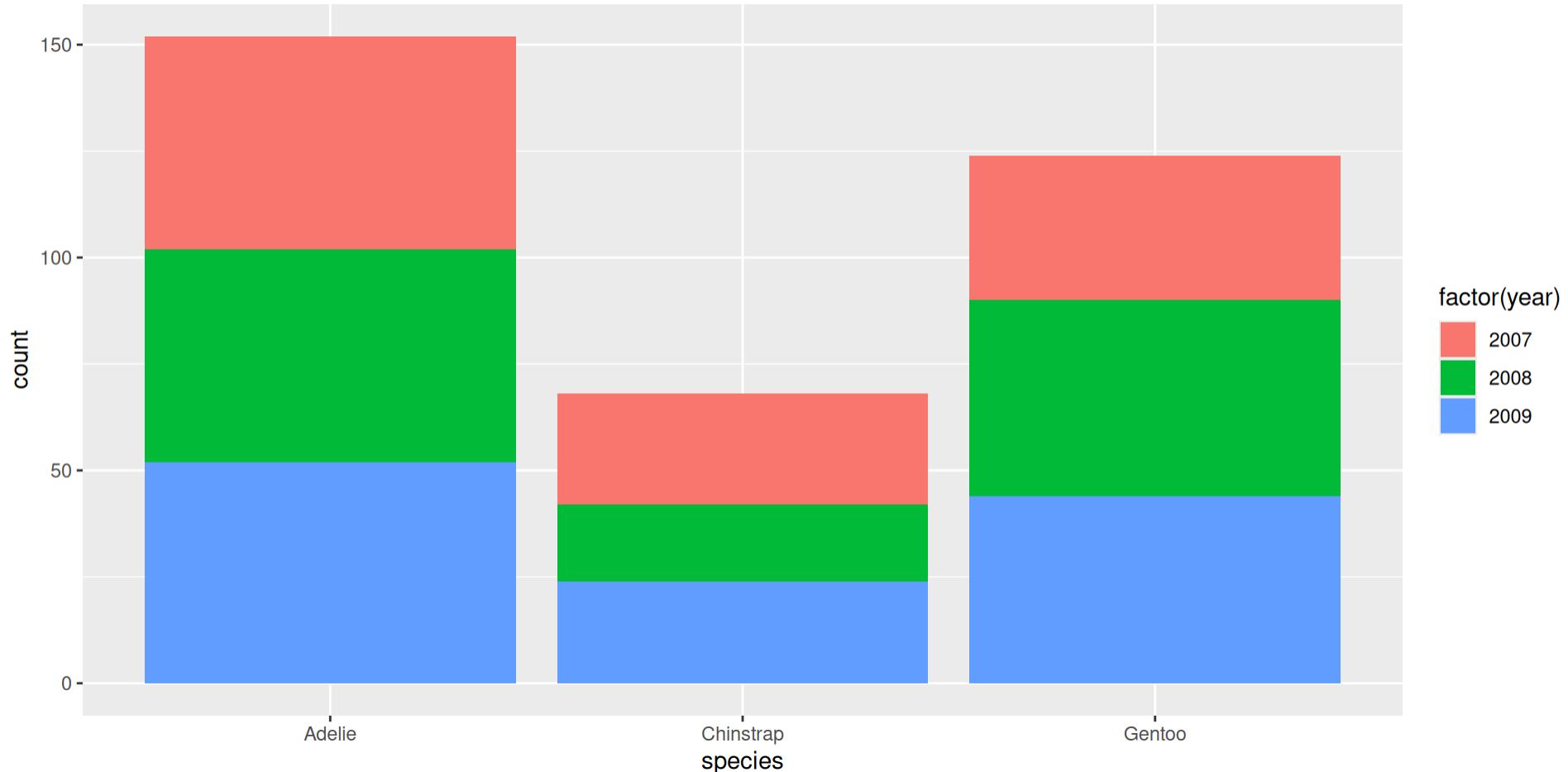
```
1 penguins |>  
2   ggplot(aes(x = species, fill = year)) +  
3     geom_bar()
```

Warning: The following aesthetics were dropped during statistical transformation:
`fill`.

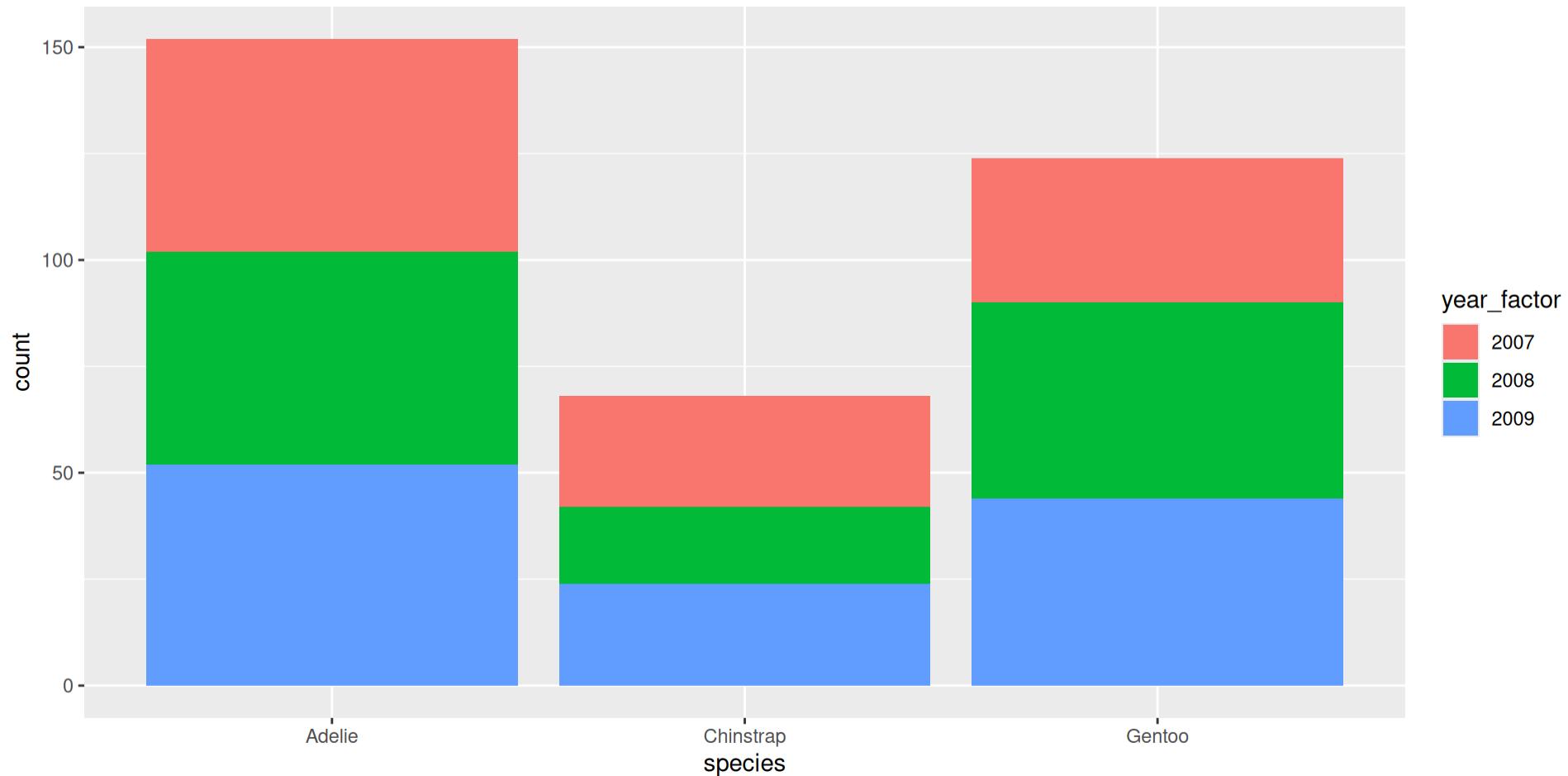
- This can happen when `ggplot` fails to infer the correct grouping structure in the data.
- Did you forget to specify a ``group`` aesthetic or to convert a numerical variable into a factor?



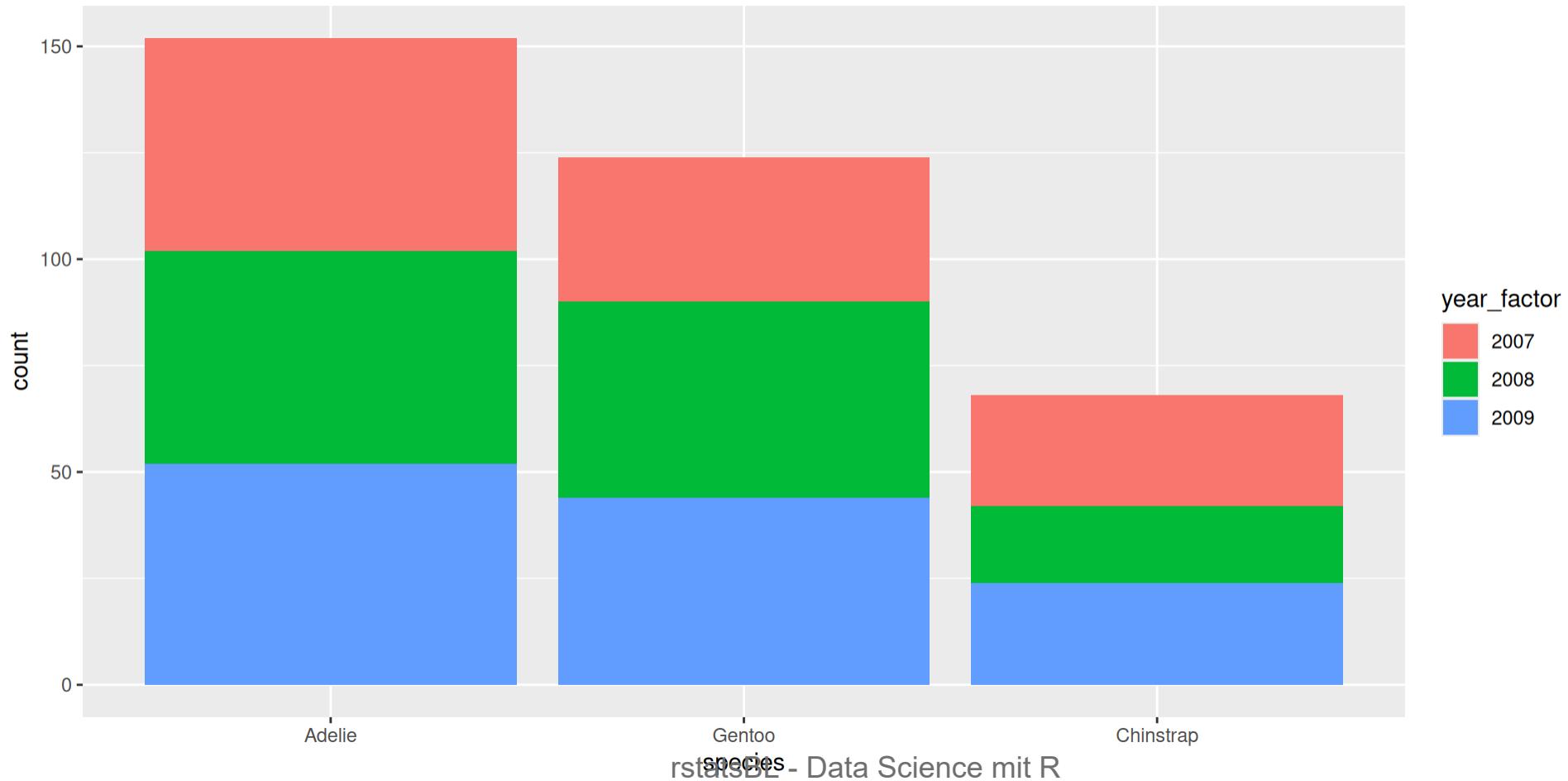
```
1 penguins |>
2   ggplot(aes(x = species, fill = factor(year))) +
3     geom_bar()
```



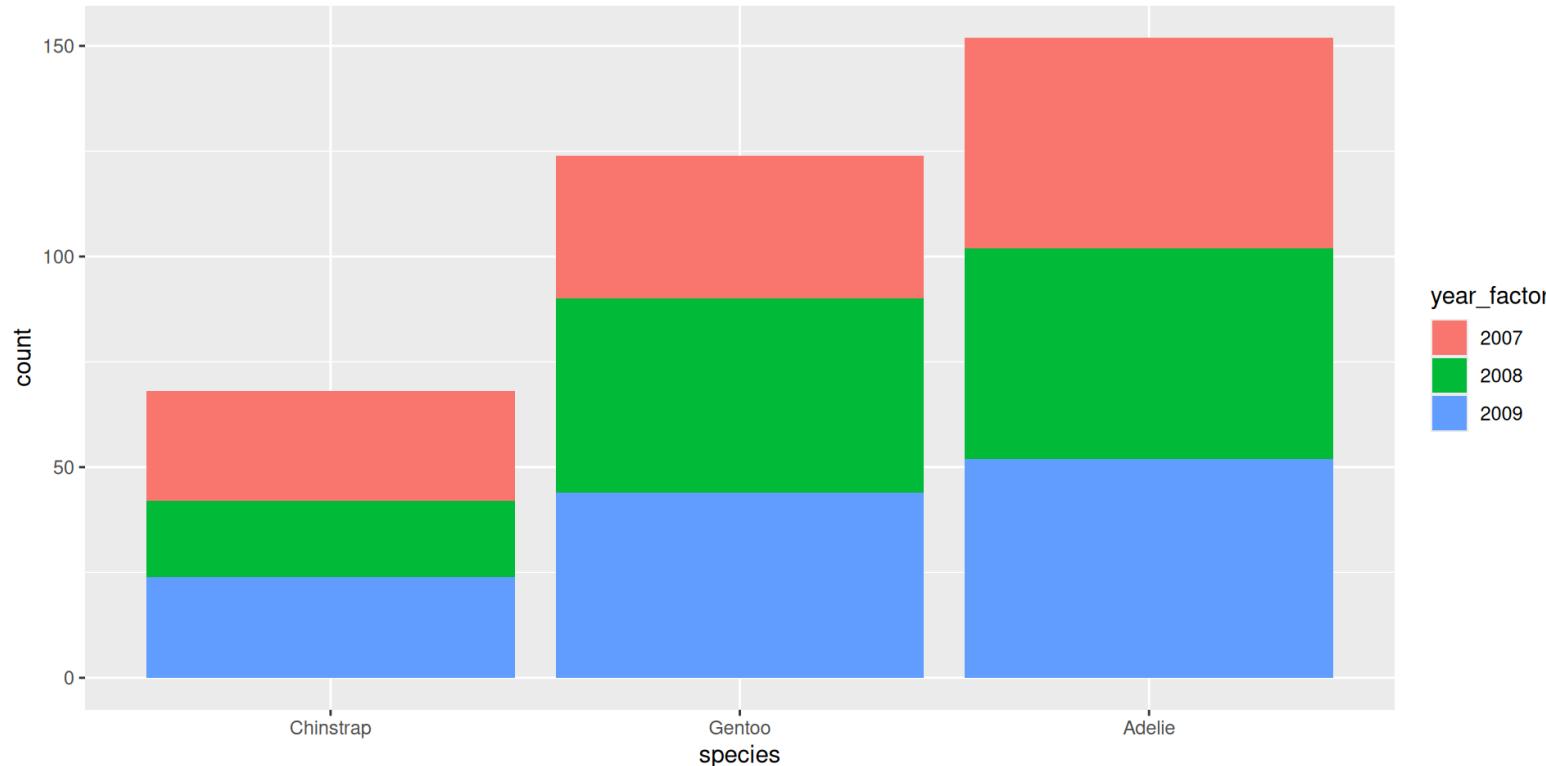
```
1 penguins |>
2   mutate(year_factor = factor(year)) |>
3   ggplot(aes(x = species, fill = year_factor)) +
4   geom_bar()
```



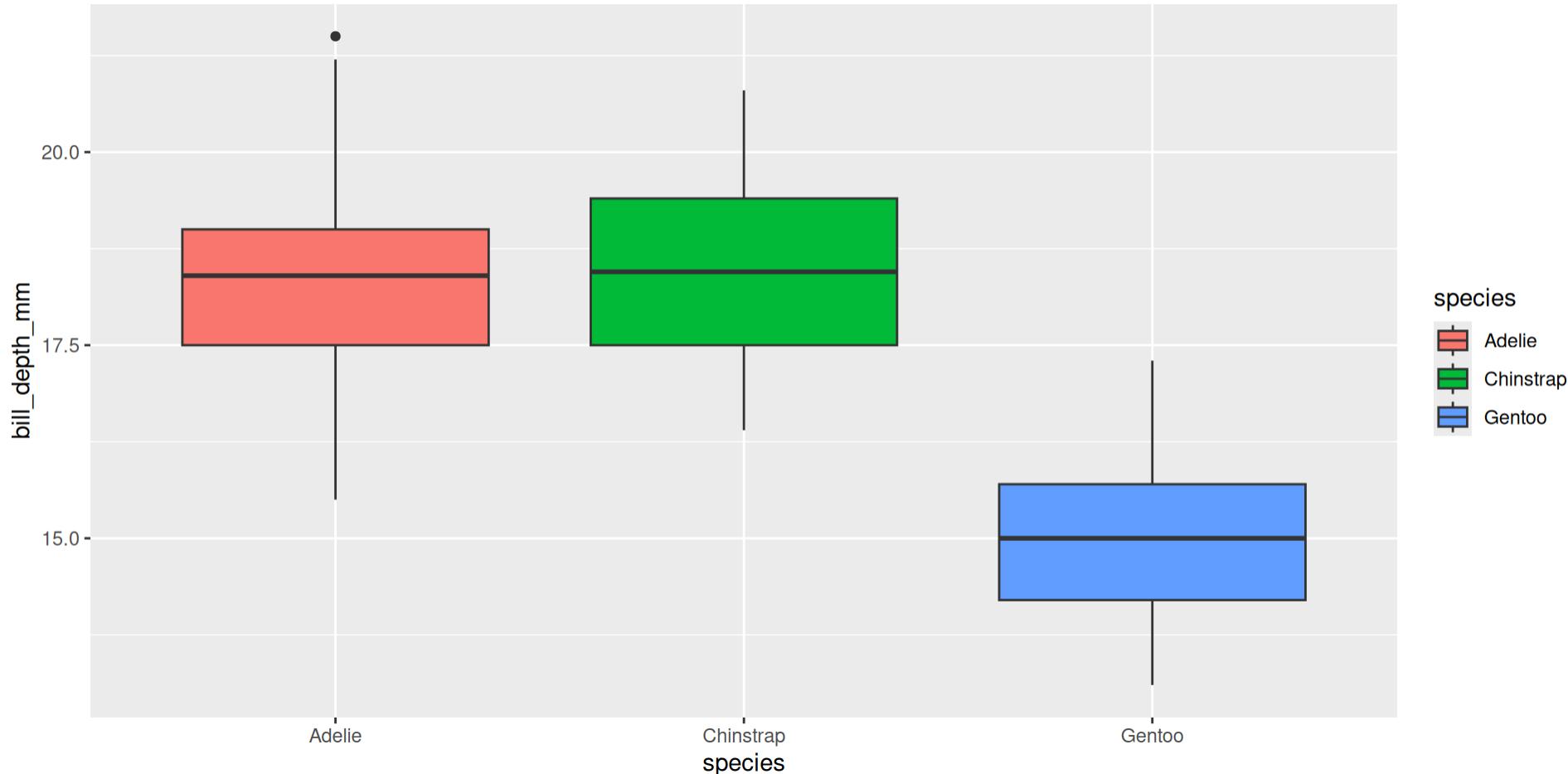
```
1 penguins |>
2   mutate(
3     year_factor = factor(year),
4     species = fct_infreq(species)
5   ) |>
6   ggplot(aes(x = species, fill = year_factor)) +
7   geom_bar()
```



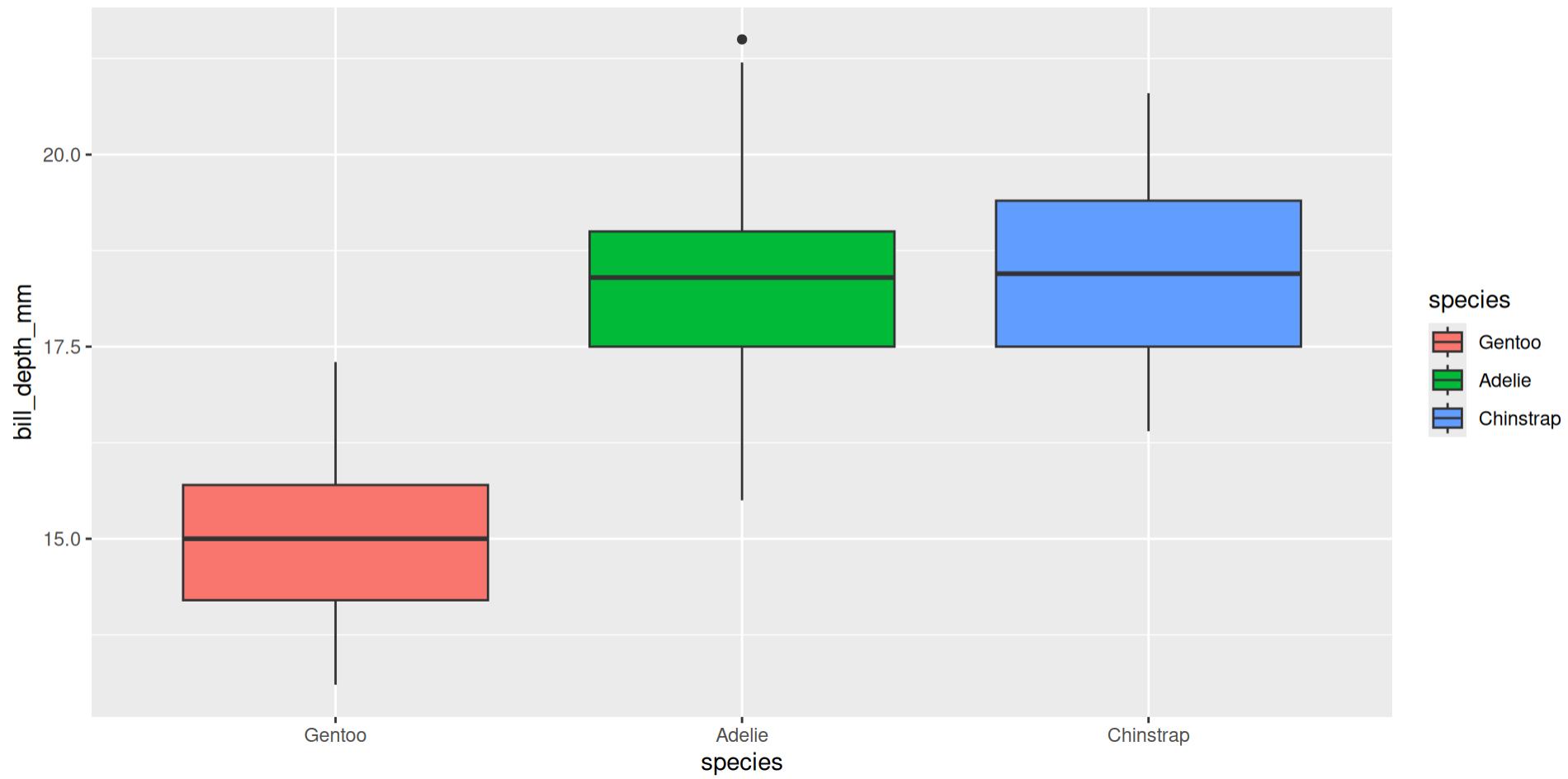
```
1 penguins |>
2   mutate(
3     year_factor = factor(year),
4     species = fct_infreq(species),
5     species = fct_rev(species)
6   ) |>
7   ggplot(aes(x = species, fill = year_factor)) +
8   geom_bar()
```



```
1 penguins |>
2   ggplot(aes(x = species, y = bill_depth_mm, fill = species)) +
3     geom_boxplot()
```



```
1 penguins |>
2   mutate(
3     species = fct_reorder(species, bill_depth_mm)
4   ) |>
5   ggplot(aes(x = species, y = bill_depth_mm, fill = species)) +
6   geom_boxplot()
```





rstatsBL - Data Science mit R

```
1 starwars |> count(species, sort = TRUE)
```

```
1 # A tibble: 38 × 2
2   species      n
3   <chr>     <int>
4   1 Human        35
5   2 Droid         6
6   3 <NA>          4
7   4 Gungan        3
8   5 Kaminoan      2
9   6 Mirialan      2
10  7 Twi'lek        2
11  8 Wookiee       2
12  9 Zabrak        2
13  10 Aleena        1
14 # i 28 more rows
```

```
1 starwars |>
2   mutate(species = fct_lump(species, n = 2)) |>
3   count(species)

1 # A tibble: 4 × 2
2   species     n
3   <fct>    <int>
4   1 Droid        6
5   2 Human       35
6   3 Other       42
7   4 <NA>         4
```

```
1 starwars |>
2   mutate(
3     species = fct_lump(species, n = 2),
4     species = fct_relevel(species, "Human")
5   ) |>
6   count(species)
```

```
1 # A tibble: 4 × 2
2   species     n
3   <fct>    <int>
4   1 Human      35
5   2 Droid       6
6   3 Other      42
7   4 <NA>        4
```

```
1 starwars |>
2   mutate(
3     species = fct_lump(species, n = 2),
4     species = fct_relevel(species, "Human"),
5     species = fct_recode(species, "Mensch" = "Human", "Anders" = "Other")
6   ) |>
7   count(species)
```

```
1 # A tibble: 4 × 2
2   species     n
3   <fct>    <int>
4   1 Mensch      35
5   2 Droid        6
6   3 Anders      42
7   4 <NA>         4
```

Praktikum: {forcats}

prak-05c-forcats-firmen.qmd

20 : 00

Break



10:00

Praktikum: `if_else()`, `case_when()`, `{forcats}`

prak-05d-cond-mutate-forcats.qmd

30 : 00

Danke! 

Slides created via [revealjs](#) and Quarto.

Access slides as [PDF](#).

All material is licensed under [Creative Commons Attribution Share Alike 4.0 International](#).