# Lab 2: Statistical inference & hypothesis testing

## Practice session covering topics discussed in Lecture 2

**M. Chiara Mimmi, Ph.D.** | Università degli Studi di Pavia

July 26, 2024

# GOAL OF TODAY'S PRACTICE SESSION

Consolidate understanding of inferential statistic, through R coding examples conducted on real biostatistics research data.

**Lecture 2: topics**

- **Purpose and foundations of inferential statistics**

- **Getting to know the "language" of hypothesis testing**

- **Hypothesis testing**
    - review examples

- **A closer look at testing assumptions**
    - more examples dealing with assumptions' violation

# R ENVIRONMENT SET UP & DATA

# Needed R Packages

- We will use functions from packages `base`, `utils`, and `stats` (pre-installed and pre-loaded)

- We will also use the packages below (specifying `package::function` for clarity).

```r
 1  # Load them for this R session
 2
 3  # General
 4  library(fs)            # file/directory interactions
 5  library(here)          # tools find your project's files, based on working directory
 6  library(janitor)       # tools for examining and cleaning data
 7  library(dplyr)         # {tidyverse} tools for manipulating and summarising tidy data
 8  library(forcats)       # {tidyverse} tool for handling factors
 9
10  # Statistics
11  library(BSDA)
12  library(rstatix)       # Pipe-Friendly Framework for Basic Statistical Tests
13  library(car)           # Companion to Applied Regression
14  library(multcomp)      # Simultaneous Inference in General Parametric Models
15
16  # Plotting
17  library(ggplot2)       # {tidyverse} tools for plotting
18  library(ggstatsplot)   # 'ggplot2' Based Plots with Statistical Details
19  library(ggpubr)        # 'ggplot2' Based Publication Ready Plots
20  library(patchwork)     # Functions for ""Grid" Graphics"composing" plots
21  library(viridis)       # Colorblind-Friendly Color Maps for R
22  library(ggthemes)      # Extra Themes, Scales and Geoms for 'ggplot2'
```

# Our dataset for today

For the most part, we will refer to a real clinical dataset (for which a *Creative Commons license* was granted) discussed in two articles (also open access) :

- Ahmad, T., Munir, A., Bhatti, S. H., Aftab, M., & Raza, M. A. (2017). **Survival analysis of heart failure patients: A case study**. PLOS ONE, 12(7), e0181001. https://doi.org/10.1371/journal.pone.0181001

- Chicco, D., & Jurman, G. (2020). **Machine learning can predict survival of patients with heart failure from serum creatinine and ejection fraction alone**. BMC Medical Informatics and Decision Making, 20(1), 16. https://doi.org/10.1186/s12911-020-1023-5

Here is the link to the dataset (or download from workshop website)

- From the UC Irvine Machine Learning Repository Heart Failure Clinical Records

# Importing from your project folder (previously downloaded file)

> **Tip**
>
> Make sure to match your own folder structure!

- The function `here` lets me specify the complete path of the destination folder

```
1   # Check my working directory location
2   # here::here()
3
4   # Use `here` in specifying all the subfolders AFTER the working directory
5   heart_failure <- read.csv(file = here::here("practice", "data_input", "02_datasets",
6                                       "heart_failure_clinical_records_dataset.csv"),
7                             header = TRUE, # 1st line is the name of the variables
8                             sep = ",", # which is the field separator character.
9                             na.strings = c("?","NA" ), # specific MISSING values
10                            row.names = NULL)
```

# INSPECTING THE "HEART FAILURE" DATASET

# What are the variables and their levels of measurement?

The data, containing the medical records of 299 heart failure patient, were collected at the Faisalabad Institute of Cardiology and at the Allied Hospital in Faisalabad (Punjab, Pakistan), during April–December 2015.

Table 1 from the second article (Chicco & Jurman, 2020, p. 3) offers a synthetic explanation of hte observed variables.

| Feature | Explanation | Measurement | Range |
|---|---|---|---|
| Age | Age of the patient | Years | [40, ..., 95] |
| Anaemia | Decrease of red blood cells or hemoglobin | Boolean | 0, 1 |
| High blood pressure | If a patient has hypertension | Boolean | 0, 1 |
| Creatinine phosphokinase (CPK) | Level of the CPK enzyme in the blood | mcg/L | [23, ..., 7861] |
| Diabetes | If the patient has diabetes | Boolean | 0, 1 |
| Ejection fraction | Percentage of blood leaving the heart at each contraction | Percentage | [14, ..., 80] |
| Sex | Woman or man | Binary | 0, 1 |
| Platelets | Platelets in the blood | kiloplatelets/mL | [25.01, ..., 850.00] |
| Serum creatinine | Level of creatinine in the blood | mg/dL | [0.50, ..., 9.40] |
| Serum sodium | Level of sodium in the blood | mEq/L | [114, ..., 148] |
| Smoking | If the patient smokes | Boolean | 0, 1 |
| Time | Follow-up period | Days | [4,...,285] |
| (target) death event | If the patient died during the follow-up period | Boolean | 0, 1 |

mcg/L: micrograms per liter. mL: microliter. mEq/L: milliequivalents per litre

# Look into the dataset just loaded in the R environment

Recall some base R functions from Lab 1

```r
1  # What variables are included in this dataset?
2  colnames(heart_failure)
```

```
 [1] "age"                     "anaemia"
 [3] "creatinine_phosphokinase" "diabetes"
 [5] "ejection_fraction"       "high_blood_pressure"
 [7] "platelets"               "serum_creatinine"
 [9] "serum_sodium"            "sex"
[11] "smoking"                 "time"
[13] "DEATH_EVENT"
```

```r
1  # How many observations & variables?
2  nrow(heart_failure)
```

```
[1] 299
```

```r
1  # How many rows & columns?
2  dim(heart_failure)
```

```
[1] 299  13
```

# Inspect the dataframe structure (base R)

```
1  # What does the dataframe look like?
2  str(heart_failure)
```

```
'data.frame':   299 obs. of  13 variables:
 $ age                     : num  75 55 65 50 65 90 75 60 65 80 ...
 $ anaemia                 : int  0 0 0 1 1 1 1 1 0 1 ...
 $ creatinine_phosphokinase: int  582 7861 146 111 160 47 246 315 157 123 ...
 $ diabetes                : int  0 0 0 0 1 0 0 1 0 0 ...
 $ ejection_fraction       : int  20 38 20 20 20 40 15 60 65 35 ...
 $ high_blood_pressure     : int  1 0 0 0 0 1 0 0 0 1 ...
 $ platelets               : num  265000 263358 162000 210000 327000 ...
 $ serum_creatinine        : num  1.9 1.1 1.3 1.9 2.7 2.1 1.2 1.1 1.5 9.4 ...
 $ serum_sodium            : int  130 136 129 137 116 132 137 131 138 133 ...
 $ sex                     : int  1 1 1 1 0 1 1 1 0 1 ...
 $ smoking                 : int  0 0 1 0 0 1 0 1 0 1 ...
 $ time                    : int  4 6 7 7 8 8 10 10 10 10 ...
 $ DEATH_EVENT             : int  1 1 1 1 1 1 1 1 1 1 ...
```

# Inspect the dataframe structure (`skimr`)

Remember the `skimr` function `skim`?

```
1  # some variables
2  heart_failure %>% skimr::skim( age, DEATH_EVENT )
3
4  # the whole dataframe
5  heart_failure %>% skimr::skim()
```

> **⊳_ You try...**
>
> Run `skimr::skim()` on your own either on the whole dataset or on any specific variable

- notice there are no (missing values) NAs in any of the variables

# Recode some variables for later ease of analysis

I may need some variables coded as `factor` (e.g. categorical variables for plotting), and, while I am at it, I can add clearer labels for the variables' levels. Here, we are:

- using tidyverse packages `dplyr` and `forcats`

- adding new (recoded) variables called *"oldname_f"*

```r
heart_failure <-heart_failure %>%
  dplyr::mutate(DEATH_EVENT_f = as.factor(DEATH_EVENT) %>%
                  forcats::fct_recode("died" = "1", "survived" = "0")) %>%
  dplyr::mutate(sex_f = as.factor(sex) %>%
                  forcats::fct_recode("male" = "1", "female" = "0"))

# check
table(heart_failure$DEATH_EVENT_f)
```

```
survived    died
     203      96
```

```r
table(heart_failure$sex_f)
```

```
female    male
   105     194
```

# Some more dummy variables recoded as factor

[Mostly for illustration: it's totally fine (if not preferable) to keep these as binary [0,1] variables]

- It's worth learning the useful function `dplyr::across`[1], which allows to iteratively transform several columns at once!

```r
1   # Recode as factor with levels "yes" (= 1), "no" (= 0)
2   fct_cols = c("anaemia", "diabetes", "high_blood_pressure", "smoking" )
3
4   heart_failure <- heart_failure  %>%
5     ## ---- 1st create new cols as "factor versions" of old cols
6     dplyr::mutate(
7       # let's introduce `across` function
8       dplyr::across(
9         # Columns to transform
10        .cols = all_of(fct_cols),
11        # Functions to apply to each col
12        .fns =  ~as.factor (.x),
13        # new name to apply where "{.col}" stands for the selected column
14        .names = "{.col}_f")) %>%
15    ## ---- 2nd create new cols as "factor versions" of old cols
16    dplyr::mutate(
17      dplyr::across(
18        # Columns to transform 2 conditions
19        .cols = ends_with("_f") & !matches(c( "DEATH_EVENT_f", "sex_f" )) ,
20        # Functions to apply to each col(different syntax)
21        .fns = ~forcats::fct_recode(.x,  yes = "1", no = "0" )))
```

1. This is a bit more advanced, but it will save a lot of typing in some situations...

# (Small digression on `dplyr::across`)

Notice how `dplyr::across(.cols = ..., .fns = ..., .names = ...)` has these arguments:

1. `.cols =` to select the columns which we want to transform (i.e. `fct_cols`)
   - with help from `tidyselect` functions: `all_of`, `ends_with`, and `matches`
2. `.fns = ~function(.x)` to specify the `function`
   - where `~function(.x)` uses the "anonymous function" syntax of the `tidyverse`
   - and `.x` inside the function is a "stand in" for *each of the columns selected*
3. [*optional*] `.names =` to name the new cols created using `{.col}` in place of each of the transformed columns

```r
1   ## ---- 1st create new cols as "factor versions" of old cols
2   heart_failure <- heart_failure  %>%
3     dplyr::mutate(
4       dplyr::across(
5         .cols = all_of(fct_cols),
6         .fns = ~as.factor (.x),
7         # (optional)
8         .names = "{.col}_f")) %>%
9   ## ---- 2nd create new cols as "factor versions" of old cols
10    dplyr::mutate(
11      dplyr::across(
12        .cols = ends_with("_f") & !matches(c( "DEATH_EVENT_f", "sex_f" )) ,
13        .fns =  ~forcats::fct_recode(.x,  yes = "1", no = "0" )))
```

# VISUAL DATA EXPLORATION FOR THE "HEART FAILURE"

**CONTINUOUS VARIABLES**

# Why is visual exploration important?

- Gaining insight on the variables (range, outliers, missing data)

- Preliminary check of assumptions for parametric hypothesis testing:

  - normally distributed outcome variables?

  - homogeneity of variance across groups?

Let's explore the **Heart failure dataset** with some data visualization…

- Following the referenced articles (which were mostly interested in predict mortality based on patients' characteristics), we will take the categorical, binary variable DEATH_EVENT_f as our main criterion to split the sample (into *survived* and *dead* patients) to explore any significant difference between groups in terms of means of known quantitative features.

- We will look at both:

  - **continuous variables** in the dataset (with the Probability Density Function (PDF))

  - **discrete variables** in the dataset (with the Probability Mass Function (PMF))

# Age

Introducing the handy R package `patchwork` which lets us compose different plots in a very simple and intuitive way

- (check it out with `??patchwork`)

```r
1   age <-ggplot(heart_failure,aes(x = age ))+
2     geom_histogram(binwidth = 5, color = "white", fill = "grey",alpha = 0.5)+
3     geom_vline(aes(xintercept = mean(age)), color = "#4c4c4c")+
4     theme_fivethirtyeight()+
5     labs(title = "Age Distribution" )+
6     scale_x_continuous(breaks = seq(40,100,5))
7
8   age2 <-ggplot(heart_failure, aes(x = age, fill = DEATH_EVENT_f))+
9     geom_histogram(binwidth = 5, position = "identity",alpha = 0.5,color = "white")+
10    geom_vline(aes(xintercept = mean(age[DEATH_EVENT == 0])), color = "#4c4c4c")+
11    geom_vline(aes(xintercept = mean(age[DEATH_EVENT==1])), color = "#d8717b")+
12    theme_fivethirtyeight()+
13    scale_fill_manual(values = c("#999999", "#d8717b"))+
14    labs(title =  "Age Distribution by group (Death Event)")+
15    scale_x_continuous(breaks = seq(40,100,5))
16
17  # patchwork
18  library(patchwork)
19  age + age2 + plot_layout(ncol = 1)
```

# Age

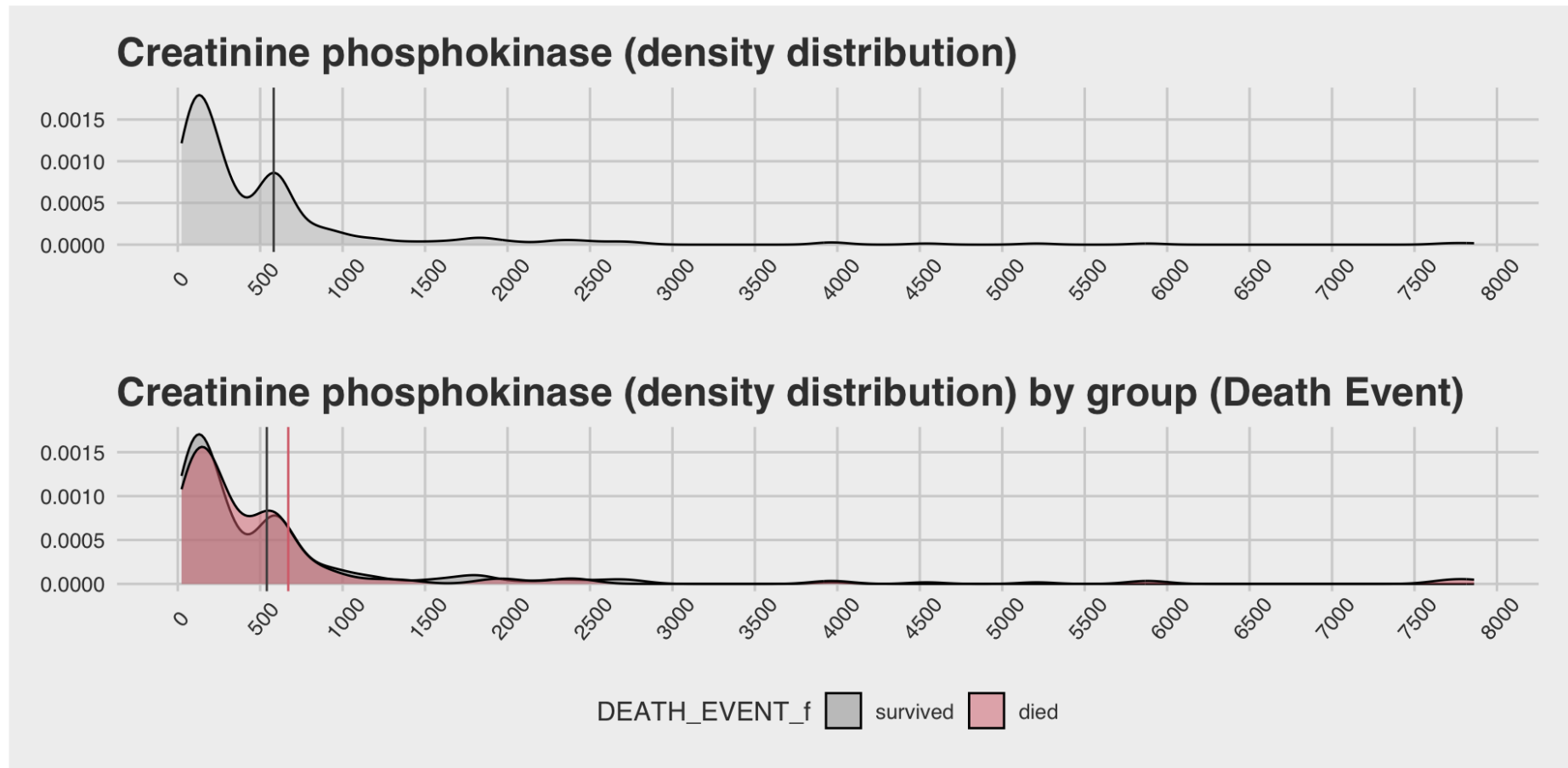As the age increases, the incidence of death event seems to increase

# Creatinine Phosphokinase (CPK)

```
 1  cpk <- ggplot(heart_failure,aes(x = creatinine_phosphokinase))+
 2    geom_density(fill = "gray", alpha = 0.5)+
 3    scale_x_continuous(breaks = seq(0,8000, 500))+
 4    geom_vline(aes(xintercept = mean(creatinine_phosphokinase)), color = "#4c4c4c")+
 5    theme_fivethirtyeight()+
 6    theme(axis.text.x = element_text(angle=50, vjust=0.75))+
 7    labs(title = "Creatinine phosphokinase (density distribution)" )+
 8    theme(plot.caption = element_text(hjust = 0.5, face = "italic"))
 9
10  cpk2 <- ggplot(heart_failure,aes(x = creatinine_phosphokinase,fill = DEATH_EVENT_f))+
11    geom_density(alpha = 0.5)+theme_fivethirtyeight()+
12    scale_fill_manual(values = c("#999999", "#d8717b"))+
13    scale_x_continuous(breaks = seq(0,8000, 500))+
14    geom_vline(aes(xintercept = mean(creatinine_phosphokinase[DEATH_EVENT == 0])),
15              color = "#4c4c4c")+
16    geom_vline(aes(xintercept = mean(creatinine_phosphokinase[DEATH_EVENT==1])),
17              color = "#d8717b")+
18    theme_fivethirtyeight()+
19    theme(axis.text.x = element_text(angle=50, vjust=0.75))+
20    labs(title =  "Creatinine phosphokinase (density distribution) by group (Death Event)")
21
22  cpk + cpk2 + plot_layout(ncol = 1)
```

# Creatinine Phosphokinase (CPK)

This definitely doesn't look like a normal distribution!
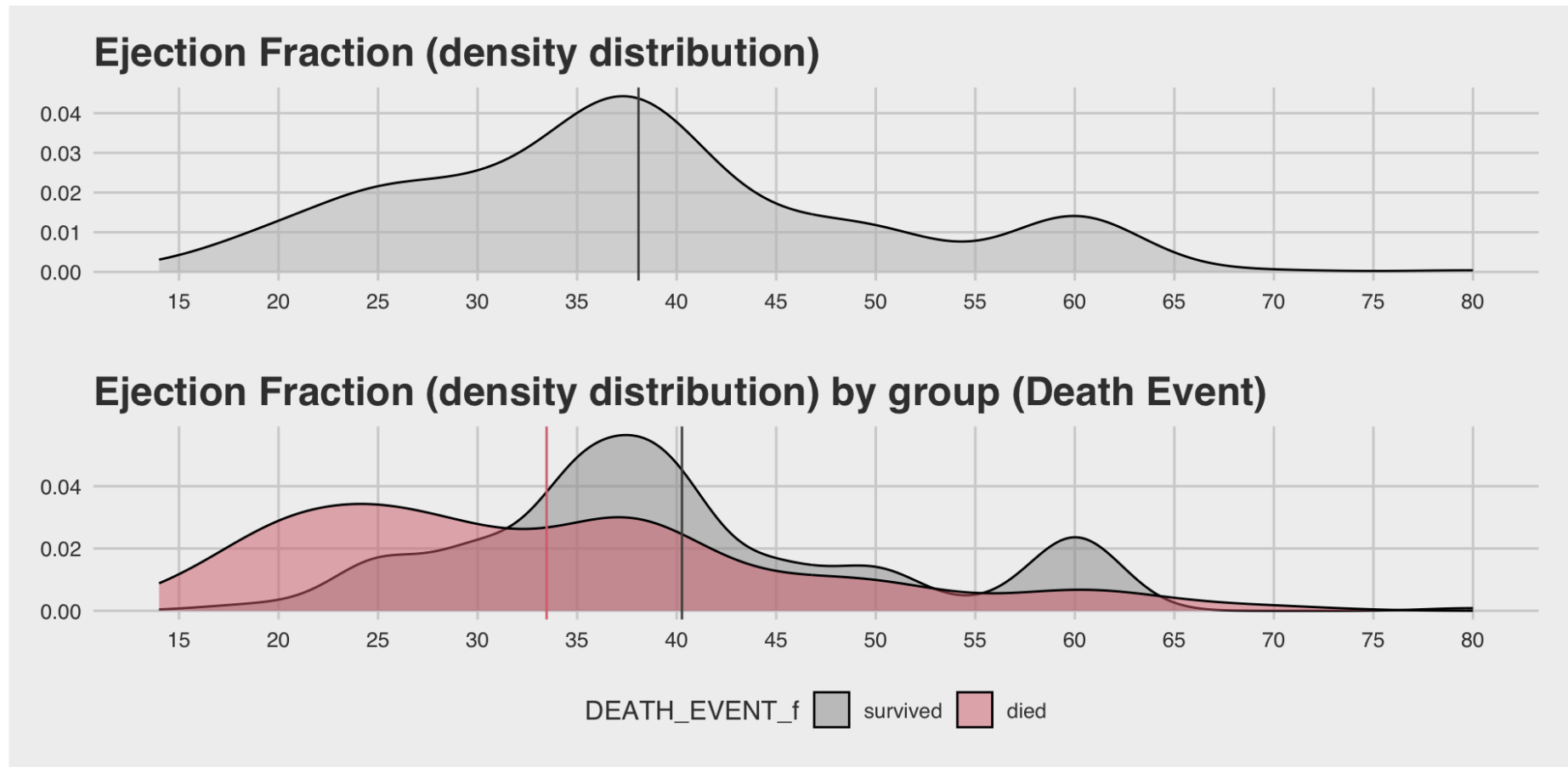
# Ejection Fraction

```r
1   ejf <- ggplot(heart_failure,aes(x = ejection_fraction))+
2     geom_density(fill = "gray", alpha = 0.5)+
3     scale_x_continuous(breaks = seq(0,100, 5))+
4     geom_vline(aes(xintercept = mean(ejection_fraction)), color = "#4c4c4c")+
5     theme_fivethirtyeight()+
6     labs(title = "Ejection Fraction (density distribution)" )+
7     theme(plot.caption = element_text(hjust = 0.5, face = "italic"))
8
9   ejf2 <- ggplot(heart_failure,aes(x = ejection_fraction,fill = DEATH_EVENT_f))+
10    geom_density(alpha = 0.5)+theme_fivethirtyeight()+
11    scale_x_continuous(breaks = seq(0,100, 5))+
12    scale_fill_manual(values = c("#999999", "#d8717b"))+
13    geom_vline(aes(xintercept = mean(ejection_fraction[DEATH_EVENT == 0])),
14              color = "#4c4c4c")+
15    geom_vline(aes(xintercept = mean(ejection_fraction[DEATH_EVENT==1])),
16              color = "#d8717b")+
17    labs(title =  "Ejection Fraction (density distribution) by group (Death Event)")+
18    theme_fivethirtyeight()
19
20  ejf + ejf2 + plot_layout(ncol = 1)
```

# Ejection Fraction

This also doesn't look like a normal distribution… and there is a remarkable change in the *probability density function* (PDF) shape when we introduce the grouping variable
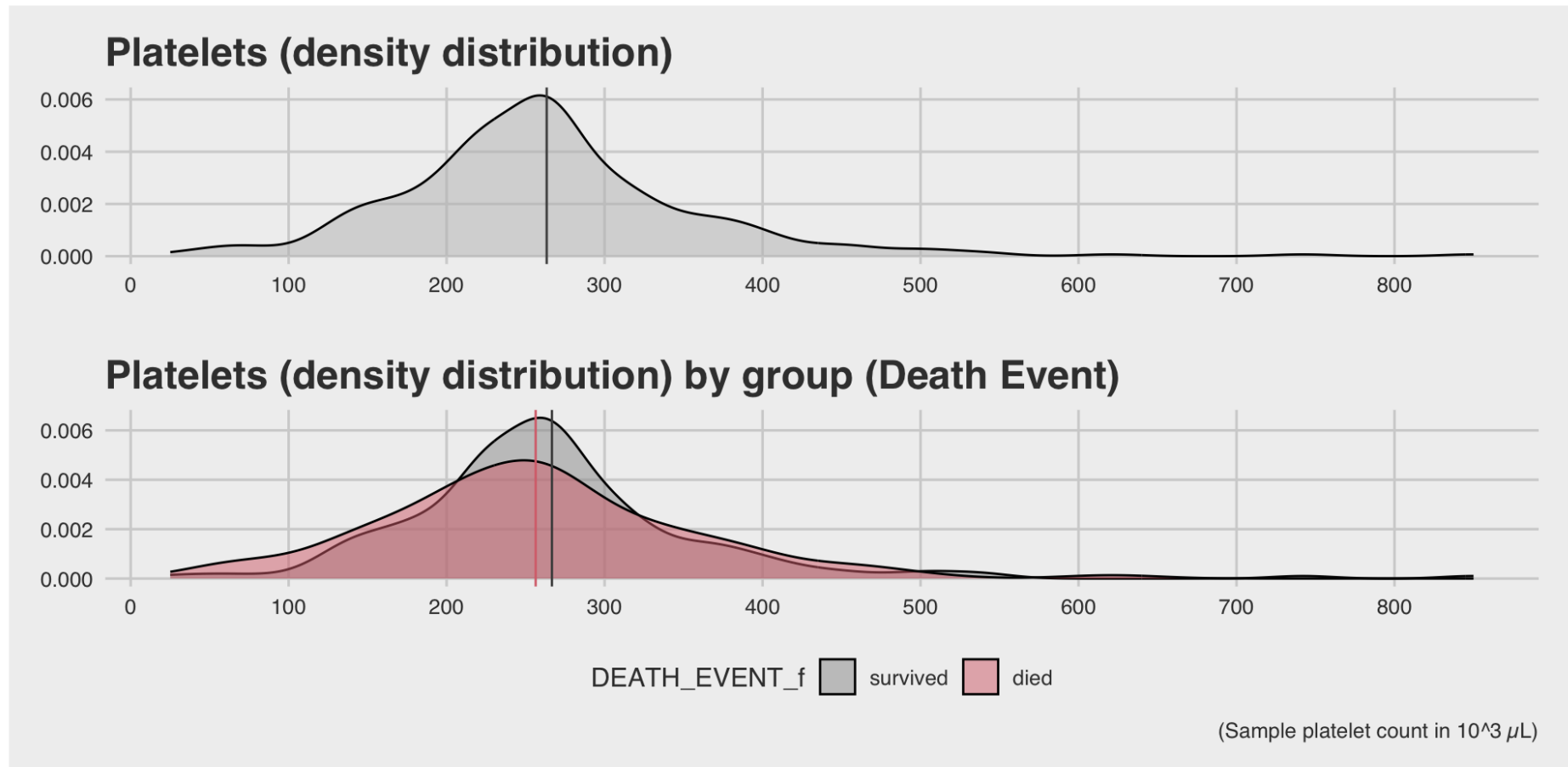
# Platelets

```r
1   # normalize the var for readability
2   heart_failure  <-  heart_failure %>%  dplyr::mutate(plat_norm = platelets/1000)
3
4   plat <- ggplot(heart_failure,aes(x = plat_norm))+
5     geom_density(fill = "gray", alpha = 0.5)+
6     scale_x_continuous(breaks = seq(0,800, 100))+
7     geom_vline(aes(xintercept = mean(plat_norm)), color = "#4c4c4c")+
8     theme_fivethirtyeight()    +
9     labs(title =  "Platelets (density distribution)",
10         y = "Density", x = "Sample platelet count (in 10^3 µL)")
11
12  plat2 <- ggplot(heart_failure,aes(x = plat_norm,fill = DEATH_EVENT_f))+
13    geom_density(alpha = 0.5)+theme_fivethirtyeight()+
14    scale_x_continuous(breaks = seq(0,800, 100))+
15    scale_fill_manual(values = c("#999999", "#d8717b"))+
16    geom_vline(aes(xintercept = mean(plat_norm[DEATH_EVENT == 0])),
17            color = "#4c4c4c")+
18    geom_vline(aes(xintercept = mean(plat_norm[DEATH_EVENT==1])),
19            color = "#d8717b")+
20    theme_fivethirtyeight()   +
21    labs(title =  "Platelets (density distribution) by group (Death Event)",
22        caption = "(Sample platelet count in 10^3 µL)")
23
24  plat + plat2 + plot_layout(ncol = 1)
```

# Platelets

Here the probability distributions resemble a Normal one and we observe more uniformity in the mean/variance across the 2 groups
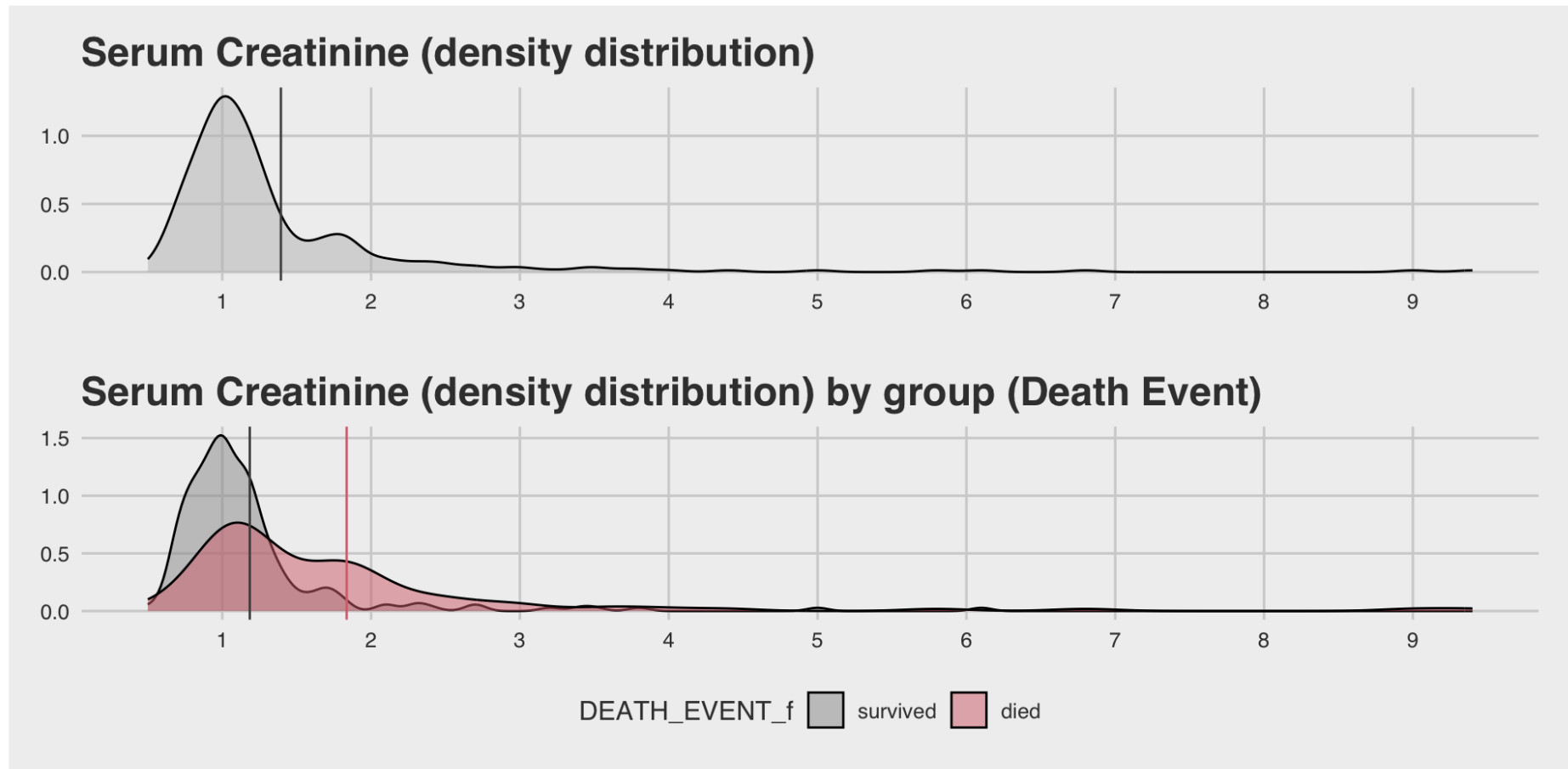
# Serum Creatinine

```r
1   ser_cr <- ggplot(heart_failure,aes(x = serum_creatinine))+
2     geom_density(fill = "gray", alpha = 0.5)+
3     scale_x_continuous(breaks = seq(0,10, 1))+
4     geom_vline(aes(xintercept = mean(serum_creatinine)), color = "#4c4c4c")+
5     theme_fivethirtyeight()+
6     labs(title = "Serum Creatinine (density distribution)" )+
7     theme(plot.caption = element_text(hjust = 0.5, face = "italic"))
8
9   ser_cr2 <- ggplot(heart_failure,aes(x = serum_creatinine,fill = DEATH_EVENT_f))+
10    geom_density(alpha = 0.5)+theme_fivethirtyeight()+
11    scale_x_continuous(breaks = seq(0,10, 1))+
12    scale_fill_manual(values = c("#999999", "#d8717b"))+
13    geom_vline(aes(xintercept = mean(serum_creatinine[DEATH_EVENT == 0])),
14              color = "#4c4c4c")+
15    geom_vline(aes(xintercept = mean(serum_creatinine[DEATH_EVENT==1])),
16              color = "#d8717b")+
17    labs(title =  "Serum Creatinine (density distribution) by group (Death Event)")+
18    theme_fivethirtyeight()
19
20  ser_cr + ser_cr2 + plot_layout(ncol = 1)
```

# Serum Creatinine

Another continuous random variable with a non-normal distribution (long right tails) and a seemingly important difference in variance between the groups.
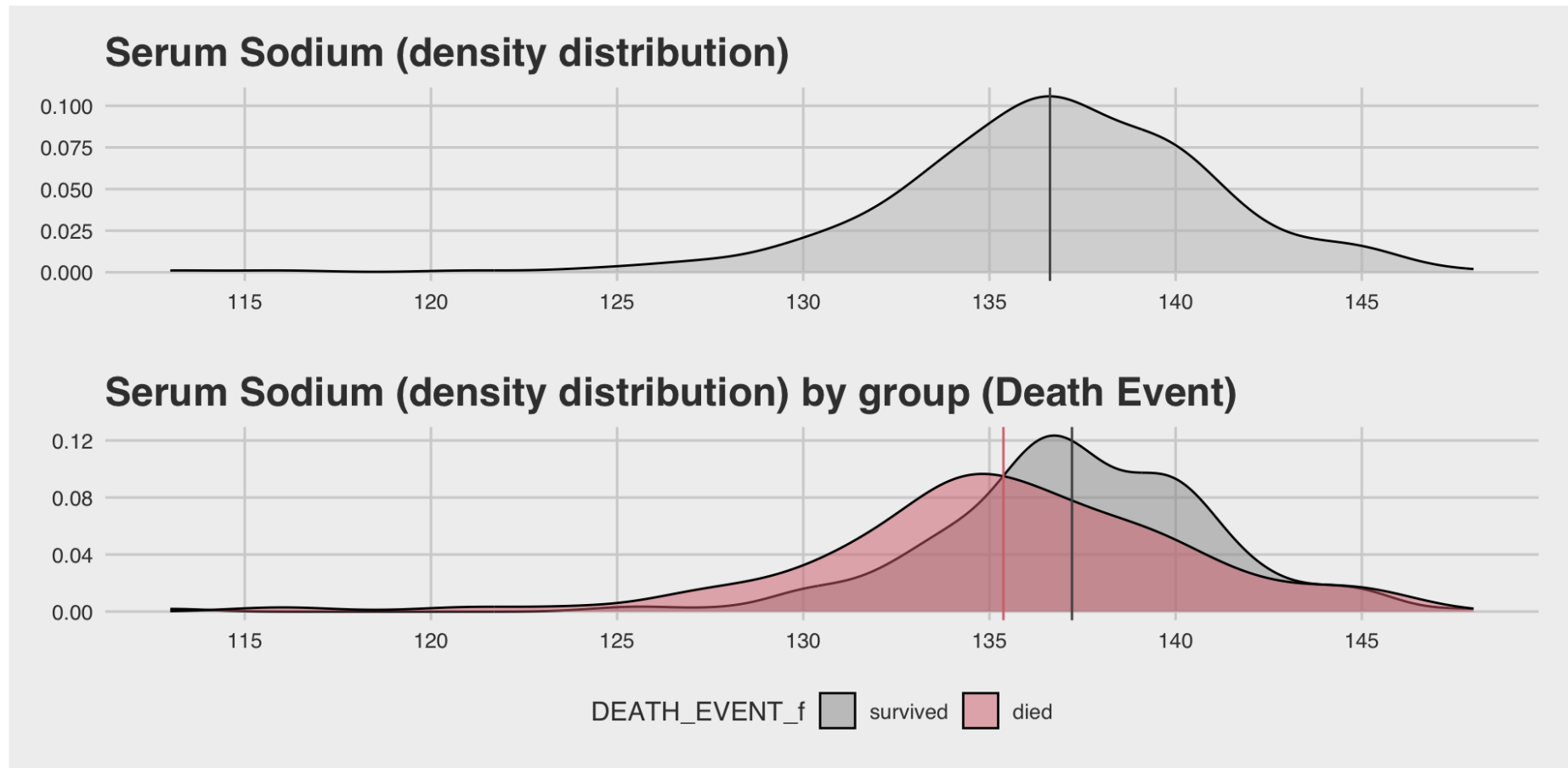
# Serum Sodium

```r
1   ser_sod <- ggplot(heart_failure,aes(x = serum_sodium))+
2     geom_density(fill = "gray", alpha = 0.5)+
3     scale_x_continuous(breaks = seq(0,150, 5))+
4     geom_vline(aes(xintercept = mean(serum_sodium)), color = "#4c4c4c")+
5     theme_fivethirtyeight()+
6     labs(title = "Serum Sodium (density distribution)" )
7
8   ser_sod2 <- ggplot(heart_failure,aes(x = serum_sodium,fill = DEATH_EVENT_f))+
9     geom_density(alpha = 0.5)+
10    scale_x_continuous(breaks = seq(0,150, 5))+
11    scale_fill_manual(values = c("#999999", "#d8717b"))+
12    geom_vline(aes(xintercept = mean(serum_sodium[DEATH_EVENT == 0])),
13              color = "#4c4c4c")+
14    geom_vline(aes(xintercept = mean(serum_sodium[DEATH_EVENT==1])),
15              color = "#d8717b")+
16    theme_fivethirtyeight()+
17    labs(title =  "Serum Sodium (density distribution) by group (Death Event)")+
18    theme_fivethirtyeight()
19
20  ser_sod + ser_sod2 + plot_layout(ncol = 1)
```

# Serum Sodium

Same as above, except for the long left tails...

# VISUAL DATA EXPLORATION FOR THE "HEART FAILURE"
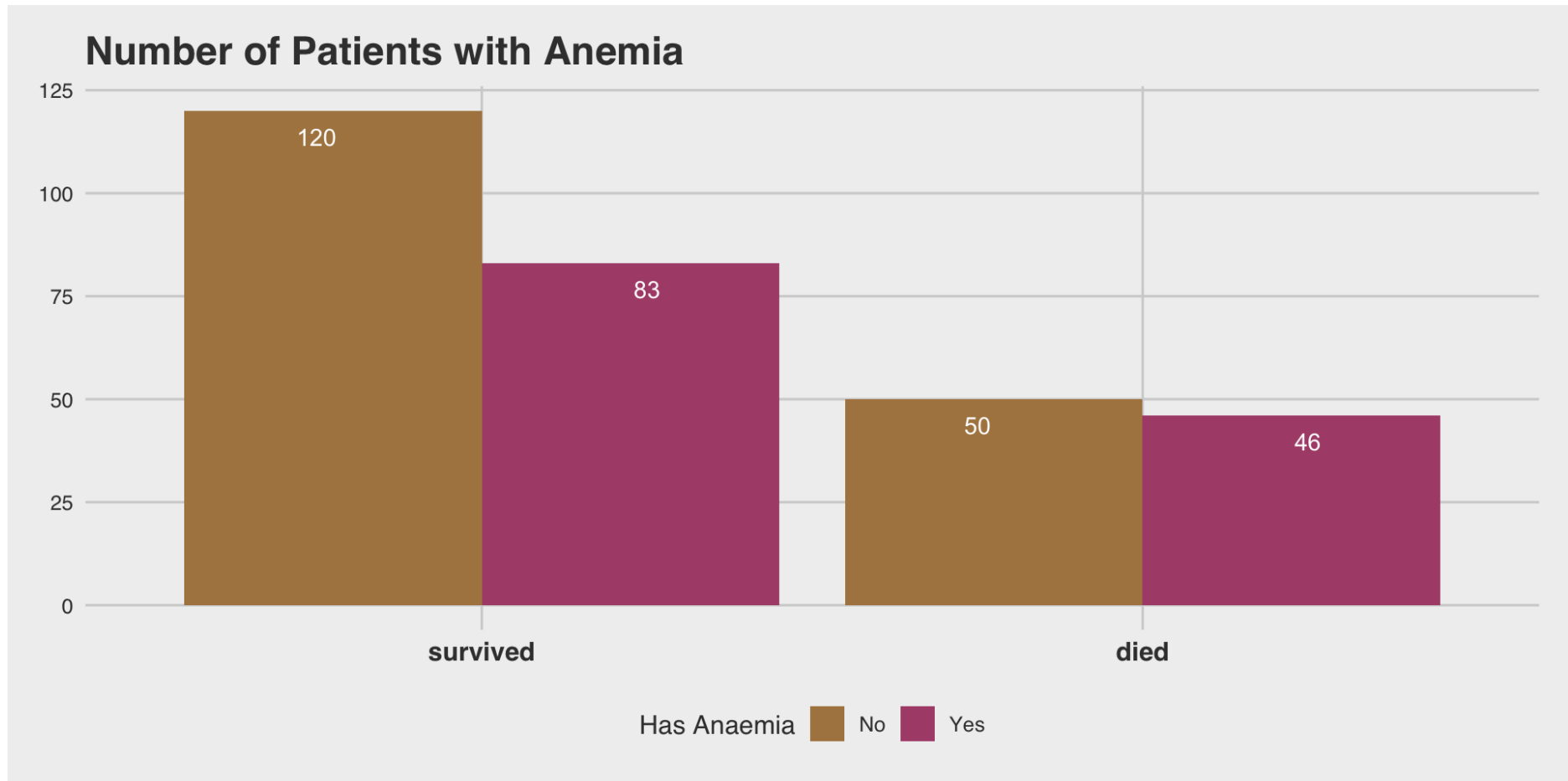
**DISCRETE VARIABLES**

# Anaemia

```r
1  anem <- ggplot(heart_failure, aes(x = forcats::fct_infreq(DEATH_EVENT_f ),
2                                     fill = anaemia_f ))+
3    geom_bar(position = "dodge")+
4    ## add count labels
5    geom_text(stat = "count", aes(label = ..count..),
6                ## make labels suit the dodged bars
7                position=position_dodge(width = 1 ),
8                hjust=0.5, vjust=2,color = "white") +
9    theme_fivethirtyeight() +
10   #scale_x_discrete(labels  = c("Death Event:No","Death Event:Yes"))+
11   scale_fill_manual(values = c("#af854f", "#af4f78"),
12                     name = "Has Anaemia",
13                     labels = c("No","Yes"))+
14   labs(title = "Number of Patients with Anemia") +
15   theme(#axis.text.x = element_text(angle=50, vjust=0.75),
16     axis.text.x = element_text(size=12,face="bold"))
17
18 anem
```
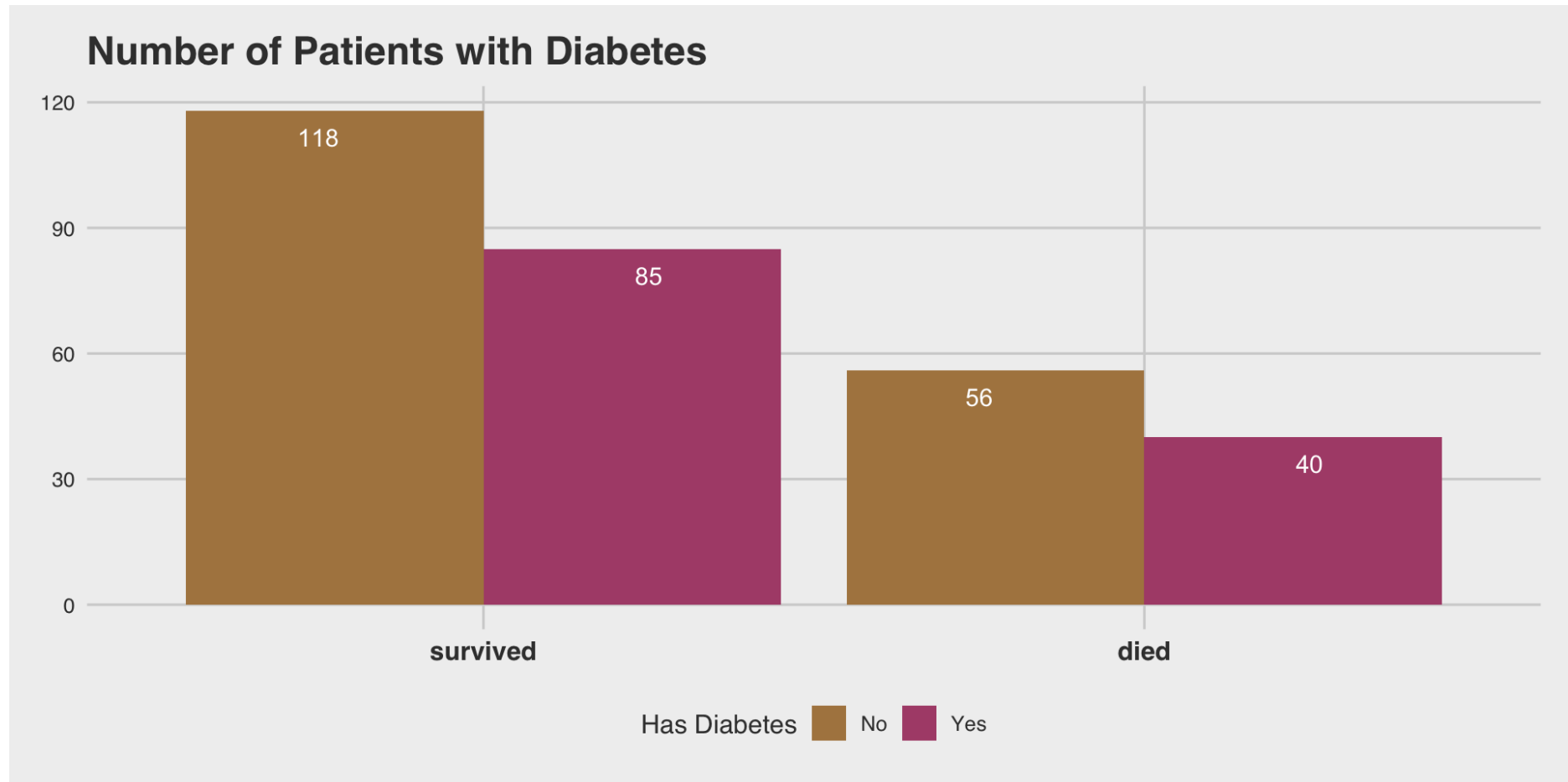
# Anaemia

There seems to be a greater incidence of anaemia in group 'died'



**Number of Patients with Anemia**

# Diabetes

```r
1  diab <- ggplot(heart_failure,
2                 aes(x = forcats::fct_infreq(DEATH_EVENT_f ), fill = diabetes_f ))+
3    geom_bar(position = "dodge")+
4    ## add count labels
5    geom_text(stat = "count", aes(label = ..count..),
6              ## make labels suit the dodged bars
7              position=position_dodge(width = 1 ),
8              hjust=0.5, vjust=2,color = "white", size =4) +
9    theme_fivethirtyeight() +
10   #scale_x_discrete(labels  = c("Death Event:No","Death Event:Yes"))+
11   scale_fill_manual(values = c("#af854f", "#af4f78"),
12                     name = "Has Diabetes",
13                     labels = c("No","Yes"))+
14   labs(title = "Number of Patients with Diabetes") +
15   theme(#axis.text.x = element_text(angle=50, vjust=0.75),
16     axis.text.x = element_text(size=12,face="bold"))
17
18 diab
```
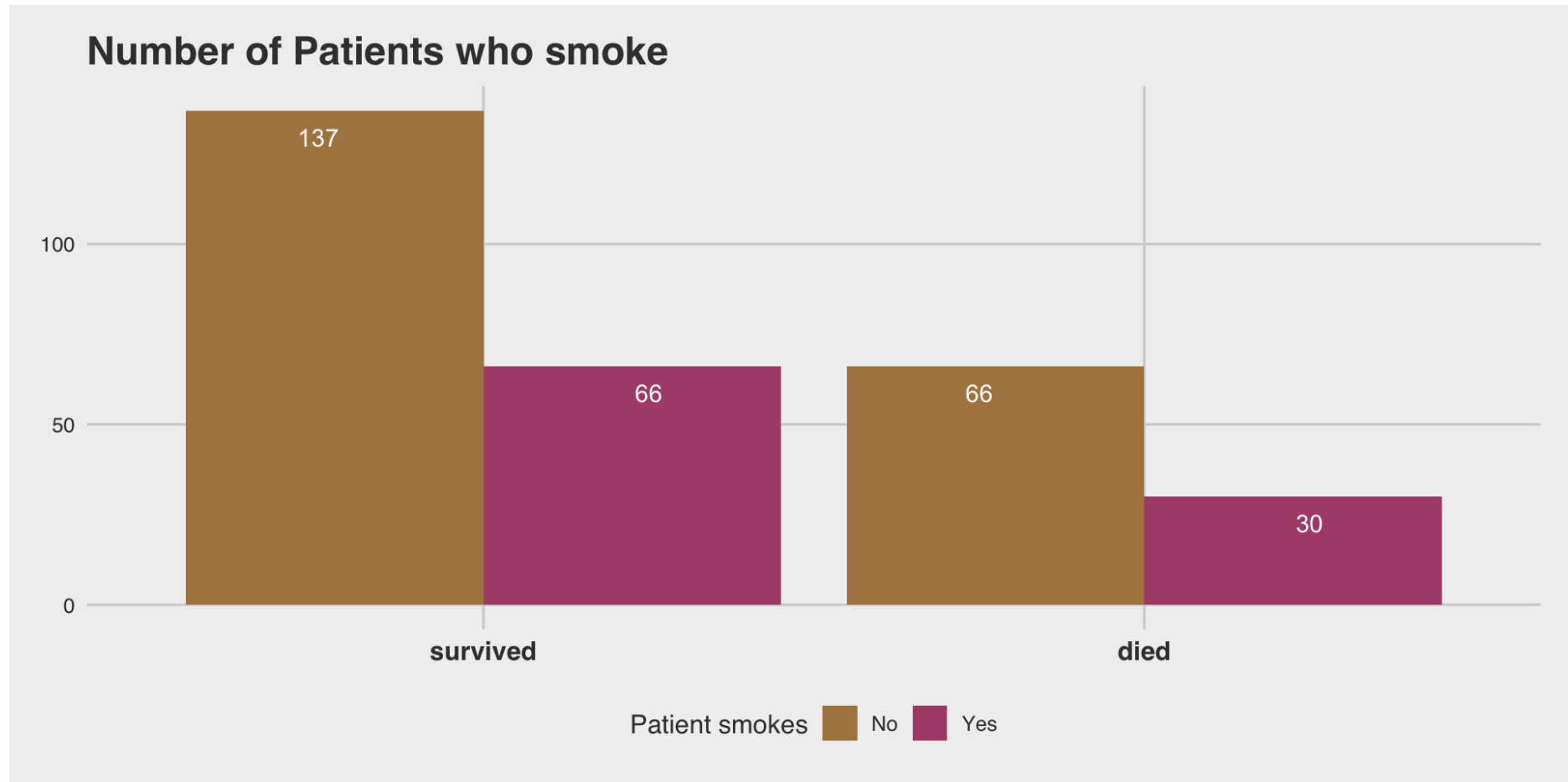
# Diabetes



Number of Patients with Diabetes

# Smoking

```
1   smok <- ggplot(heart_failure, aes(x = forcats::fct_infreq(DEATH_EVENT_f ),
2                                      fill = smoking_f ))+
3     geom_bar(position = "dodge")+
4     ## add count labels
5     geom_text(stat = "count", aes(label = ..count..),
6               ## make labels suit the dodged bars
7               position=position_dodge(width = 1 ),
8               hjust=0.5, vjust=2,color = "white", size =4) +
9     theme_fivethirtyeight() +
10    #scale_x_discrete(labels  = c("Death Event:No","Death Event:Yes"))+
11    scale_fill_manual(values = c("#af854f", "#af4f78"),
12                      name = "Patient smokes",
13                      labels = c("No","Yes"))+
14    labs(title = "Number of Patients who smoke") +
15    theme(#axis.text.x = element_text(angle=50, vjust=0.75),
16      axis.text.x = element_text(size=12,face="bold"))
17
18  smok
```

# Smoking



**Number of Patients who smoke**

survived: No = 137, Yes = 66
died: No = 66, Yes = 30

Patient smokes — No — Yes
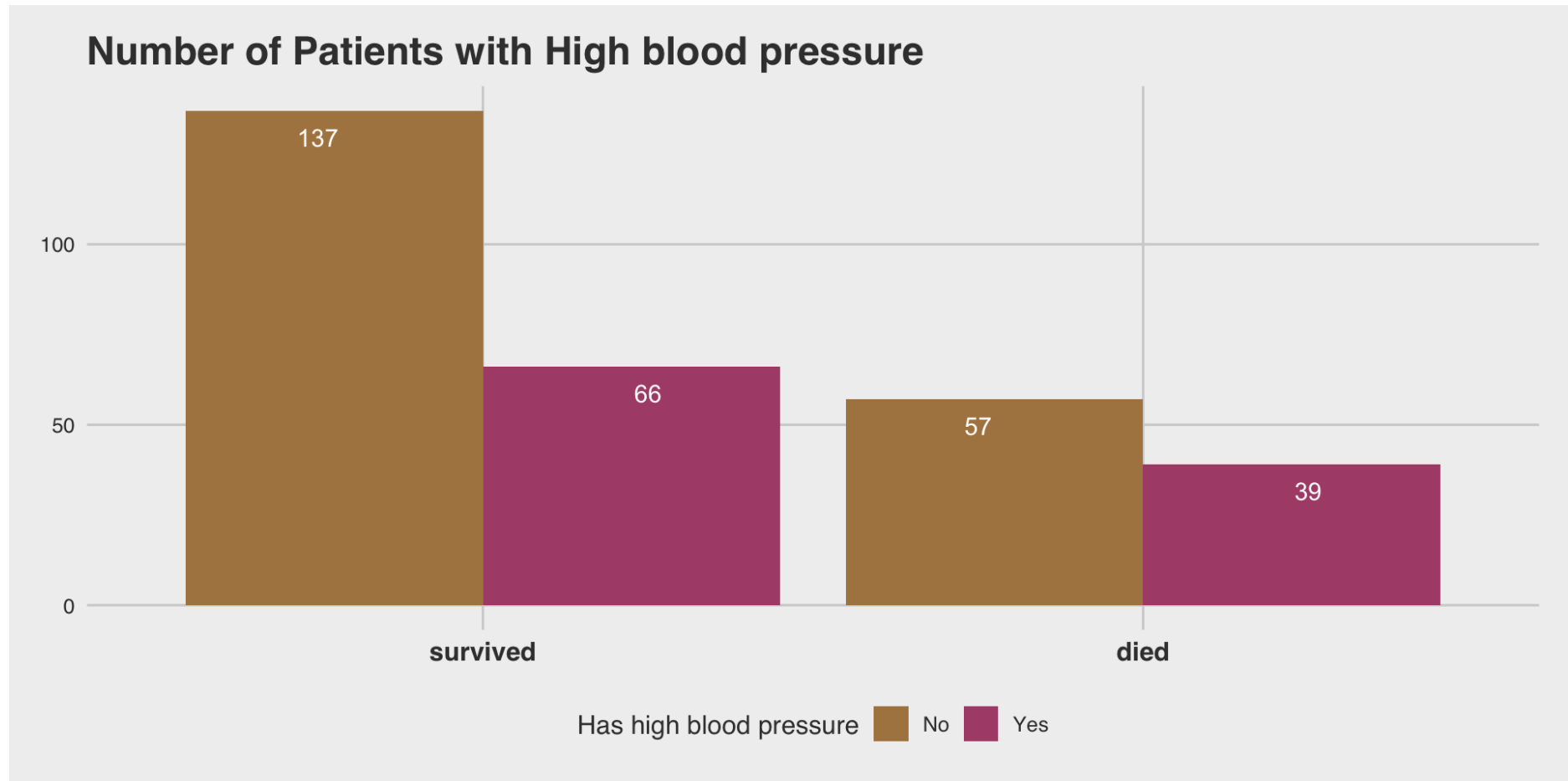
# High blood pressure

```
1   hbp <- ggplot(heart_failure, aes(x = forcats::fct_infreq(DEATH_EVENT_f ),
2                                     fill = high_blood_pressure_f ))+
3     geom_bar(position = "dodge")+
4       ## add count labels
5     geom_text(stat = "count", aes(label = ..count..),
6               ## make labels suit the dodged bars
7               position=position_dodge(width = 1 ),
8               hjust=0.5, vjust=2,color = "white", size =4) +
9     theme_fivethirtyeight() +
10    #scale_x_discrete(labels  = c("Death Event:No","Death Event:Yes"))+
11    scale_fill_manual(values = c("#af854f", "#af4f78"),
12                      name = "Has high blood pressure",
13                      labels = c("No","Yes"))+
14    labs(title = "Number of Patients with High blood pressure") +
15    theme(#axis.text.x = element_text(angle=50, vjust=0.75),
16      axis.text.x = element_text(size=12,face="bold"))
17
18  hbp
```

# High blood pressure

There is also a greater incidence of high blood pressure in group 'died'

# HYPOTHESIS TESTNG
# - some examples -

Let's continue to explore data from the **heart failure patients' dataset**, but this time using **hypothesis testing** as we learned in Lecture 2. We will do two types of test:

1. Comparing a sample **against a *hypothetical* general population**

2. Testing if mean variables' **differences between the two groups of patients** (those who survived after heart failure event and those who didn't) is statistically significant

# — EXAMPLE A —

**(1 sample | n > 30 | Z test)**

**Comparing sample mean to a hypothesized population mean (with Z test)**

Stating the above hypotheses more formally:

**What is the population Total Platelet Count (TPC) mean for all people who suffered of heart failure ($\mu_{HF}$)?**
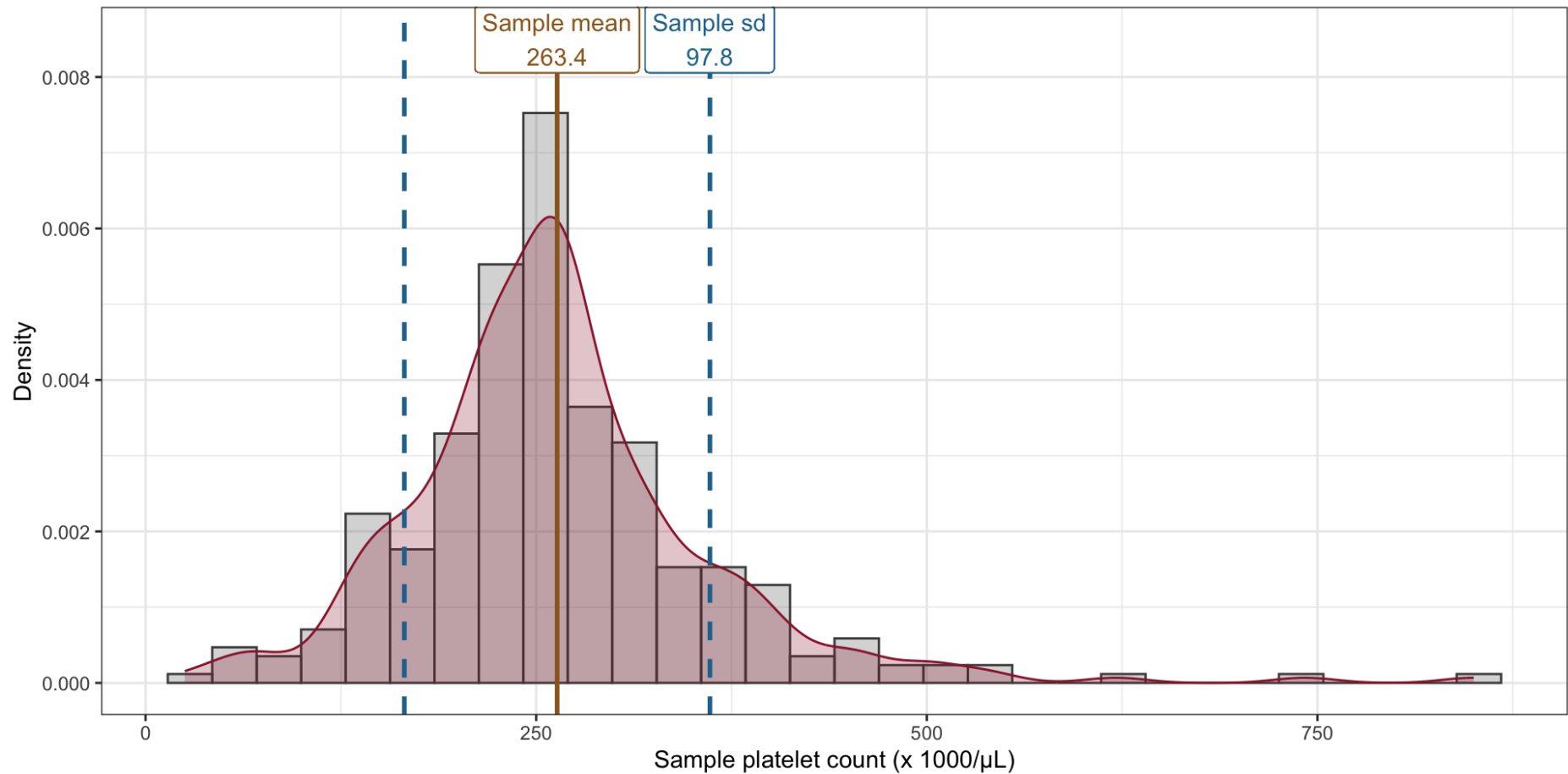
- $H_0$ : there is no difference in mean TPC between patients who suffered heart failure and the general population
  - $\mu_{HF}$ = 236 -> hypothesis of no effect or ("no difference")
- $H_a$ : there is a difference in mean TPC between patients who have suffered heart failure and the general population ("some effect"). This can be formalized as either:
  - $\mu_{HF}$ < 236 (one-sided test), or
  - $\mu_{HF}$ > 236 (one-sided test), or
  - $\mu_{HF}$ ≠ 236 (two-sided test)

# 1. Question: How does the mean platelets count in the patients' sample compare against a r

```r
1   # compute mean & sd for plot
2   mean_plat_p <- round(mean(heart_failure$plat_norm), digits = 1)
3   sd_plat_p <- round(sd(heart_failure$plat_norm), digits = 1)
4
5   heart_failure %>%
6     ggplot(aes(x = plat_norm))+
7     geom_histogram(aes(y = ..density..), bins=30, alpha=0.25, colour = "#4c4c4c") +
8     geom_density(colour ="#9b2339", alpha=0.25, fill = "#9b2339") +
9     # add mean vertical line
10    geom_vline(xintercept = mean_plat_p, na.rm = FALSE,size = 1,color= "#9b6723") +
11    # add also +/- 1sd
12    geom_vline(aes(xintercept = mean_plat_p + sd_plat_p),
13               color = "#23749b", size = 1, linetype = "dashed") +
14    geom_vline(aes(xintercept = mean_plat_p - sd_plat_p),
15               color = "#23749b", size = 1, linetype = "dashed") +
16    # add annotations with the mean value
17    geom_label(aes(x=mean_plat_p,  y=0.0085, label=paste0("Sample mean\n",mean_plat_p)),
18               color = "#9b6723") +
19    geom_label(aes(x=361,  y=0.0085, label=paste0("Sample sd\n",sd_plat_p)),
20               color = "#23749b") +
21    theme_bw() +  labs(y = "Density", x = "Sample platelet count (x 1000/µL)")
```

General population data taken from the literature (See Wongsaengsak, Dennis, Arevalo, Ball, & Nugent, 2019).

# 1. Question: How does the mean platelets count in the patients' sample compare against a r...

For a general population, the Total Platelet Count (TPL) has $\mu$=236 (1000 /µL) and $\sigma$= 59 (1000 /µL). Below is the sample distribution:

# 2.a Computation of the test statistic

In this case, we have:

- a large sample ($n > 100$)

- a known $\sigma^2$ (of the reference population)

- the observed sample mean $\bar{x}$ and sample sd $s$.

So we can compute:

$$Z_{calc} = \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}}$$

- ✍️ Let's do it "by hand" first to see the steps

```r
1  # General Population of reference
2  mu <- 236
3  sigma  <- 59
4  # Sample of HF patients
5  n <- 299
6  x_HF <- mean(heart_failure$plat_norm)      #    263.358
7  s_HF <- sd(heart_failure$plat_norm)        #    97.80424
8  # IF large sample & KNOWN pop variance
9  std_err_HF <- sigma /sqrt(n)               # 3.412058
10 z_calc_HF <-  (x_HF - mu) / std_err_HF     # 8.018043
```

# 2.b Computation of the p-value associated to the test statistic

To find the **p-value** associated with a z-score in R, we can use the `pnorm()` function, which uses the following syntax:

- `q`: The z-score

- `mean`: The mean of the normal distribution. Default is 0.

- `sd`: The standard deviation of the normal distribution. Default is 1.

- `lower.tail`:

  - If TRUE, the probability to the left of q in the normal distribution is returned

  - If FALSE, the probability to the right is returned. Default is TRUE.

```r
1  # Left-tailed test
2  p_value_l <- stats::pnorm(z_calc_HF, mean = 0, sd = 1, lower.tail = TRUE)
3  # Right-tailed test
4  p_value_r <- stats::pnorm(z_calc_HF, mean = 0, sd = 1,lower.tail = FALSE)
5  # Two-tailed test  (our case)
6  p_value_two <- 2*stats::pnorm(z_calc_HF, mean = 0, sd = 1, lower.tail = FALSE)
```

# 2.c Computation of the p-value associated to the test statistic

- 🧑‍💻 Let's see how this could be done using an R function `BSDA::z.test`

```r
1  z_test_summary <- BSDA::z.test(x = heart_failure$plat_norm,
2                    alternative='two.sided',
3                    mu=236,
4                    sigma.x=59,
5                    conf.level=.95)
6  z_test_summary
```

```
    One-sample z-Test

data:  heart_failure$plat_norm
z = 8.018, p-value = 1.074e-15
alternative hypothesis: true mean is not equal to 236
95 percent confidence interval:
 256.6705 270.0455
sample estimates:
mean of x
  263.358
```

Same results!

# 3. Results and interpretation

1. Based on the critical region, the calculated test statistic z_calc_HF = 8.0180 falls in the CRITICAL REGION (well beyond the critical point)

```
1  # given
2  z_critical  <- c(-1.96, +1.96) # (Z score corresponding to α  = 0.05)
3  # Check
4  z_calc_HF > z_critical
```

```
[1] TRUE TRUE
```

2. Based on the p-value, p_value_two = 1.07443e-15 is much much smaller than $\alpha$

```
1  # Check
2  p_value_two <   0.05
```

```
[1] TRUE
```

**DECISION**: we reject the Null Hypothesis (basically we conclude that it is extremely unlikely that the sample we drew could have occurred just by chance). So the test indicates that, indeed, there is a difference between heart failure patients and the general population in terms of average platelets count.

# — EXAMPLE B —

**(1 sample | n < 30 | t test)**

# Comparing sample mean to a hypothesized population mean (with t test)

Same question, but with a *smaller sample* to work on (this varies, but generally it means $n < 30$). Imagine the patients were only observed over a **follow-up period of 21 days**, and also let's assume we don't know the population's variance

Stating the hypothesis more formally:

**What is the population Total Platelet Count (TPC) mean for all people who suffered of heart failure ($\mu_{HF21d}$) in the past 21 days or less?**
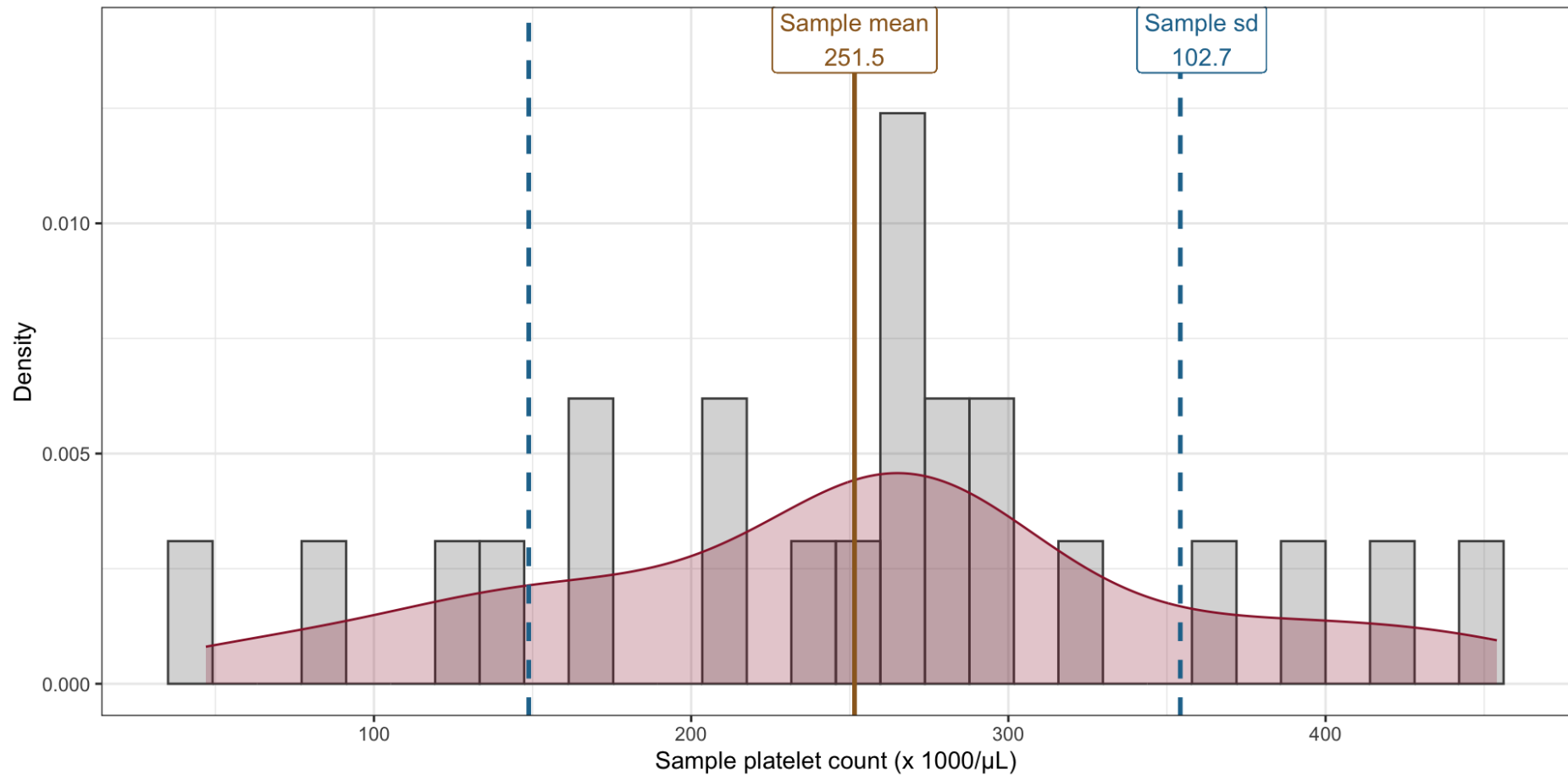
- $H_0$ : there is no difference in mean TPC between patients who suffered heart failure (visited in 21 days) and the general population

  - $\mu_{HF21d}$ = 236 -> hypothesis of no effect or ("no difference")

- $H_a$ : there is a difference in mean TPC between patients who have suffered heart failure and the general population ("some effect"). This can be formalized as:

  - $\mu_{HF21d} \neq$ 236 (two-sided test)

# 1. Question: How does the mean platelets count in the patients' sample compare against a r

```r
1   # normalize the var for readability
2   heart_21d  <-  heart_failure %>%  dplyr::mutate(plat_norm = platelets/1000) %>%
3     filter(time <= 21)                              # 23 obs
4   # compute mean & sd for plot
5   mean_plat_p <- round(mean(heart_21d$plat_norm), digits = 1)
6   sd_plat_p <- round(sd(heart_21d$plat_norm), digits = 1)
7
8   heart_21d %>%
9     ggplot(aes(x = plat_norm))+
10    geom_histogram(aes(y = ..density..), bins=30, alpha=0.25, colour = "#4c4c4c") +
11    geom_density(colour ="#9b2339", alpha=0.25, fill = "#9b2339") +
12    # add mean vertical line
13    geom_vline(xintercept = mean_plat_p, na.rm = FALSE,size = 1,color= "#9b6723") +
14    # add also +/- 1sd
15    geom_vline(aes(xintercept = mean_plat_p + sd_plat_p),
16              color = "#23749b", size = 1, linetype = "dashed") +
17    geom_vline(aes(xintercept = mean_plat_p - sd_plat_p),
18              color = "#23749b", size = 1, linetype = "dashed") +
19    # add annotations with the mean value
20    geom_label(aes(x=mean_plat_p,  y=0.014, label=paste0("Sample mean\n",mean_plat_p)),
21              color = "#9b6723") +
22    geom_label(aes(x=361,  y=0.014, label=paste0("Sample sd\n",sd_plat_p)),
23              color = "#23749b") +
24    theme_bw() +  labs(y = "Density", x = "Sample platelet count (x 1000/$\mu$L)")
```

# 1. Question: How does the mean platelets count in the patients' sample compare against a r

For a general population, the Total Platelet Count (TPL) has $\mu$=236 (1000 /µL) and $\sigma$= 59 (1000 /µL). Below is the smaller sample distribution:

# 2.a Picking the suitable test

In this case, we have:

- a "small" sample $n = 23$

- an unknown $\sigma^2$ (of the reference population)

- We obtained the sample mean $\bar{x}$ and sample sd $s$.

So we can compute:

$$t_{calc} = \frac{\bar{x} - \mu}{\frac{s_{\bar{x}}}{\sqrt{n-1}}}$$

# 2.b Computation of the test statistic

- Option 1: Let's compute the t test "by hand" ✍🏻

```r
1   # General Population of reference
2   mu_pop <- 236
3
4   # SAMPLE HF patients follow up less 21 days
5   heart_21d <- heart_failure %>% filter(time <= 21)
6
7   n_21d <- nrow(heart_21d)                          # 23
8   x_HF_21d <- mean(heart_21d$plat_norm)             # 251.5094
9   s_HF_21d <- sd(heart_21d$plat_norm)               # 102.7341
10  df_HF_21d <- n_21d-1                              # 22
11
12  # IF SMALL sample UNKNOWN sigma
13  std_err_HF_21d <- s_HF_21d /sqrt(n_21d -1)        # 21.90298
14  t_calc <-  (x_HF_21d - mu_pop) / std_err_HF_21d   # 0.7080951
```

- Option 2: Let's compute the t test with `stats::t.test` 🧑🏻‍💻

```r
1   t_stat_HF_21d_v2 <- stats::t.test(x = heart_21d$plat_norm,
2                                     mu = mu_pop,
3                                     alternative = "two.sided")
4   # extract t_calc from results df
5   t_calc_v2  <- t_stat_HF_21d_v2[["statistic"]][["t"]] # 0.7240093
```

There is a small difference in the `t_calc` ≠ `t_calc_v2` due to the fact that I use the Bessel's correction $n-1$ in the

# 2.c Computation of the p-value associated to the test statistic

- Option 1: "by hand" ✍🏻

To find the **p-value** associated with a t-score in R, we can use the `pt(q, df, lower.tail = TRUE)` function, which uses the following syntax:

- `q`: The t-score

- `df`: The degrees of freedom

- `lower.tail`:

    - TRUE to calculate the probability to the left of q which is called as left-tailed test

    - FALSE as right-tailed test.

```
1  # ---- Option 1
2  # -- Left-tailed test
3  #pt(t_stat_HF_21d, df_HF_21d, lower.tail = TRUE)
4
5  # -- Right-tailed test
6  #pt(t_stat_HF_21d, df_HF_21d, lower.tail = FALSE)
7
8  # -- Two-tailed test  (our case)
9  p_value_t_test <- 2*pt(t_calc, df_HF_21d, lower.tail = FALSE) # 0.4863214
```

- Option 2: from results of `stats::t.test` 🧑🏻

```
1  # ---- Option 2
2  # extract  p_value from results df
3  p_value_v2  <- t_stat_HF_21d_v2[["p.value"]] # 0.4766892
```

# 3. Results and interpretation

1. Based on the critical region, `t_calc ≃ 0.71` is smaller than the t critical value, i.e. it falls within the region of acceptance, so he null hypothesis is not rejected

```
1  #find two-tailed t critical values
2
3  t_crit_two <- qt(p=.05/2, df=22, lower.tail=FALSE)    # 2.073873
4  # Compare t score against t critical
5  t_calc > t_crit_two   # FALSE
```

```
[1] FALSE
```

2. Based on the p-value, `p_value ≃ 0.48` is larger than $\alpha$, i.e. the probability of observing a test statistic (assuming $H_0$ is true) is quite large

```
1  # Check
2  p_value_t_test <  0.05   # FALSE
```

```
[1] FALSE
```

**DECISION**: we FAIL to reject $H_0$. So the test indicates that there is not a statistically significant difference between heart failure patients visited within 21 days and the general population in terms of average platelets count.

> ⓘ **Note**
>
> What changed testing a sample with smaller n, instead of a large one?

# — EXAMPLE C —

**(2 samples | t test)**

# Comparing two independent sample means (t test)

This time, we investigate if there might be an actual difference in the Platelet Count means **between the patients who died and the patients who survived** heart failure.

Stating the above hypotheses more formally:

**Is there a statistically significant difference between the mean values of two groups?**

- $H_0$ : The two population means are equal

    - $\mu_1 = \mu_0 \iff \mu_1 - \mu_0 = 0$

- $H_a$ : There is a mean difference between the two groups in the population. Possible directional difference formulation (two-tailed, left-tailed, right-tailed)

    - $\mu_1 \neq \mu_0 \iff \mu_1 - \mu_0 \neq 0$ (the two population means are not equal)

    - $\mu_1 < \mu_0 \iff \mu_1 - \mu_0 < 0$ (population 1 mean is less than population 0 mean)

    - $\mu_1 > \mu_0 \iff \mu_1 - \mu_0 > 0$ (population 1 mean is greater than population 0 mean)
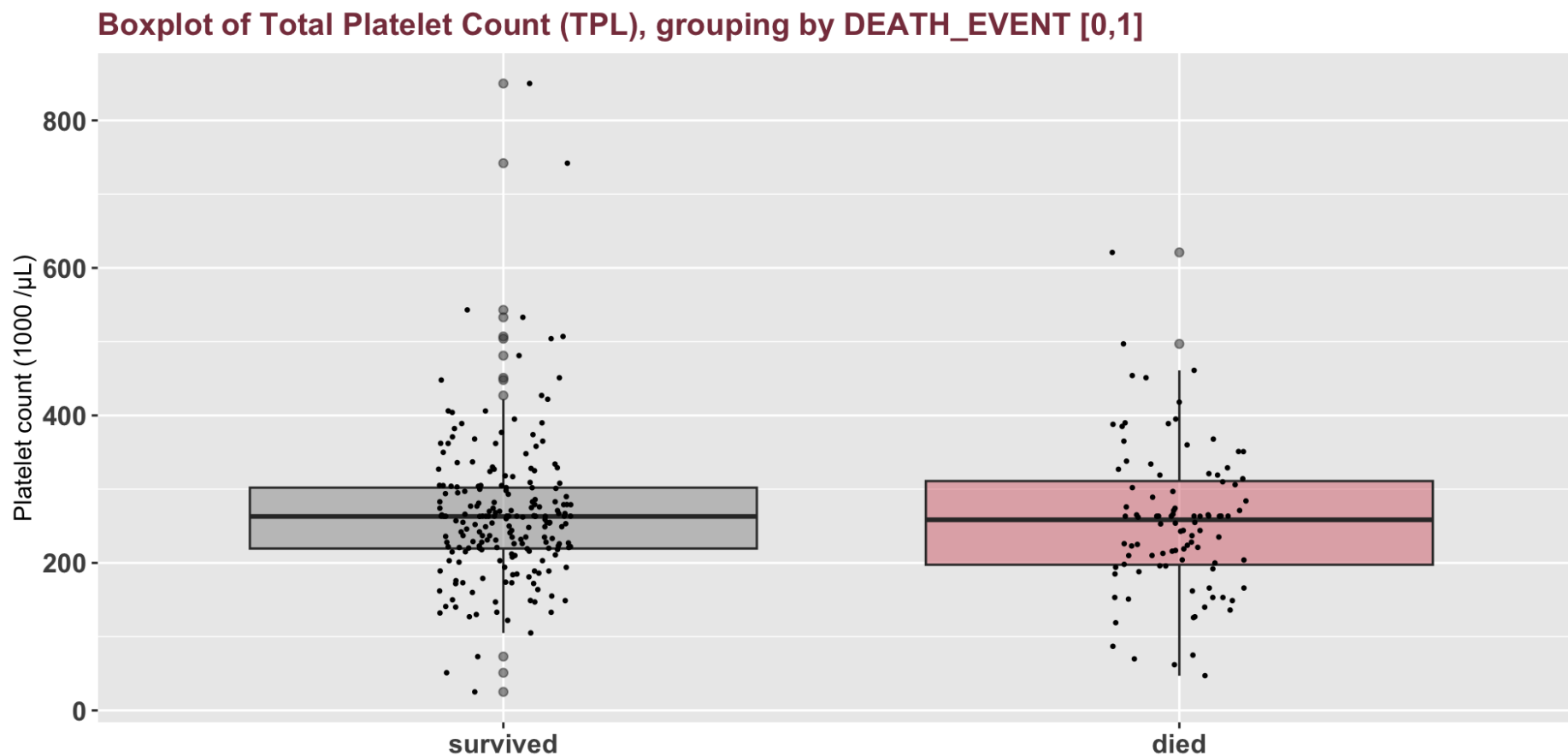
# Comparing two independent sample means (t test) (cont.)

1. Question: Is there a statistically significant difference between the Platelet Counts in the patients

```r
1   # boxplot by group
2   heart_failure %>%
3     ggplot(mapping = aes(y = plat_norm, x = DEATH_EVENT_f, fill = DEATH_EVENT_f)) +
4     geom_boxplot(alpha=0.5)+
5     #geom_violin(alpha=0.5) +
6     geom_point(position = position_jitter(width = 0.1), size = 0.5)+
7     scale_fill_manual(values = c("#999999", "#d8717b"))  +
8     # drop legend and Y-axis title
9     theme(plot.title = element_text(size = 14,face="bold", color = "#873c4a"),
10          legend.position = "none",
11          axis.text.x = element_text(size=12,face="bold"),
12          axis.text.y = element_text(size=12,face="bold")) +
13    labs(title = "Boxplot of Total Platelet Count (TPL), grouping by DEATH_EVENT [0,1]",
14         x = "", y  = "Platelet count (1000 /µL)")
```

# Comparing two independent sample means (t test) (cont.)

There seems to be no major difference in the two groups



Boxplot of Total Platelet Count (TPL), grouping by DEATH_EVENT [0,1]

# 2. Verify the assumptions for independent t-test

1. The 2 samples ("died" and "survived") must be independent ✅

2. The dependent variable is scaled in intervals (Platelets Count in 10^3 "/µL") ✅

3. The dependent variable is normally distributed (Platelets Count in 10^3 "/µL") ✅

- (If not, use *non parametric* test)

4. The variance within the 2 groups should be similar  **?**

- (If not, perform Welch's t-test)

# Preliminary Fisher's F test to check for variance equality

- We can compute the Fisher test "by hand" ✍🏻

```r
1   ## -- data by group
2   n_died <- nrow(heart_failure[heart_failure$DEATH_EVENT == 1 ,])
3   mean_died <- mean(heart_failure [ heart_failure$DEATH_EVENT == 1,  "plat_norm"])
4   sd_died <- sd(heart_failure [heart_failure$DEATH_EVENT == 1 ,  "plat_norm"])
5   var_died <- var(heart_failure [heart_failure$DEATH_EVENT == 1 ,  "plat_norm"])
6
7   n_survived <- nrow(heart_failure[heart_failure$DEATH_EVENT == 0, ])
8   mean_survived <- mean(heart_failure [ heart_failure$DEATH_EVENT == 0,  "plat_norm"])
9   sd_survived <- sd(heart_failure [heart_failure$DEATH_EVENT == 0 ,  "plat_norm"])
10  var_survived <- var(heart_failure [heart_failure$DEATH_EVENT == 0 ,  "plat_norm"])
11
12  ## -- F TEST
13  F_ratio <- var_died / var_survived
14  F_ratio  # 1.020497
```

```
[1] 1.020497
```

# Preliminary Fisher's F test to check for variance equality (.cont)

```r
1  ## -- Define the critical value of F distribution for a risk of alpha = 0.05
2  # qf(p=.05, df1 = n_died-1, df2 = n_survived-1, lower.tail = FALSE) # RIGHT-Tailed
3  # qf(0.95, df1 = n_died-1, df2 = n_survived-1, lower.tail = FALSE) # LEFT- Tailed
4  qf(c(0.025, 0.975), df1 = n_died-1, df2 = n_survived-1) # TWO-Tailed
```

```
[1] 0.6994659 1.3987233
```

```r
1  ## --Compute the exact p-value (two-tailed )
2  p_value_f <- 2 * (1 - pf(F_ratio, df1 = (n_died-1), df2 = (n_survived-1)))
3  p_value_f
```

```
[1] 0.8914982
```

A test statistic (F) of 1.02 is obtained, with degrees of freedom 95 and 202.

The p-value is 0.89, greater than the p-value threshold of 0.05. This suggests **we can not reject the null hypothesis of equal variances**.

The variance within the 2 groups should be similar ✅ –> we can run a t-test.

# 3.a Computation of t test statistic

Since we verified the required assumptions, the test method is the independent (two-sample) t-test. In this case, we have:

- a large sample ($\mathbf{n}_1 + \mathbf{n}_2 > 100$)

- the population variance(s) are unknown, but we can assume = variances in 2 groups

- *standard error* of the means' difference is obtained as **pooled estimate standard deviation of the sampling distribution of the difference**

So we can compute: $t_{calc} = \dfrac{Difference\, Between\, Sample\, means}{Std.\, Err.\, of\, the\, difference} = \dfrac{\bar{x}_1 - \bar{x}_2}{\sqrt{\dfrac{s_1^2}{n_1} + \dfrac{s_2^2}{n_2}}}$

```
 1  # Step 1 - compute difference of sample means
 2  mean_diff <- (mean_died - mean_survived) # -10.27645
 3
 4  # Step 2 - Compute associated t-statistics
 5  # pooled std error
 6  pooled_stderror <- sqrt(sd_died^2/(n_died ) + sd_survived^2/(n_survived ))
 7  # pooled std error corrected
 8  pooled_stderror_corr <- sqrt(sd_died^2/(n_died-1) + sd_survived^2/(n_survived-1))
 9
10  ###  t statistic
11  t_calc <- (mean_died - mean_survived) / pooled_stderror_corr
```

# 3.b Computation of the p-value associated to the t statistic

```r
1  # Step 3 - degrees of freedom
2  # n1 + n2 - number of estimated parameters (2 means)
3  d_f <- n_died + n_survived - 1 - 1 # 297
4
5  # Step 4 - Deduced p-value
6  p_value <- 2 * pt(t_calc, df = d_f) # 0.4009635
7  p_value
```

```
[1] 0.4009635
```

# 4. Results and interpretation

1. Looking at the confidence interval of the difference, the `sample mean_diff` is well inside the 95% CI of = population mean

```
1  mean_diff
```

```
[1] -10.27645
```

```
1  # CI of the means difference
2  CI_lower <- mean_diff + qt(.025, sum(n_died + n_survived) - 2) * pooled_stderror_corr
3  CI_lower
```

```
[1] -34.32074
```

```
1  CI_upper <- mean_diff + qt(.975, sum(n_died + n_survived) - 2) * pooled_stderror_corr
2  CI_upper
```

```
[1] 13.76785
```

2. As for the p-value, `p_value = 0.40` is bigger than threshold probability $\alpha$

```
1  # Check
2  p_value
```

```
[1] 0.4009635
```

```
1  p_value <  0.05  # FALSE
```

```
[1] FALSE
```

**DECISION**: So, we fail to reject the null hypothesis of equal populations means of TPC. So the test indicates that we do not have sufficient evidence to say that the mean counts of platelets in between these two populations is different.

# — EXAMPLE D —

**(3+ samples | ANOVA test)**

# Comparing sample means from 3 or more groups (ANOVA)

In this example, we adopt the ANOVA ("Analysis Of Variance") test, i.e. an extension of the previous test, but examined how means of a variable differ across 3 or more groups. We will use **'one- way' ANOVA**, which serves when there is only one explanatory variable ("treatment") with 3 or more levels, and only one level of treatment is applied for a given subject.

For this particular case, we use another realistic dataset showing the survival times of 33 laboratory mice with thymic leukemia who were randomly divided into 3 groups:

- 1st group received Treatment 1

- 2nd group received Treatment 2

- 3rd group as Control

```r
1  # load new dataset
2  mice <- readxl::read_excel(here::here("practice","data_input",
3                                        "02_datasets","mice_exe_ANOVA.xlsx"))
```

**1. Question: Is there a statistically significant difference between the mean values of the k**
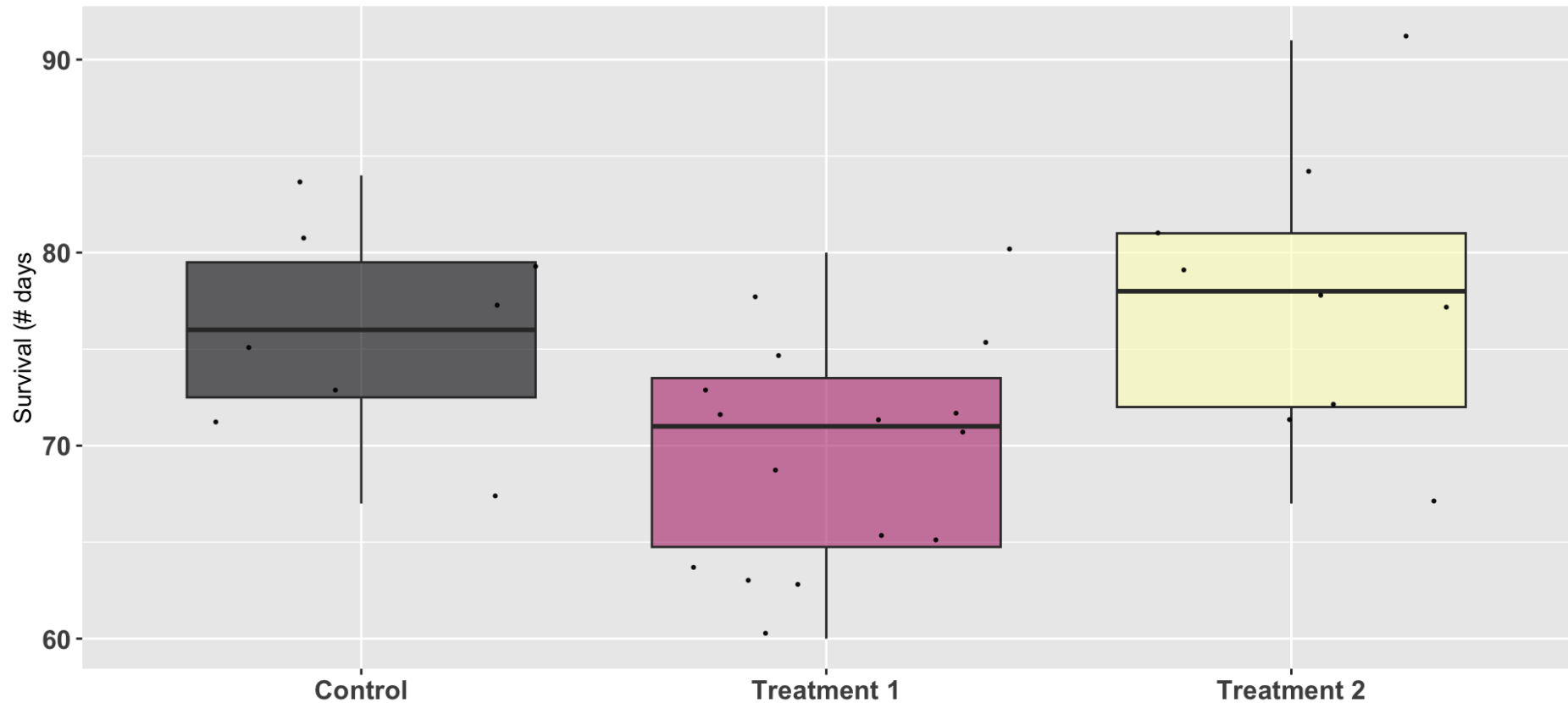
Defining the question formally:

- $H_0 : \mu_1 = \mu_2 = \mu_3$ all 3 population means are equal

- $H_a$ : at least one of $(\mu_1, \mu_2, \mu_3)$ is not equal to the other means

```r
1   # boxplot by group
2   mice %>%
3   ggplot(., aes(x = group, y = surv_days, fill = group)) +
4       geom_boxplot() +
5       scale_fill_viridis(discrete = TRUE, alpha=0.6, option="A") +
6       geom_jitter(color="black", size=0.4, alpha=0.9) +
7       # theme_minimal() +
8       # drop legend and Y-axis title
9       theme(plot.title = element_text(size = 14,face="bold", color = "#873c4a"),
10          axis.text.x = element_text(size=12,face="bold"),
11          axis.text.y = element_text(size=12,face="bold"),
12          legend.position = "none",
13          ) +
14      labs(title = "Visually check mean and variance in populations' samples" ) +
15      ylab(label = "Survival (# days") + xlab(label = "")
```

# 1. Question: Is there a statistically significant difference between the mean values of the k p

The boxplot suggests that the 3 groups might have some fairly different distributions



**Visually check mean and variance in populations' samples**
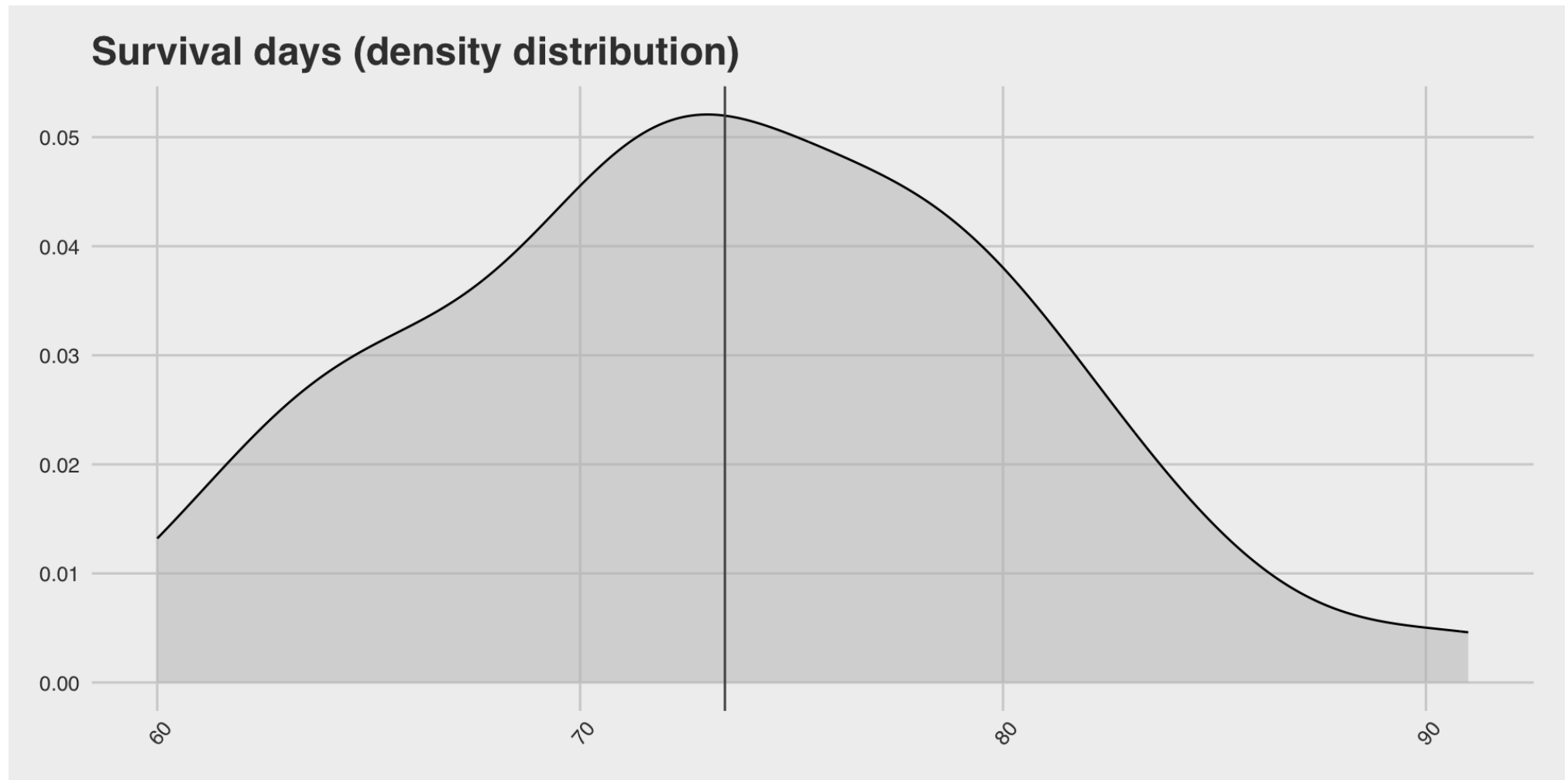
# 2. Verify the assumptions for one-way ANOVA

The dependent variable is on a metric scale. In the case of the analysis of variance, the independent variable (factor) has at least three levels.

Assumptions for the results of a one-way ANOVA to be valid:

1. **Independence of observations** – The observations in each group are independent of each other and the observations within groups were obtained by a random sample. ✅

2. **Normally-distributed response variable** – The values of the dependent variable follow a normal distribution. **?**

3. **Homogeneity of variance** – The variances of the populations that the samples come from are equal. **?**

# Preliminary check for normality (visual)

2. **Normally-distributed response variable** ✅ • (confirmed by visual inspection )



**Survival days (density distribution)**

# Preliminary check for normality (test) with `stats::shapiro.test`

```
1  # Shapiro-Wilk Normality Test to verify normality
2  # option 1
3  stats::shapiro.test(mice[mice$group == "Control", "surv_days", drop=TRUE])
```

```
    Shapiro-Wilk normality test

data:  mice[mice$group == "Control", "surv_days", drop = TRUE]
W = 0.99374, p-value = 0.9989
```

```
1  stats::shapiro.test(mice[mice$group == "Treatment 1", "surv_days", drop=TRUE])
```

```
    Shapiro-Wilk normality test

data:  mice[mice$group == "Treatment 1", "surv_days", drop = TRUE]
W = 0.95716, p-value = 0.6106
```

```
1  stats::shapiro.test(mice[mice$group == "Treatment 2", "surv_days", drop=TRUE])
```

```
    Shapiro-Wilk normality test

data:  mice[mice$group == "Treatment 2", "surv_days", drop = TRUE]
W = 0.97921, p-value = 0.9601
```

# Preliminary check for normality (test) with rstatix::shapiro_test

(same thing, but using a different R function)

2. **Normally-distributed response variable** – ✅

- (confirmed by Shapiro-Wilk normality test)

[The null hypothesis of this test is $H_0$ = **"sample distribution is normal"** ]

```
1  # Shapiro-Wilk Normality Test to verify normality
2  # option 2 (all 3 groups at once)
3  mice %>%
4    dplyr::group_by(group) %>%
5    rstatix::shapiro_test(surv_days)
```

```
# A tibble: 3 × 4
  group       variable   statistic     p
  <chr>       <chr>          <dbl> <dbl>
1 Control     surv_days      0.994 0.999
2 Treatment 1 surv_days      0.957 0.611
3 Treatment 2 surv_days      0.979 0.960
```

# Preliminary check variance equality

3. **Homogeneity of variance** – ✅

- (Besides visual inspection, confirmed by Levene test for variance equality)

[The null hypothesis $H_0$ = **several groups have the same variance** (possible variance differences occur only by chance, since there are small differences in each sampling)]

```
1   # Levene test for variance equality
2   levene <- mice %>%                              # name of the data
3     car::leveneTest(surv_days ~ as.factor(group),  # continuous DV ~  group IV
4                  data = .,                 # pipe the data from above
5                  center = mean)        # default is median
6   levene
```

```
Levene's Test for Homogeneity of Variance (center = mean)
      Df F value Pr(>F)
group  2  0.1721 0.8427
      30
```

No evidence of violations of HOV were found, since the p-value for the Levene test (= 0.8427157) is greater than .05, then the variances are not significantly different from each other (i.e., the homogeneity assumption of the variance is met).

# 3 Computation of ANOVA F-ratio

ANOVA in R can be done in several ways.

Since it's quite straightforward, let's do all the steps by hand first. We need to obtain the needed "ingredients" to calculate the F-ratio:

$$F_{calc} = \frac{Mean\ Square\ Between}{Mean,\ Square\ Within} = \frac{MSB}{MSW} = \frac{\frac{SSB}{df1}}{\frac{SSW}{df2}}$$

# 3.a Computation of ANOVA F-ratio ("by hand")

- Option 1: Let's compute the ANOVA test "by hand" ✍🏻

```r
1   # Summary statistics
2   mice_calc <- mice %>%
3     dplyr::mutate(mean_all = mean(surv_days),
4             sd_all = sd (surv_days),
5             dfw = 33-3, # df1 = n-k
6             dfb = 3-1, # df2 = K-1
7             group_f = as.factor(group)
8             ) %>%
9     dplyr::group_by(group) %>%
10    dplyr::mutate(n_group = n(),
11            mean_group = mean(surv_days),
12            sd_group = sd (surv_days)) %>%
13    ungroup() %>%
14    mutate (ST = (surv_days - mean_all)^2,
15            SW = (surv_days - mean_group)^2,
16            SB = (mean_group - mean_all)^2)
17
18  # Sum of Squares
19  SST <- sum(mice_calc$ST)
20  SSB <- sum(mice_calc$SB)
21  SSW <- sum(mice_calc$SW)
22  dfw <- 33-3   # df2
23  dfb <- 3-1 # df1
24
25  # calculated F statistic
26  F_calc <- (SSB/dfb)/(SSW/dfw) # 5.65
27  # F critical value
28  F_crit <- qf(p = 0.01, df1 = 2, df2 = 30, lower.tail = FALSE) # 5.390346
```

# 3.b Computation of ANOVA F-ratio (with R functions)

That was just to show how to build it step-by-step (🤓), but we don't have to! We have alternative R functions that can do ANOVA for us:

- Option 2: With the `stats::aov` followed by the command `summary` 💇🏻‍♀️

```
1  aov_1 <- stats::aov(surv_days ~ group_f,
2                        data = mice_calc)
3  summary(aov_1)
```

```
            Df Sum Sq Mean Sq F value  Pr(>F)
group_f      2  434.6  217.32   5.652 0.00826 **
Residuals   30 1153.4   38.45
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

- Option 3: With the `stats::oneway.test()` function 💇🏻‍♀️

```
1  aov_2 <- stats::oneway.test(surv_days ~ group_f,
2                  data = mice_calc,
3                  # assuming equal variances
4                  var.equal = TRUE)
5  aov_2
```

```
    One-way analysis of means

data:  surv_days and group_f
F = 5.6522, num df = 2, denom df = 30, p-value = 0.008258
```

# 4. Results and interpretation

All 3 options have given the same results, i.e., `F-ratio = 5.652` and a `p-value = 0.00826`

**DECISION**: Given that the p-value is smaller than 0.05, we reject the null hypothesis, so we reject the hypothesis that all means are equal. Therefore, we can conclude that *at least one* group is different than the others in mean number of survival days.

> **(i) Note**
>
> Have you seen the kind of notation `Pr(>F) 0.00826 **` before (as in the output of the `stats::aov` function)?

# A CLOSER LOOK AT TESTING ASSUMPTIONS

# — EXAMPLE E —
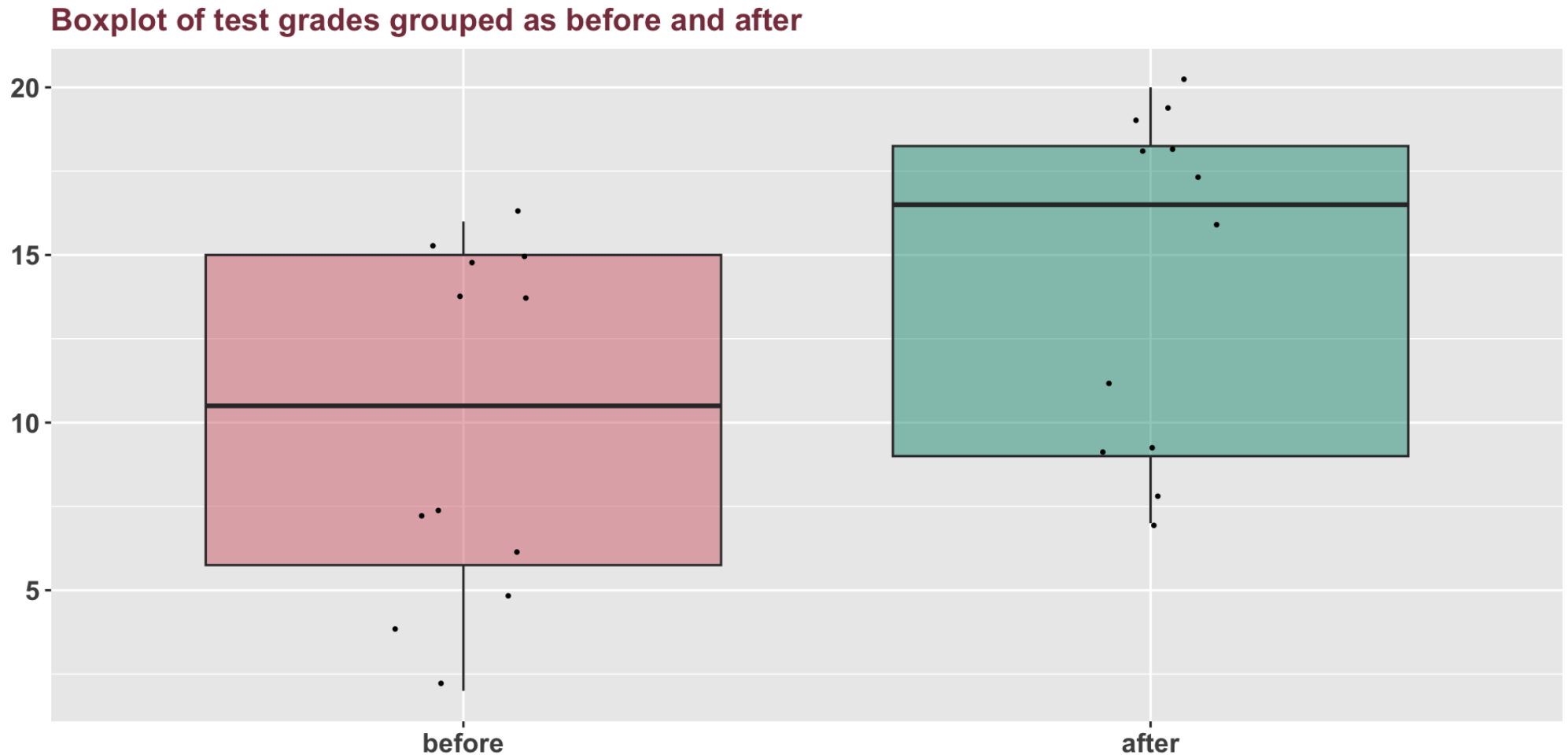
# Testing two groups that are *not* independent

Let's introduce another toy dataset just for demonstration purposes: imagine a statistics test is administered to the *same* group of 12 students **before and after** attending a workshop 😉.

```r
1  # toy dataset for paired groups
2  grades <- data.frame(
3    before = c(16, 5, 15, 2, 14, 15, 4, 7, 15, 6, 7, 14),
4    after = c(19, 18, 9, 17, 8, 7, 16, 19, 20, 9, 11, 18)
5  )
```

We may reshape the dataframe into the long form using `tidyr::pivot_longer` (for plotting)

```r
1   # reshape into long form
2   grades_long <- grades %>%
3     dplyr::mutate(id = row_number()) %>%
4     tidyr::pivot_longer(cols = before:after,
5                         names_to = "time",
6                         values_to = "grade") %>%
7     dplyr::group_by(id) %>%
8     # recode time as factor
9     dplyr::mutate(time_f = as_factor(time )) %>%
10    # reorder time_ levels
11    dplyr::mutate(time_f = fct_relevel(time_f, "after", after = 1))
```

# 1. Question: Is the difference between two PAIRED samples statistically significant?



**Boxplot of test grades grouped as before and after**

What a successful workshop! 😁

# 2 Hypotehsis for the PAIRED t-test for dependent samples

In this example, it is clear that the two samples are not independent since the same 12 students took the test before and after the workshop.

Given that the normality assumption is NOT violated (and given the small sample size), we use the **paired t-test**, with the following hypotheses:

- $H_0$ : mean grades before and after the workshop are equal
- $H_a$ : mean grades before and after the workshop are different

# 2 Computation of the PAIRED t-test for dependent samples

```r
1  t_stat_paired <- stats::t.test(x = grades$before,
2                                 y = grades$after,
3                                 mu = 0,
4                                 alternative = "two.sided",
5                                 paired = TRUE
6  )
7  t_stat_paired
```

```
    Paired t-test

data:  grades$before and grades$after
t = -1.8777, df = 11, p-value = 0.08718
alternative hypothesis: true mean difference is not equal to 0
95 percent confidence interval:
 -9.2317713  0.7317713
sample estimates:
mean difference
          -4.25
```

```r
1  # extract t_calc from results df
2  t_calc_pair   <- t_stat_paired[["statistic"]][["t"]] # -1.877683
3  p_value_pair  <- t_stat_paired[["p.value"]] # -1.877683
```

# 3. Results and interpretation

We obtain the test statistic, the p-value and a reminder of the hypothesis tested.

The calculated **t value** is -1.8776829 The **p-value** is 0.087177. Therefore, at the 5% significance level, **we do not reject the null hypothesis** that the statistics' grades are similar before and after the workshop (😭).
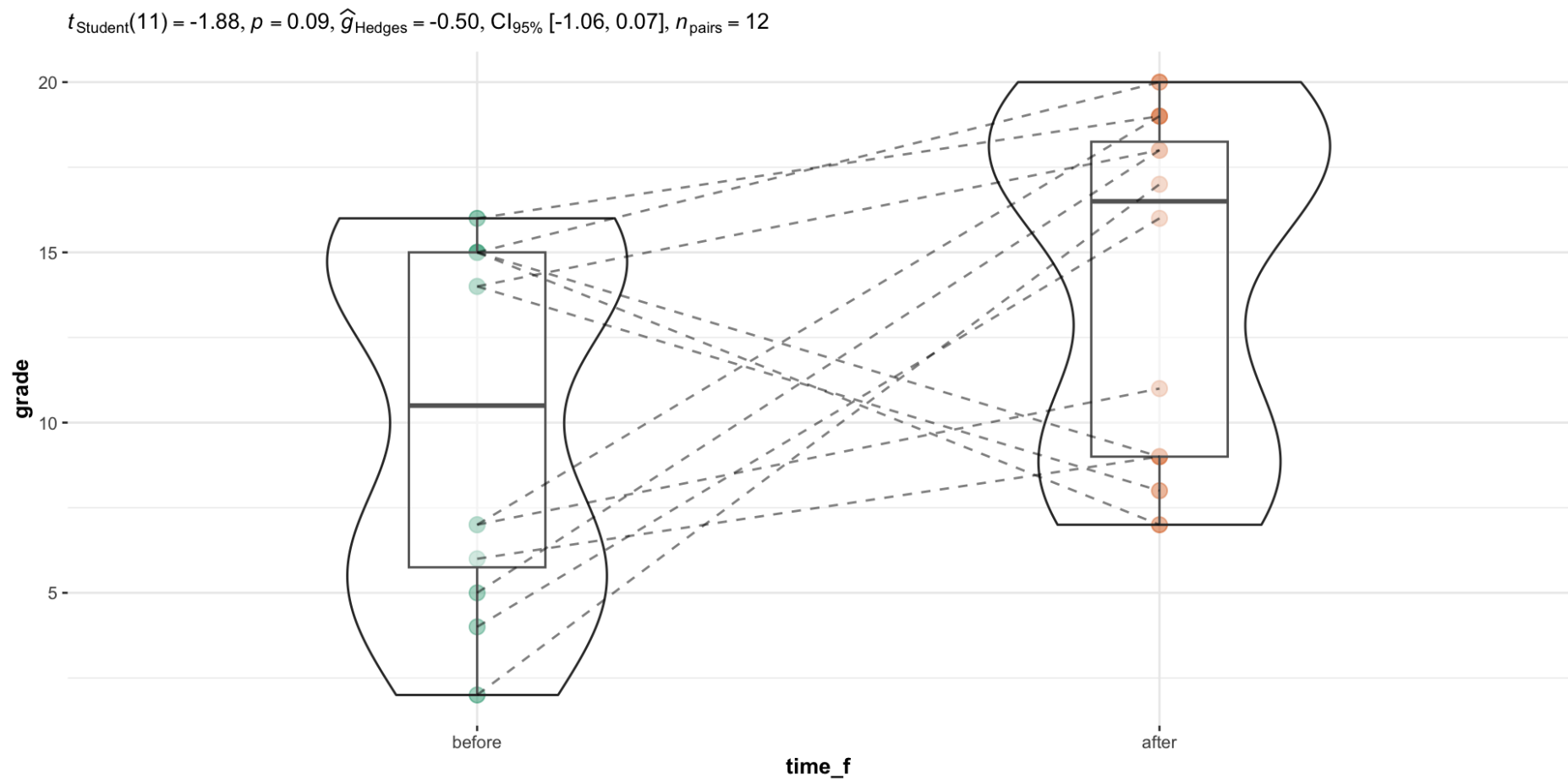
# Bonus function!

It is worth mentioning the `ggstatsplot` package, which combines plots representing the distribution for each group—and the results of the statistical test displayed in the subtitle of the plot.

Below we check out the `ggwithinstats()` function for *paired samples*.

```r
1   # load package
2   library(ggstatsplot)
3
4   # plot with statistical results
5   grades_long %>%
6     # must ungroup the dataframe or it will give an error
7     ungroup () %>%
8     ggstatsplot::ggwithinstats(.,
9                                 x = time_f ,
10                                y = grade ,
11                                type = "parametric", # for wilcoxon
12                                centrality.plotting = FALSE # remove median
13   )
```

# Bonus function!

The test results are rendered with the plot!



$t_{Student}(11) = -1.88, p = 0.09, \widehat{g}_{Hedges} = -0.50, CI_{95\%}$ [-1.06, 0.07], $n_{pairs} = 12$

# — EXAMPLE F —

**(2 samples no normal | Wilcoxon Rank Sum Test)**

# Testing samples *without* normality assumption

Let's go back to the HEART FAILURE dataset but looking at the levels of **Creatinine Phosphokinase (CPK)** in the blood, an enzyme that might indicate a heart failure or injury

**1. Question: Is there a statistically significant difference between CPK levels in the blood of**
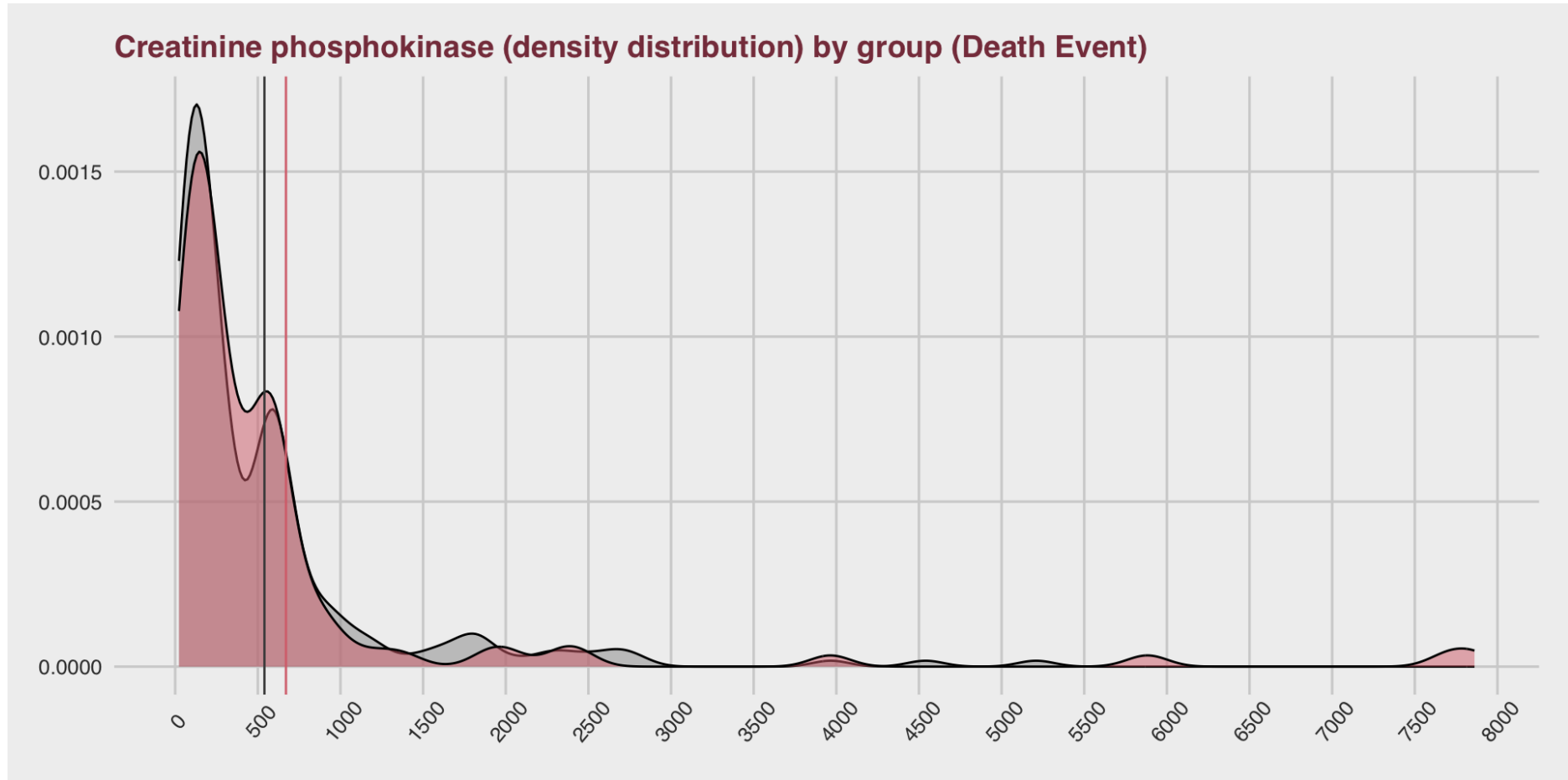
Defining the question formally:

- $H_0 : \mu_{CPK-died} = \mu_{CPK-surv}$ there is no difference in mean CPK between patients who suffered heart failure and died versus patients who survived after heart failure

- $H_a : \mu_{CPK-died} \neq \mu_{CPK-surv}$ there is a difference in mean CPK between patients who suffered heart failure and died versus patients who survived after heart failure (two-sided test)

```
1  ggplot(heart_failure,aes(x = creatinine_phosphokinase,fill = DEATH_EVENT_f))+
2    geom_density(alpha = 0.5)+theme_fivethirtyeight()+
3    scale_fill_manual(values = c("#999999", "#d8717b"))+
4    guides(fill = "none") +
5    scale_x_continuous(breaks = seq(0,8000, 500))+
6    geom_vline(aes(xintercept = mean(creatinine_phosphokinase[DEATH_EVENT == 0])),
7              color = "#4c4c4c")+
8    geom_vline(aes(xintercept = mean(creatinine_phosphokinase[DEATH_EVENT==1])),
9              color = "#d8717b")+
10   theme_fivethirtyeight()+
11   theme(axis.text.x = element_text(angle=50, vjust=0.75))+
12   labs(title =  "Creatinine phosphokinase (density distribution) by group (Death Event)") +
13   theme(plot.title = element_text(size = 14,face="bold", color = "#873c4a"))
```

# 1. Question: Is there a statistically significant difference between CPK levels in the blood of

The density plot suggests non normality of the variable distribution



Creatinine phosphokinase (density distribution) by group (Death Event)

# Preliminary check for normality (visual)

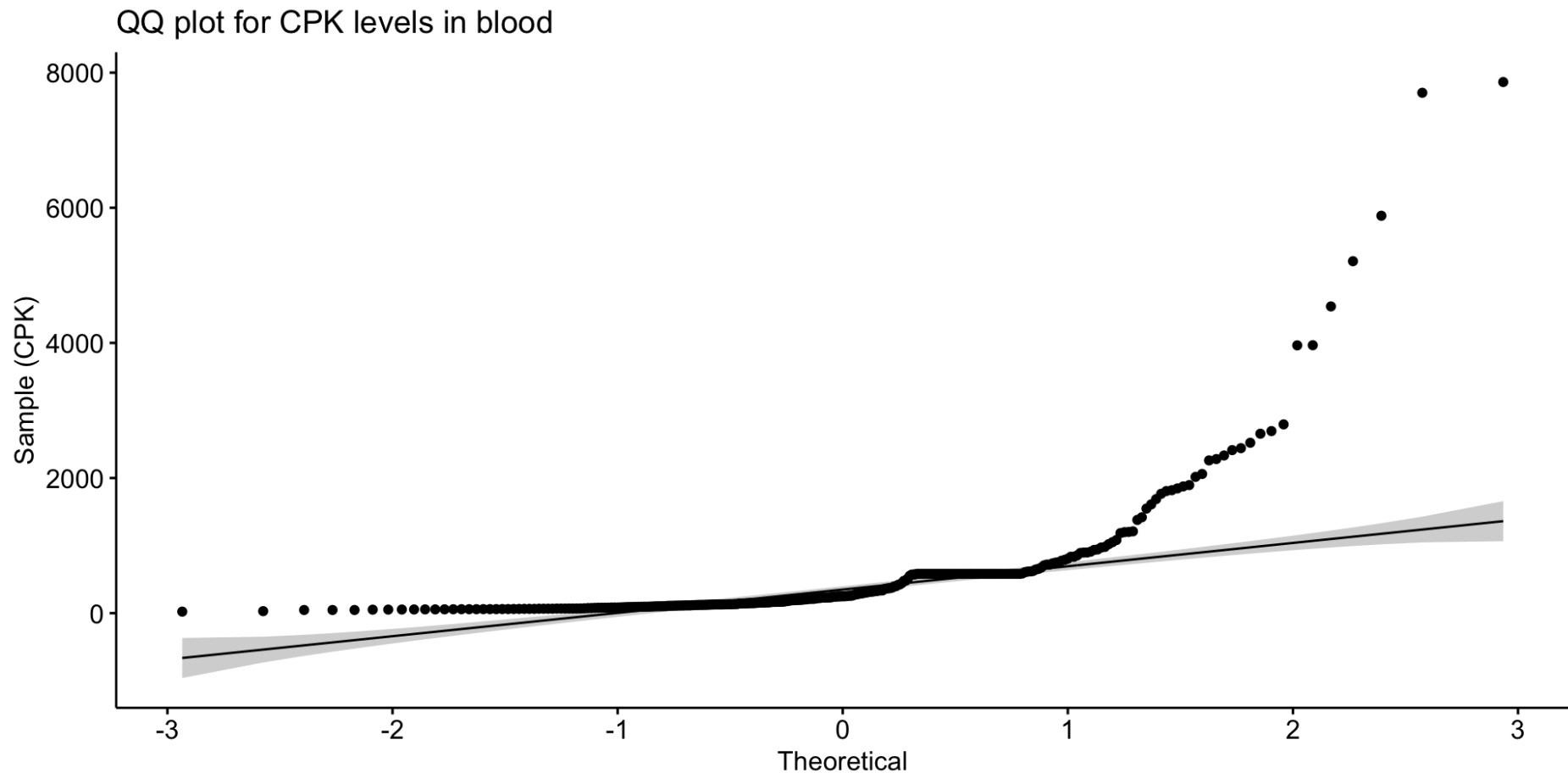- **Normally-distributed response variable** - ❌

**QQ plot** (or quantile-quantile plot) draws the correlation between a given sample and the normal distribution. A 45-degree reference line is also plotted. In a QQ plot, each observation is plotted as a single dot.

- If the data are normal, the dots should form a straight line.

```r
1  # visual verification with QQ plot
2  ggpubr::ggqqplot(
3    heart_failure$creatinine_phosphokinase,
4    title = "QQ plot for CPK levels in blood",
5    xlab ="Theoretical", ylab = "Sample (CPK)")
```

# Preliminary check for normality (visual)

In a QQ plot, if the data are normal, the dots should follow a straight line.



QQ plot for CPK levels in blood

# Preliminary check for normality (test) with **rstatix::shapiro_test**

(same thing, but using a different R function)

- **Normally-distributed response variable** - ❌

    - (NOT normality confirmed by Shapiro-Wilk normality test)

[The null hypothesis of this test is $H_0$ = "sample distribution(s) is/are normal" ]

Given the p-value we reject the null hypothesis

```
1  # Shapiro-Wilk Normality Test to verify normality
2  heart_failure %>%
3    dplyr::group_by(DEATH_EVENT_f) %>%
4    rstatix::shapiro_test(creatinine_phosphokinase)
```

```
# A tibble: 2 × 4
  DEATH_EVENT_f variable                statistic       p
  <fct>         <chr>                       <dbl>   <dbl>
1 survived      creatinine_phosphokinase    0.628 8.51e-21
2 died          creatinine_phosphokinase    0.439 1.99e-17
```

# 3. Computation of the Wilcoxon Rank Sum test statistic

The **Wilcoxon Rank Sum test** is considered to be the nonparametric equivalent to the **two-sample independent t-test**

Its ASSUMPTIONS are:

- Ordinal or Continuous dependent variable: e.g. CPK levels ✅

- Independence: All of the observations from both groups are independent of each other ✅

- Shape: The shapes of the distributions for the two groups are roughly the same ✅

```
1  wrs_res <- wilcox.test(creatinine_phosphokinase ~ DEATH_EVENT, # immagino 0, 1
2                          data = heart_failure ,
3                          exact = FALSE,
4                          alternative = "two.sided" )
5  wrs_res
```

```
        Wilcoxon rank sum test with continuity correction

data:  creatinine_phosphokinase by DEATH_EVENT
W = 9460, p-value = 0.684
alternative hypothesis: true location shift is not equal to 0
```

The **Wilcoxon Rank Sum test** is equivalent to the **Mann-Whitney U test** to compare two independent samples. Different

# 4. Results and interpretation

RESULTS: since the test statistic is `W = 9460` and the corresponding `p-value is 0.684 > 0.05`, we fail to reject the null hypothesis.

INTERPRETATION: We do not have sufficient evidence to say that CPK levels for dead patients is different than that of survived patients $\mu_{CPK-died} \neq \mu_{CPK-surv}$ at some statistically significant level)

# — EXAMPLE G —

**(2 samples no HOV | t test with the Welch correction )**

# Testing samples *without* homogeneous variance of observations assumption

**1. Question: Is there a statistically significant difference between serum sodium levels in th**
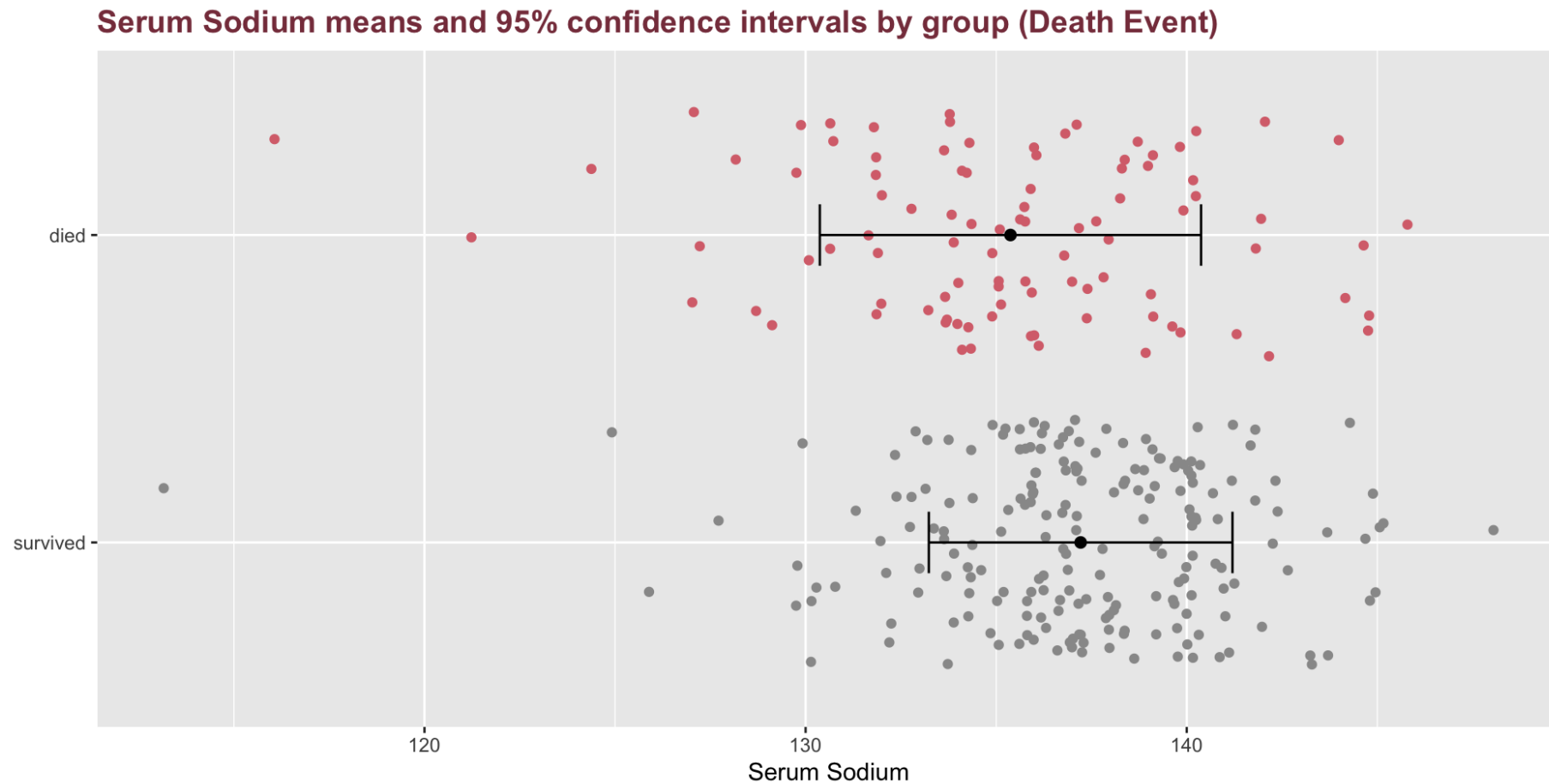
Defining the question formally:

- $H_0 : \mu_{sersod-died} = \mu_{sersod-surv}$ there is no difference in mean serum sodium between patients who suffered heart failure and died versus patients who survived after heart failure

- $H_a : \mu_{sersod-died} \neq \mu_{sersod-surv}$ there is a difference in mean serum sodium between patients who suffered heart failure and died versus patients who survived after heart failure (two-sided test)

# Preliminary check "HOV" assumption (visual)

- **Homogeneity of Variance assumption** - ❌ Plotting the data offers some graphical intuition that the variance of observations in the two groups seem not homogenous

```r
1  #Compute means and 95% confidence intervals
2  swstats <- heart_failure %>%
3    group_by(DEATH_EVENT_f) %>%
4    summarise(count = n(),
5      mean = mean(serum_sodium,na.rm=TRUE),
6      stddev = sd(serum_sodium, na.rm=TRUE),
7      meansd_l = mean - stddev,
8      meansd_u = mean + stddev)
9
10 #The complete script with some styling added
11 ggplot(swstats, aes(x=DEATH_EVENT_f, y=mean)) +
12   geom_point(colour = "black" , size = 2) +
13   #Now plotting the individual data points before the mean values
14   geom_point(data=heart_failure, aes(x=DEATH_EVENT_f, y=serum_sodium, colour = DEATH_EVENT_f),
15            position = position_jitter() ) +
16   scale_colour_manual(values = c("#999999","#d8717b") ) +
17   #Add the error bars
18   geom_errorbar(aes(ymin = meansd_l, ymax = meansd_u), width=0.2, color = "black") +
19   labs(title = "Mean (-/+SD) serum sodium (mEq/L) by group", x = "", y = "Serum Sodium") +
20   guides(fill = "none")  +
21   coord_flip() +
22   labs(title =  "Serum Sodium means and 95% confidence intervals by group (Death Event)") +
23   theme(legend.position="none",plot.title = element_text(size = 14,face="bold", color = "#873c4a"))
```

# Preliminary check "HOV" assumption (visual)



Serum Sodium means and 95% confidence intervals by group (Death Event)

# Preliminary check "HOV" assumption (test)

It is always best to use an actual test, so we use also the **Fisher's F test** to verify equal variances of Serum Sodium concentration in the two groups. [In this test $H_0 =$ **"the ratio of variances is equal to 1"**]

```
1  f_test_res <- stats::var.test(heart_failure$serum_sodium[heart_failure$DEATH_EVENT == 1] ,
2                                heart_failure$serum_sodium[heart_failure$DEATH_EVENT == 0])
3  f_test_res
```

```
    F test to compare two variances

data:  heart_failure$serum_sodium[heart_failure$DEATH_EVENT == 1] and
heart_failure$serum_sodium[heart_failure$DEATH_EVENT == 0]
F = 1.5769, num df = 95, denom df = 202, p-value = 0.007646
alternative hypothesis: true ratio of variances is not equal to 1
95 percent confidence interval:
 1.127401 2.254466
sample estimates:
ratio of variances
        1.576922
```

Given the `p-value = 0.007646` (smaller than $\alpha$) we reject the null hypothesis, hence the HOV assumption for the t test does not hold.

# 2 Computation of the t test with the Welch correction

We can still run the **t test but with Welch correction**, i.e. the unequal variance condition is compensated by lowering the df. In fact the documentation (`?t.test`), reads:

- If `var.equal = TRUE`, then the pooled variance is used to estimate the variance

- Otherwise (`var.equal = FALSE`), the Welch approximation to the degrees of freedom is used.

```
1  # With Welch correction (on by default) Unequal variance is compensated by lowering df
2  t_test_w <- t.test(heart_failure$serum_sodium[heart_failure$DEATH_EVENT == 1],
3                     heart_failure$serum_sodium[heart_failure$DEATH_EVENT == 0],
4                     # here we specify the situation
5                     var.equal = FALSE,
6                     paired = FALSE, alternative = "two.sided")
7
8  t_test_w
```

```
	Welch Two Sample t-test

data:  heart_failure$serum_sodium[heart_failure$DEATH_EVENT == 1] and
heart_failure$serum_sodium[heart_failure$DEATH_EVENT == 0]
t = -3.1645, df = 154.01, p-value = 0.001872
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
 -2.9914879 -0.6920096
sample estimates:
mean of x mean of y
 135.3750  137.2167
```

# 3. Results and interpretation

RESULTS: since the test statistic is `t = -3.1645 (with df = 154.01)` and the corresponding `p-value is 0.001872 < 0.05`, we reject the null hypothesis.

INTERPRETATION: We therefore have sufficient evidence to say that the level of serum sodium levels for dead patients is significantly different than that of survived patients $\mu_{sersod-died} \neq \mu_{sersod-surv}$

# Final thoughts/recommendations

- There are often **many ways to do the same thing in R** (which is both a blessing and a curse in open source software). *Which should you choose?* It depends on the situation, but you may want to consider:

  - how recent/popular/well maintained is a `{package}` (this affects its stability)

  - the more a function abstracts away complexity, the easier it is to use interactively, but the harder it gets to handle inside your own custom functions

  - different function outputs may be more/less suitable for your analysis/publication requirements (check out your peers' choices!)

  - *(Always **read the documentation** to assess all of the above)*

- With easy equations, breaking them down "by hand" (at least once!) can really help you understand them

- It may seem a lot of work to write R code the first time 🥵 (e.g. for a publication-ready plot), but the good news is **once you wrote a script, you will be able to easily re-use it in many more instances** 🙌🏻 😃

- **Sample size $n$ has a very powerful impact** on classical hypothesis testing results! More on this later...