

# Datenvisualisierung 3

## Darstellung der zusammenfassenden Statistik

Daniela Palleschi

2023-12-18

### Inhaltsverzeichnis

|  |           |
|--|-----------|
| <b>Learning objectives</b>                                 | <b>2</b>  |
| <b>Ressourcen</b>  | <b>2</b>  |
| <b>Einrichten</b>  | <b>2</b>  |
| Pakete . . . . .   | 2         |
| Daten . . . . .  | 2         |
| <b>1 Rückblick: Visualisierung von Verteilungen</b>        | <b>3</b>  |
| <b>2 Darstellung von zusammenfassenden Statistiken</b>     | <b>3</b>  |
| 2.1 Boxplot . . . . .                                      | 4         |
| 2.1.1 geom_boxplot() . . . . .                             | 5         |
| 2.1.2 Gruppiertes Boxplot . . . . .                        | 8         |
| <b>3 Visualisierung des Mittelwerts</b>                    | <b>9</b>  |
| 3.1 Fehlerbalkenplots . . . . .                            | 9         |
| 3.1.1 Berechnung der zusammenfassenden Statistik . . . . . | 9         |
| 3.1.2 Plotting mean . . . . .                              | 10        |
| 3.1.3 Hinzufügen von Fehlerbalken . . . . .                | 11        |
| <b>4 Barplot von Mittelwerten: Bleiben Sie weg!</b>        | <b>13</b> |
| Hausaufgabe . . . . .                                      | 14        |
| Boxplot mit Facette . . . . .                              | 14        |
| Errorbar plot . . . . .                                    | 14        |
| Patchwork . . . . .  | 15        |
| <b>Session Info</b>  | <b>15</b> |

## Learning objectives

Today we will learn to...

- create and interpret boxplots
- visualize mean values and standard deviations

## Ressourcen

- [Abschnitt 2.5 \(Visualisierung von Beziehungen\)](#) in Wickham et al. (2023)
- [Kapitel 4 \(Darstellung von zusammenfassenden Statistiken\)](#) in Nordmann et al. (2022)
- Abschnitte 3.5-3.9 in Winter (2019)

## Einrichten

### Pakete

```
pacman::p_load(tidyverse,  
               here,  
               janitor,  
               ggthemes,  
               patchwork)
```

### Daten

```
df_eng <- read_csv(  
  here(  
    "daten",  
    "languageR_english.csv"  
  )  
) |>  
  clean_names() |>  
  rename(  
    rt_lexdec = r_tlexdec,  
    rt_naming = r_tnaming  
  )
```

# 1 Rückblick: Visualisierung von Verteilungen

- Betrachten Sie jede Abbildung in Abbildung 1
  - Wie viele Variablen werden in jeder Abbildung dargestellt?
  - welche *Typen* von Variablen sind es?
  - Welche zusammenfassende(n) Statistik(en) wird/werden in jedem Diagramm dargestellt?

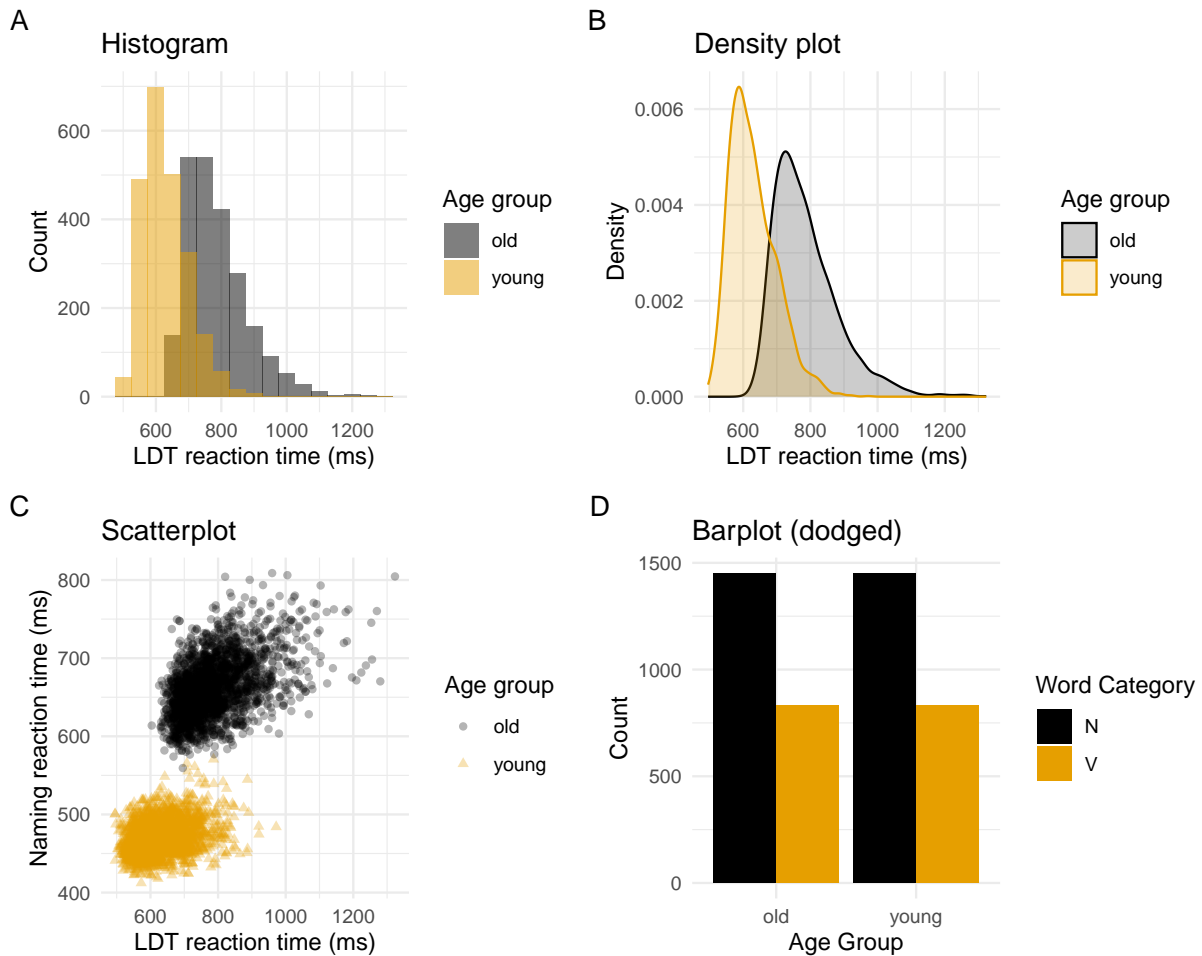


Abbildung 1: Different plots types

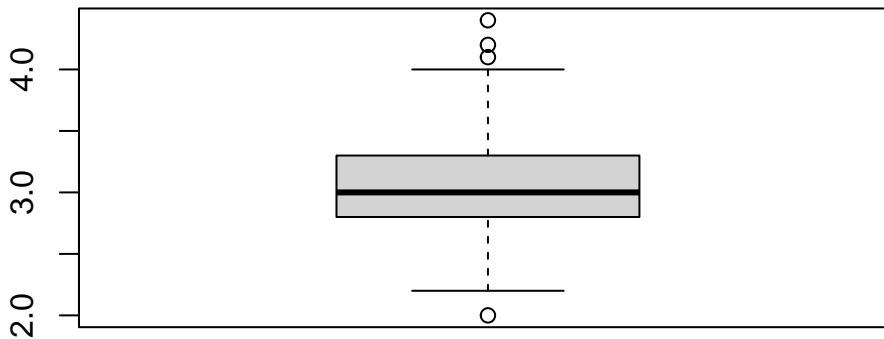
# 2 Darstellung von zusammenfassenden Statistiken

- Modus und Bereich werden in Histogrammen und Dichteplots visualisiert

- die Anzahl der Beobachtungen wird in Balkendiagrammen visualisiert

## 2.1 Boxplot

- auch bekannt als Box-and-Whisker-Plots, enthalten
  - eine Box
  - eine Linie in der Mitte der Box
  - Linien, die an beiden Enden der Box herausragen (die ‘Whisker’)
  - manchmal Punkte



- Betrachten Sie Abbildung 2
  - identifiziere jeden dieser 4 Aspekte des Plots
  - können Sie erraten, was jeder dieser Aspekte darstellen könnte und wie Sie die Darstellung interpretieren sollten?

- Boxplots vermitteln eine Menge Informationen in einer einzigen Visualisierung
  - Die Box selbst stellt den *Interquartilsbereich* (IQR; der Bereich der Werte, der zwischen den mittleren 50% der Daten liegt) dar.
    - \* Die Grenzen der Box repräsentieren Q1 (1. Quartil, unter dem 25% der Daten liegen) und Q3 (3. Quartil, über dem 25% der Daten liegen)
  - die Linie in der Mitte des Boxplots stellt den *Median* dar
    - \* auch Q2 genannt (2. Quartil; der mittlere Wert, über/unter dem 50% der Daten liegen)
  - Die Whisker repräsentieren  $1,5 \cdot \text{IQR}$  von Q1 (unterer Whisker) oder Q3 (oberer Whisker)

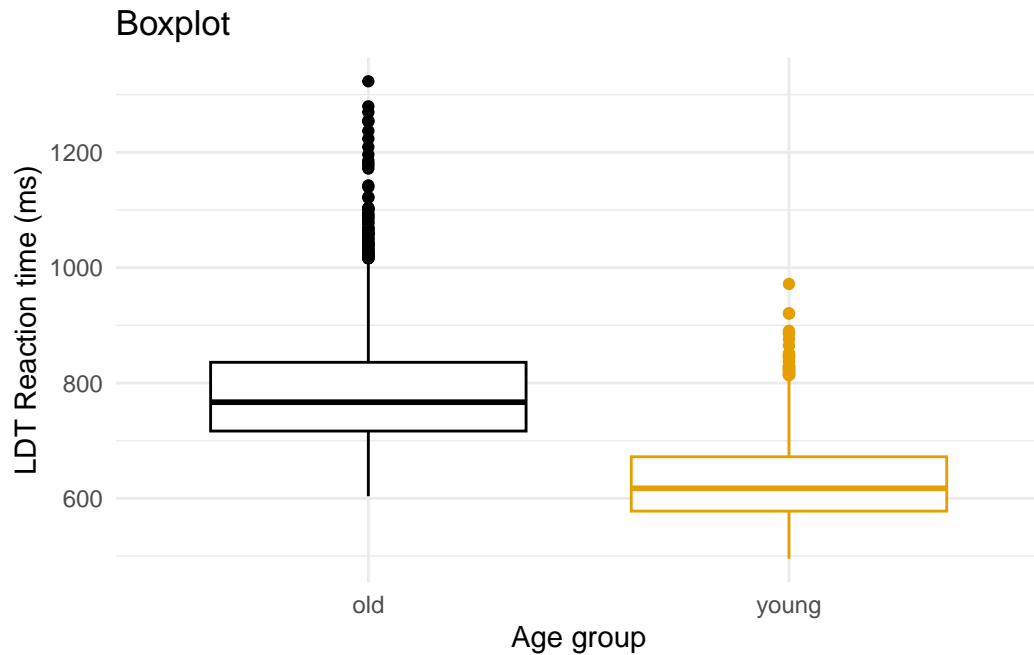


Abbildung 2: Boxplot of `df_eng` (body mass by `age_subject`)

- Punkte, die außerhalb der Whisker liegen, stellen Ausreißer dar (d. h. Extremwerte, die außerhalb des IQR liegen).

- 
- `?@fig-winter-boxplot-hist` zeigt die Beziehung zwischen einem Histogramm und einem Boxplot

- 
- `?@fig-wickham-boxplot-hist` hat einen ähnlichen Vergleich, einschließlich eines Streudiagramms

### 2.1.1 `geom_boxplot()`

- Die Funktion `geom_boxplot()` von `ggplot2` erzeugt Boxplots
  - sie benötigt eine numerische Variable als `x` oder `y` Achse (Abbildung 5)

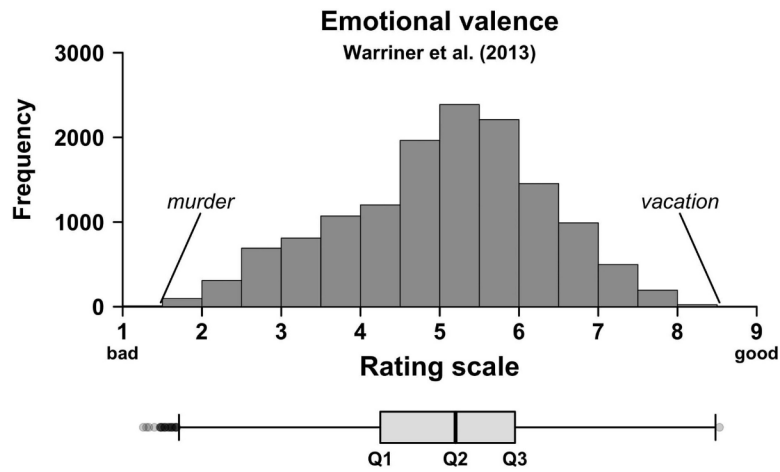


Figure 3.4. A histogram of the emotional valence rating data

Abbildung 3: Image source: Winter (2019) (all rights reserved)

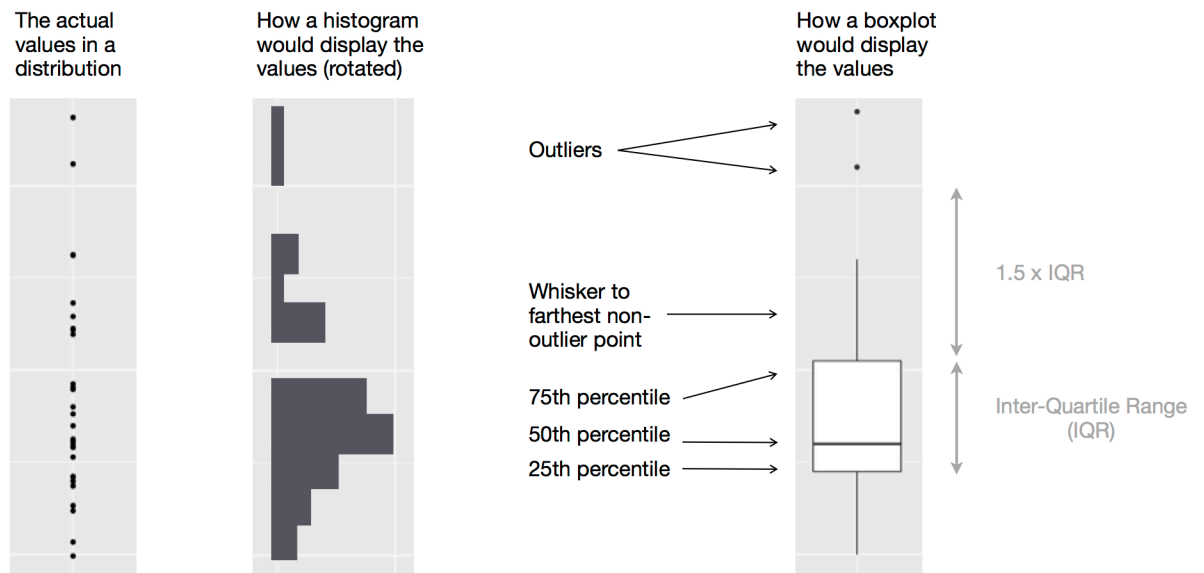


Abbildung 4: Image source: Wickham et al. (2023) (all rights reserved)

```
df_eng |>
  ggplot(aes(y = rt_lexdec)) +
  geom_boxplot()
```

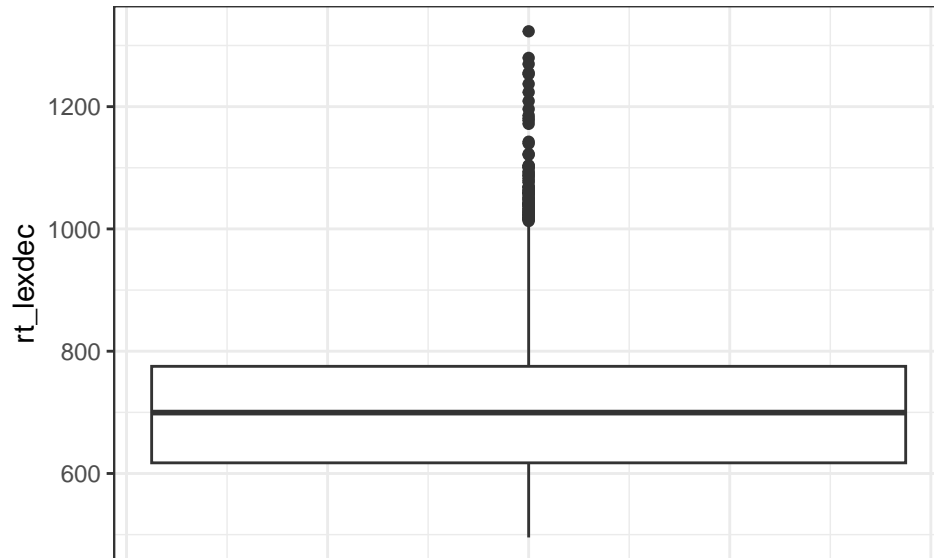


Abbildung 5: A boxplot for all observations of a continuous variable

- 
- für Boxplots verschiedener Gruppen: eine kategorische Variable entlang der anderen Achse (Abbildung 6)

```
df_eng |>
  ggplot(aes(x = age_subject, y = rt_lexdec)) +
  geom_boxplot() +
  theme_bw()
```

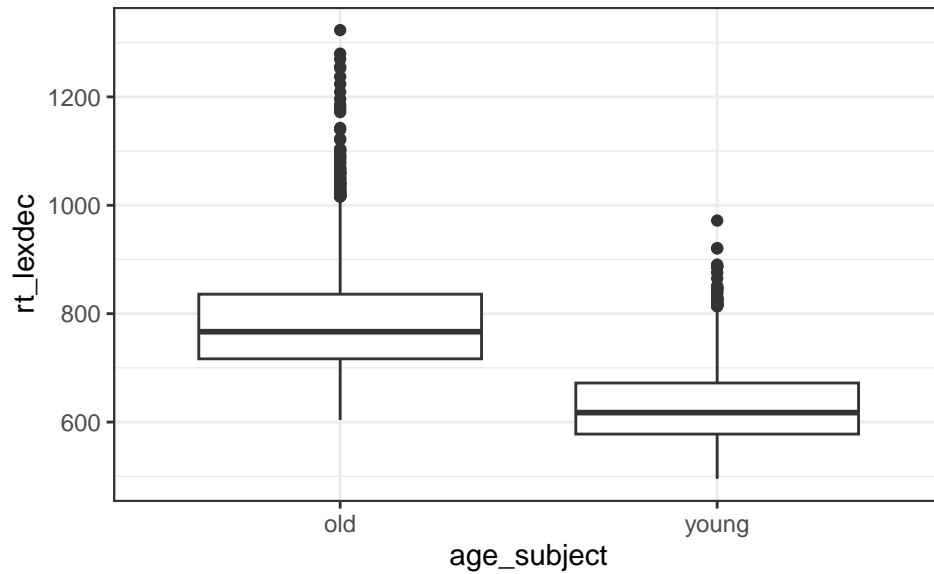


Abbildung 6: A boxplot for two groups

### 2.1.2 Gruppiertes Boxplot

- Wir können gruppierte Boxplots erstellen, um mehr Variablen zu visualisieren
  - einfach eine neue Variable mit `colour` oder `fill` ästhetisch zuordnen

```
df_eng |>
  ggplot(aes(x = age_subject, y = rt_lexdec, colour = word_category)) +
  geom_boxplot() +
  labs(
    x = "Age group",
    y = "LDT reaction time (ms)",
    color = "Word type"
  ) +
  scale_colour_colorblind() +
  theme_bw()
```



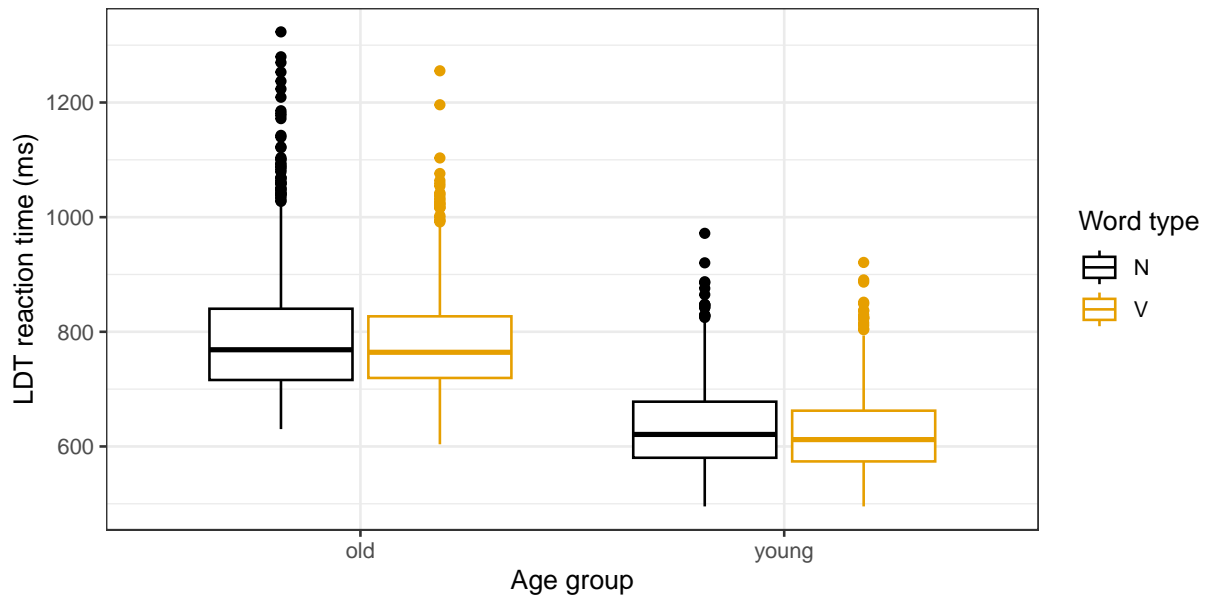


Abbildung 7: A grouped boxplot

### 3 Visualisierung des Mittelwerts

- In der Regel wollen wir auch den Mittelwert mit der Standardabweichung darstellen.
  - Wie können wir das tun?

#### 3.1 Fehlerbalkenplots

- Diese Diagramme bestehen aus 2 Teilen:
  - der Mittelwert, visualisiert mit `geom_point()`
  - ein Maß für die Streuung, visualisiert mit “`geom_errorbar()`”.
- für diesen Kurs werden wir die Standardabweichung verwenden
- Abbildung 8 ist das, was wir heute erzeugen werden

##### 3.1.1 Berechnung der zusammenfassenden Statistik

- müssen wir zunächst den Mittelwert und die Standardabweichung berechnen
  - gruppiert nach den Variablen, die wir visualisieren wollen

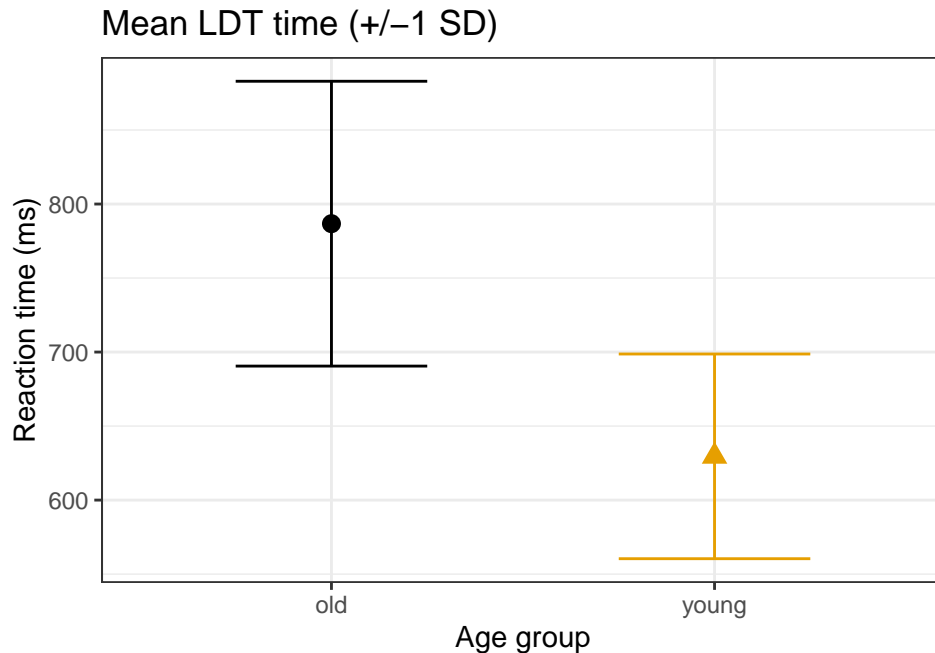


Abbildung 8: Errorbar plot of `df_eng` (body mass by `age_subject`)

- Wie kann man den Mittelwert und die Standardabweichung von `rt_lexdec` nach `age_subject` berechnen?

```
sum_eng <- df_eng |>
  summarise(mean = mean(rt_lexdec),
            sd = sd(rt_lexdec),
            N = n(),
            .by = age_subject) |>
  arrange(age_subject, age_subject)
```

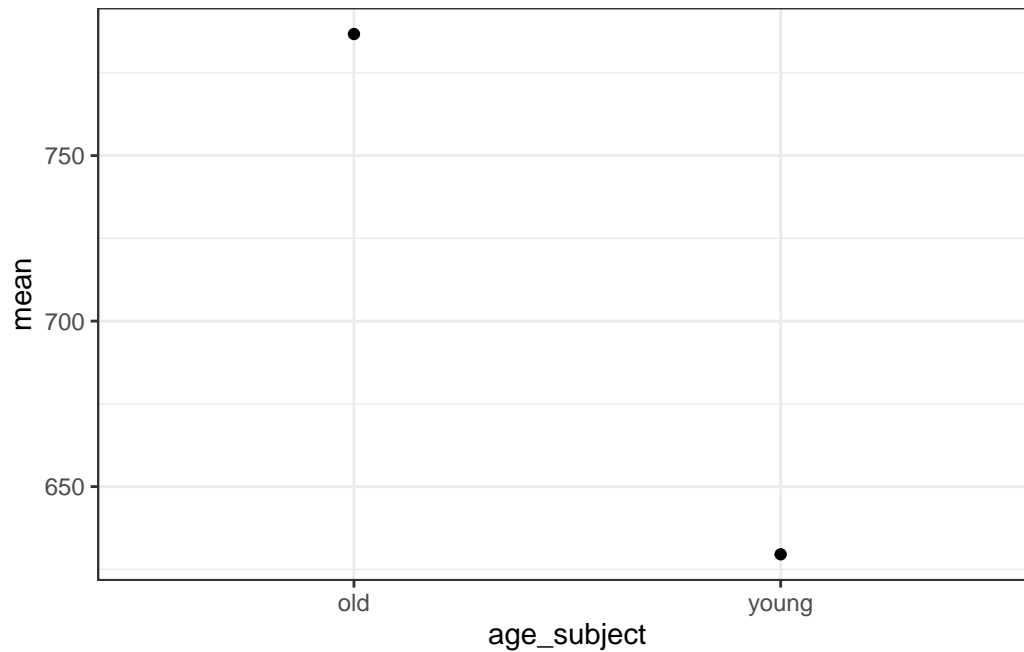
- Diese Zusammenfassung können wir dann in `ggplot()` mit den entsprechenden ästhetischen Zuordnungen und Geomen einfügen

### 3.1.2 Plotting mean

- Zunächst werden die Mittelwerte mit “`geom_point()`” dargestellt.

```
1 sum_eng |>
2   ggplot() +
3   aes(x = age_subject, y = mean) +
```

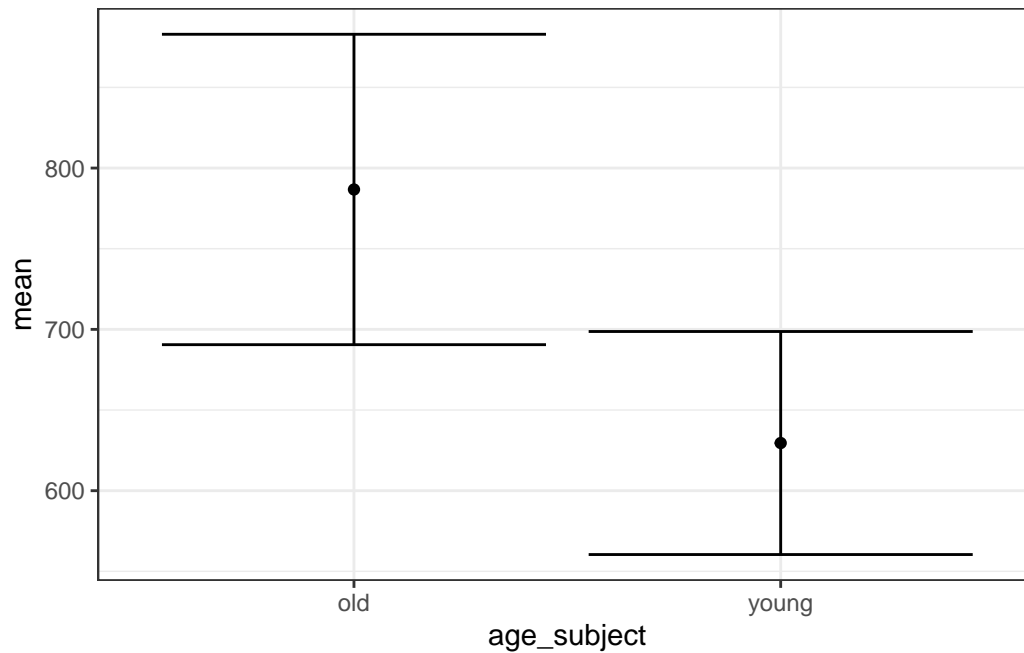
```
4 geom_point()
```



### 3.1.3 Hinzufügen von Fehlerbalken

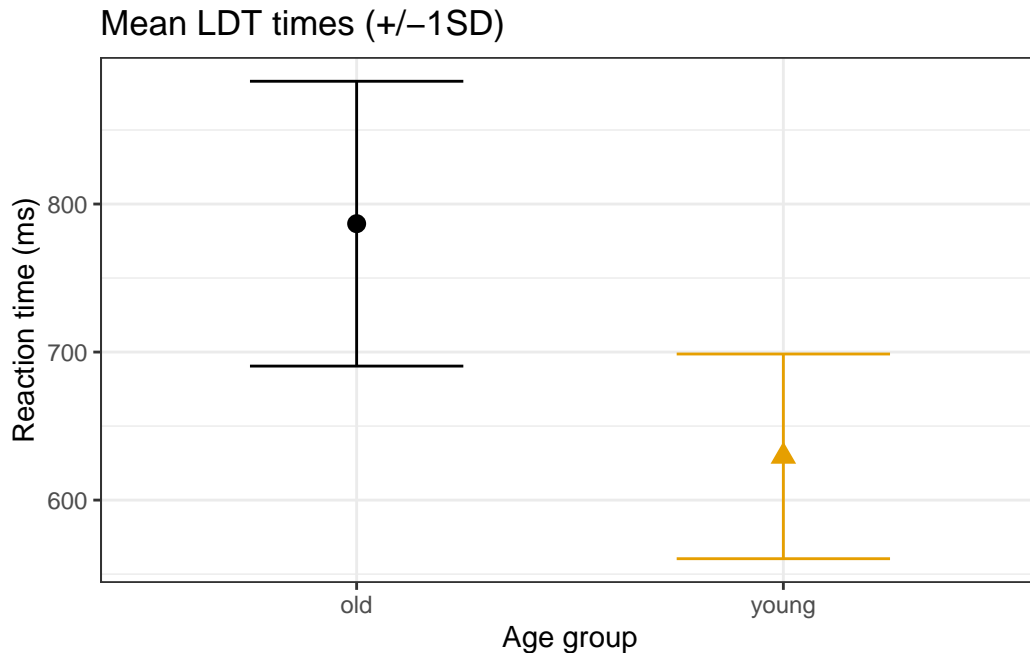
- Fügen wir nun unsere Fehlerbalken hinzu, die 1 Standardabweichung über und unter dem Mittelwert darstellen
- wir tun dies mit `geom_errorbar()`
  - nimmt `ymin` und `ymax` als Argumente
  - In unserem Fall sind dies `mean-/+sd`.

```
1 sum_eng |>
2   ggplot() +
3   aes(x = age_subject, y = mean) +
4   geom_point() +
5   geom_errorbar(aes(ymin = mean-sd,
6                     ymax = mean+sd))
```



- Wenn wir weitere Anpassungen hinzufügen, erhalten wir `?@fig-errorbar-custom`

```
sum_eng |>
  ggplot(aes(x = age_subject, y = mean, colour = age_subject, shape = age_subject)) +
  # geom_point(data = df_eng, alpha = .4, position = position_jitterdodge(.5), aes(x = age
  geom_point(size = 3) +
  geom_errorbar(width = .5, aes(ymin=mean-sd, ymax=mean+sd)) +
  labs(title = "Mean LDT times (+/-1SD)",
        x = "Age group",
        y = "Reaction time (ms)",
        color = "Age group"
  ) +
  scale_color_colorblind() +
  theme_bw() +
  theme(
    legend.position = "none"
  )
```



#### 4 Barplot von Mittelwerten: Bleiben Sie weg!

- Sie werden sehr oft Balkendiagramme von Mittelwerten sehen
  - aber es gibt viele Gründe, warum dies eine schlechte Idee ist!!
- Der Balkenplot hat ein schlechtes Daten-Tinten-Verhältnis, d.h. die Menge der Datentinte geteilt durch die Gesamtinte, die zur Erstellung der Grafik benötigt wird
  - Was ist, wenn es nur sehr wenige oder gar keine Beobachtungen in der Nähe von Null gibt? Wir verbrauchen eine Menge Tinte, wo es keine Beobachtungen gibt!
  - Außerdem deckt der Balken nur den Bereich ab, in dem die untere  *Hälfte*  der Beobachtungen liegt; ebenso viele Beobachtungen liegen über dem Mittelwert!
- Fehlerbalken allein sind keine Lösung: auch hier wird eine Menge Information verborgen
  - ein guter Grund, die Rohdatenpunkte *immer* zu visualisieren, unabhängig davon, welche zusammenfassende Darstellung Sie erstellen

#### Learning objects

In this section we learned how to...

- produce and interpret boxplots
- produce and interpret errorbar plots

## Hausaufgabe

### Boxplot mit Facette

1. Erzeugen Sie einen Plot namens `fig_boxplot`, der ein Boxplot der `df_eng` Daten ist, mit:
  - `age_subject` auf der x-Achse
  - `rt_naming` auf der y-Achse
  - `age_subject` als `colour` oder `fill` (wähle eine, es gibt keine falsche Wahl)
  - `Wort_Kategorie` in zwei Facetten mit `facet_wrap()` aufgetragen
  - die von Ihnen gewählte `theme_-`Einstellung (z.B. `theme_bw()`; für weitere Optionen siehe [hier](#))

### Errorbar plot

2. Versuchen Sie, Abbildung 9 zu reproduzieren. Hinweis: Sie werden die Variable `rt_naming` aus `df_eng` verwenden.

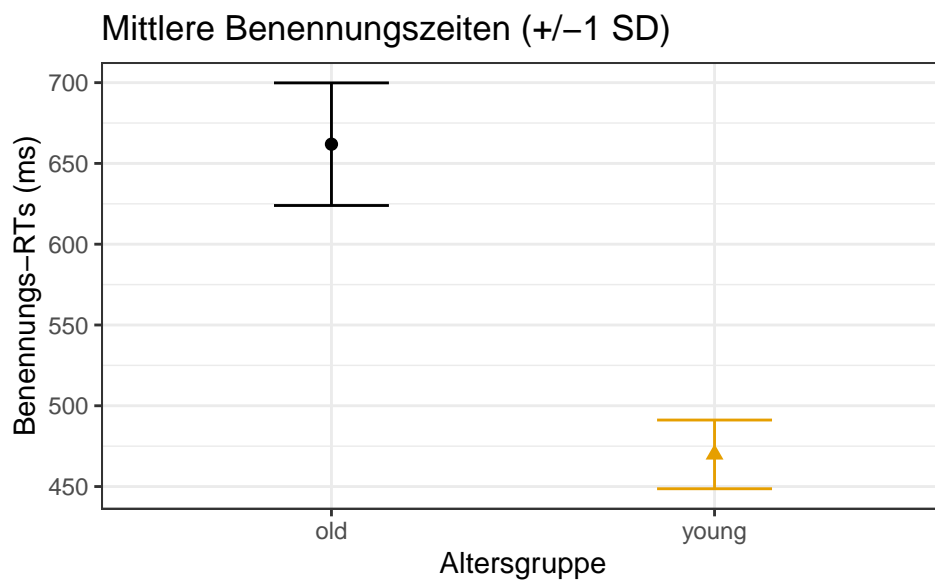


Abbildung 9: Plot to be reproduced

## Patchwork

- Verwenden Sie das Paket `patchwork`, um Ihren Boxplot und Ihre Fehlerbalkenplots nebeneinander darzustellen. Es sollte ungefähr so aussehen wie Abbildung 10. Hinweis: Wenn Sie die “tag-level” (“A” und “B”) zu den Plots hinzufügen möchten, müssen Sie `+ plot_annotation(tag_level = "A")` aus `patchwork` hinzufügen.

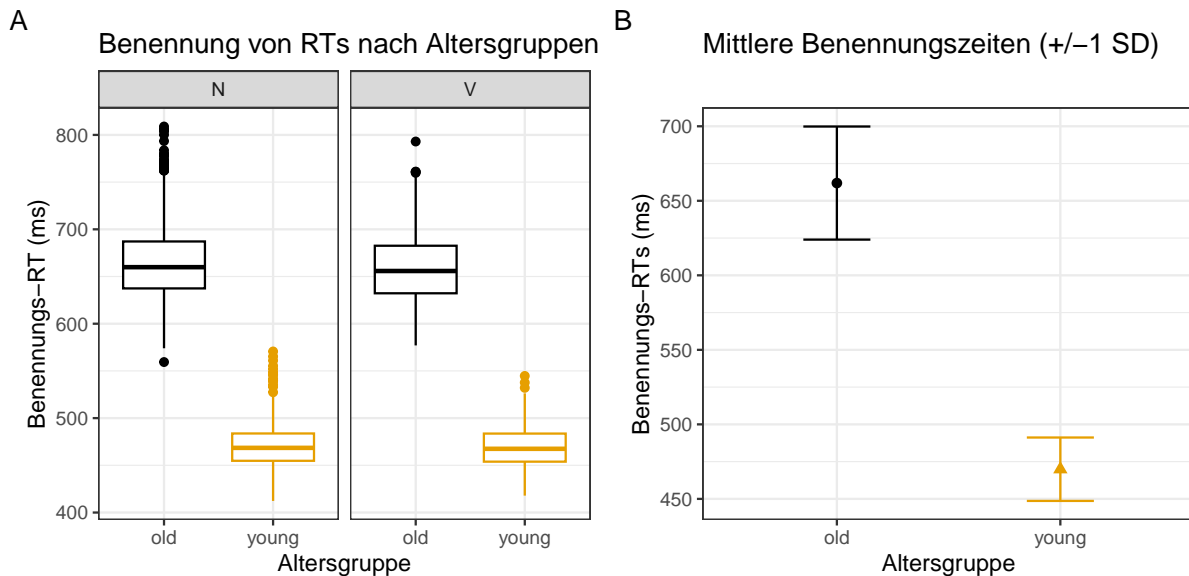


Abbildung 10: Combined plots with `patchwork`

## Session Info

Hergestellt mit R version 4.3.0 (2023-04-21) (Already Tomorrow) und RStudioversion 2023.9.0.463 (Desert Sunflower).

```
print(sessionInfo(), locale = F)
```

```
R version 4.3.0 (2023-04-21)
Platform: aarch64-apple-darwin20 (64-bit)
Running under: macOS Ventura 13.2.1
```

```
Matrix products: default
```

```
BLAS: /Library/Frameworks/R.framework/Versions/4.3-arm64/Resources/lib/libRblas.0.dylib
LAPACK: /Library/Frameworks/R.framework/Versions/4.3-arm64/Resources/lib/libRlapack.dylib;
```

attached base packages:

```
[1] stats      graphics  grDevices  utils      datasets  methods   base
```

other attached packages:

```
[1] magick_2.7.4    patchwork_1.1.3 ggthemes_4.2.4  janitor_2.2.0  
[5] here_1.0.1      lubridate_1.9.2 forcats_1.0.0   stringr_1.5.0  
[9] dplyr_1.1.3     purrr_1.0.2     readr_2.1.4     tidyr_1.3.0  
[13] tibble_3.2.1    ggplot2_3.4.3   tidyverse_2.0.0
```

loaded via a namespace (and not attached):

```
[1] utf8_1.2.3      generics_0.1.3  stringi_1.7.12  hms_1.1.3  
[5] digest_0.6.33   magrittr_2.0.3  evaluate_0.21    grid_4.3.0  
[9] timechange_0.2.0 fastmap_1.1.1   rprojroot_2.0.3  jsonlite_1.8.7  
[13] fansi_1.0.4     scales_1.2.1    cli_3.6.1        crayon_1.5.2  
[17] rlang_1.1.1     bit64_4.0.5     munsell_0.5.0    withr_2.5.0  
[21] yaml_2.3.7      parallel_4.3.0  tools_4.3.0      tzdb_0.4.0  
[25] colorspace_2.1-0 pacman_0.5.1     png_0.1-8        vctrs_0.6.3  
[29] R6_2.5.1        lifecycle_1.0.3 snakecase_0.11.0 bit_4.0.5  
[33] vroom_1.6.3     pkgconfig_2.0.3 pillar_1.9.0     gtable_0.3.4  
[37] glue_1.6.2      Rcpp_1.0.11     xfun_0.39        tidyselect_1.2.0  
[41] rstudioapi_0.14 knitr_1.44       farver_2.1.1     htmltools_0.5.5  
[45] labeling_0.4.3  rmarkdown_2.22  compiler_4.3.0
```

## Literaturverzeichnis

- Nordmann, E., McAleer, P., Toivo, W., Paterson, H., & DeBruine, L. M. (2022). Data Visualization Using R for Researchers Who Do Not Use R. *Advances in Methods and Practices in Psychological Science*, 5(2), 251524592210746. <https://doi.org/10.1177/25152459221074654>
- Wickham, H., Çetinkaya-Rundel, M., & Golemund, G. (2023). *R for Data Science* (2. Aufl.).
- Winter, B. (2019). Statistics for Linguists: An Introduction Using R. In *Statistics for Linguists: An Introduction Using R*. Routledge. <https://doi.org/10.4324/9781315165547>