

Data Wrangling 2

Data Tidying

Daniela Palleschi

Humboldt-Universität zu Berlin

2024-06-25

Lernziele

Heute werden wir lernen...

- über breite versus lange Daten
- wie man breite Daten länger macht
- wie man lange Daten breiter macht

Ressourcen

Die vorgeschlagenen Ressourcen für dieses Thema sind

- Kurs-Website: [Kapitel 9 \(Data Wrangling 2\)](#)
- [Kapitel 6 \(Data Tidying\)](#) in Wickham et al. (2023)
- [Kapitel 8 \(Data Tidying\)](#) in Nordmann & DeBruine (2022)

Einrichtung

Pakete

```
1 pacman::p_load(tidyverse,  
2                 here,  
3                 janitor)
```

Daten

- Wir verwenden den Datensatz `languageR_english.csv` (im Ordner `daten`)

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

①

②

③

④

⑤

⑥

- ① Bereinigen (d.h. *tidy*) von Variablennamen (von `janitor`)
- ② Zeilen nach `word` in ansteigender Reihenfolge anordnen (A-Z)
- ③ Variablen umbenennen...
- ④ `r_tlexdec` in `rt_lexdec` umbenennen
- ⑤ `r_tlexdec` in `rt_lexdec` umbenennen
- ⑥ nur die genannten Spalten behalten

‘Tidy’ Arbeitsablauf

- [Abbildung 1](#) zeigt einen Überblick über den typischen Data-Science-Prozess
 - Wir importieren unsere Daten, bereinigen sie und durchlaufen dann einen Zyklus aus Umwandlung, Visualisierung und Modellierung, bevor wir schließlich unsere Ergebnisse kommunizieren

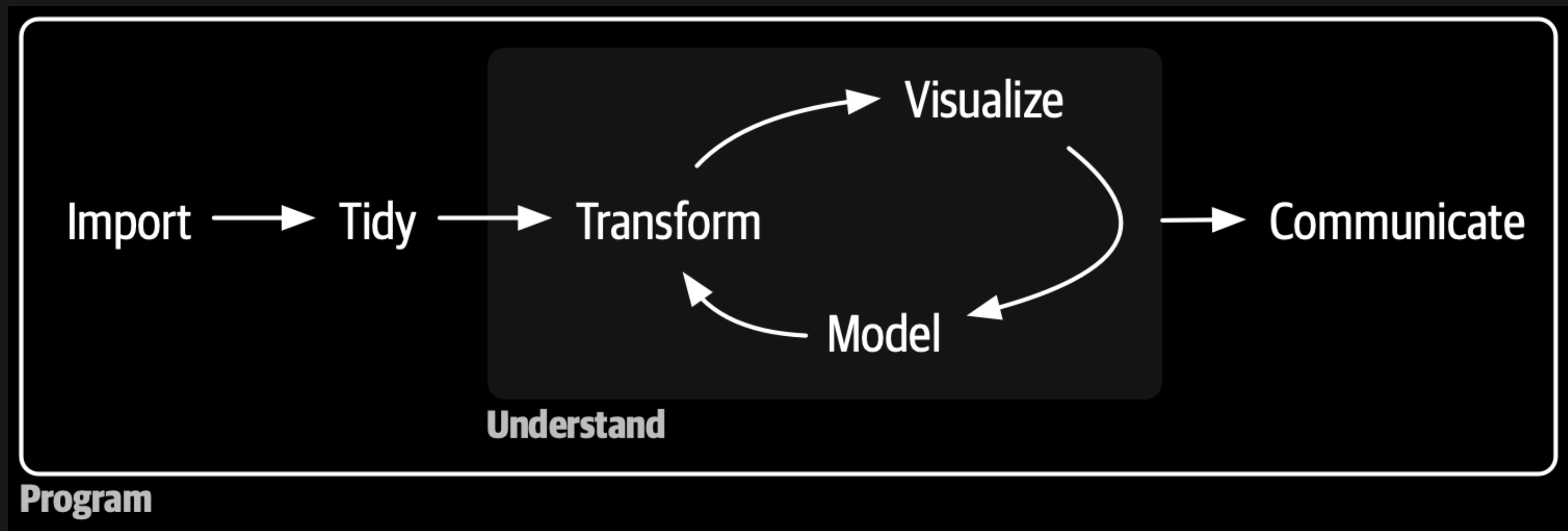


Abbildung 1: A model of the data science process from Wickham et al. (2023) (all rights reserved)

- Bisher haben wir gelernt, wie man
 - unsere Daten importieren (`readr::read_csv`)
 - Daten transformieren (Paket `dplyr`)
 - Daten zu visualisieren (Paket `ggplot`)
 - unsere Ergebnisse mit dynamischen Berichten zu kommunizieren (Quarto)
- aber wir haben bis jetzt nur aufgeräumte Daten gesehen
 - daher mussten wir den Schritt des “tidy” (Paket `tidyr`) noch nicht durchführen

‘Tidy’ Daten

- dieselben Daten können auf verschiedene Weise dargestellt werden
- Wir werden uns 3 Tabellen ansehen, die genau dieselben Daten in verschiedenen Formaten darstellen
- Die Tabellen zeigen die gleichen Werte von vier Variablen:
 - Land (**country**)
 - Jahr (**year**)
 - Bevölkerung (**population**)
 - Anzahl der Tuberkulosefälle (**cases**)
- Jeder Datensatz ordnet die Werte anders an
- überlegen Sie, welche Tabelle für Sie am einfachsten zu lesen ist

Tabelle 1: Version 1

| country | year | cases | population |
|-------------|------|--------|------------|
| Afghanistan | 1999 | 745 | 19987071 |
| Afghanistan | 2000 | 2666 | 20595360 |
| Brazil | 1999 | 37737 | 172006362 |
| Brazil | 2000 | 80488 | 174504898 |
| China | 1999 | 212258 | 1272915272 |
| China | 2000 | 213766 | 1280428583 |

Tabelle 2: Version 2

| country | year | type | count |
|-------------|------|------------|------------|
| Afghanistan | 1999 | cases | 745 |
| Afghanistan | 1999 | population | 19987071 |
| Afghanistan | 2000 | cases | 2666 |
| Afghanistan | 2000 | population | 20595360 |
| Brazil | 1999 | cases | 37737 |
| Brazil | 1999 | population | 172006362 |
| Brazil | 2000 | cases | 80488 |
| Brazil | 2000 | population | 174504898 |
| China | 1999 | cases | 212258 |
| China | 1999 | population | 1272915272 |
| China | 2000 | cases | 213766 |
| China | 2000 | population | 1280428583 |

Tabelle 3: Version 3

| country | year | rate |
|-------------|------|-------------------|
| Afghanistan | 1999 | 745/19987071 |
| Afghanistan | 2000 | 2666/20595360 |
| Brazil | 1999 | 37737/172006362 |
| Brazil | 2000 | 80488/174504898 |
| China | 1999 | 212258/1272915272 |
| China | 2000 | 213766/1280428583 |

Regeln für 'tidy' Daten

- Wahrscheinlich ist [Tabelle 1](#) für Sie am einfachsten zu lesen
 - sie folgt den drei Regeln für aufgeräumte Daten (visualisiert in [Abbildung 2](#)):

1. Jede Variable ist eine Spalte, jede Spalte ist eine Variable
2. Jede Beobachtung ist eine Zeile, jede Zeile ist eine Beobachtung
3. Jeder Wert ist eine Zelle, jede Zelle ist ein Einzelwert

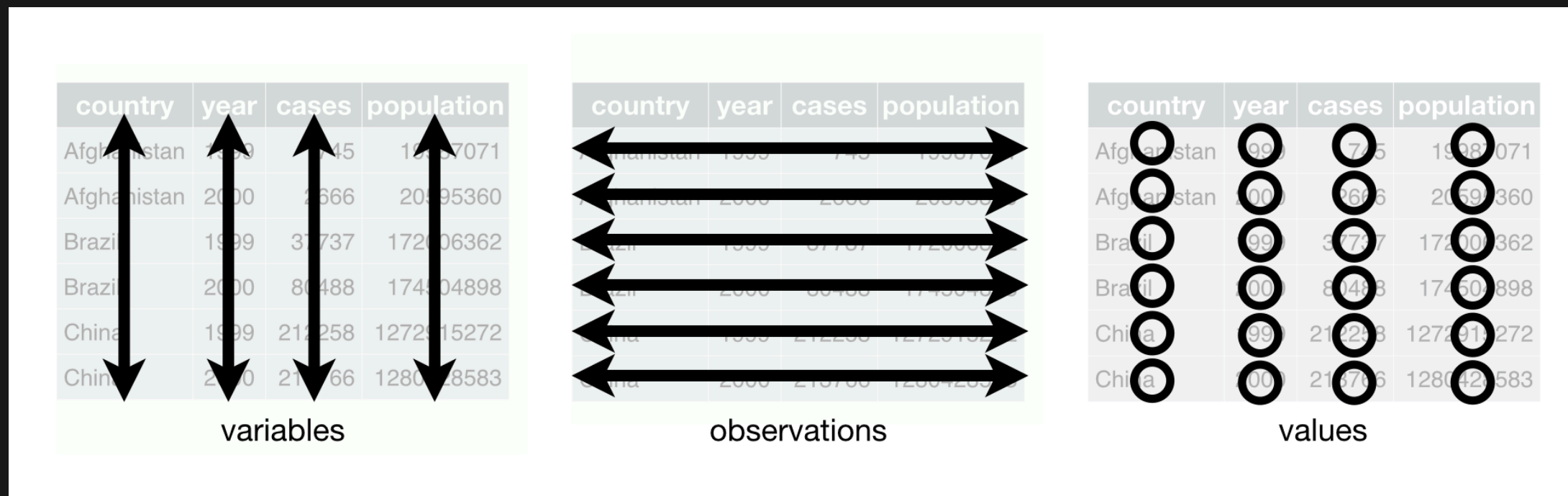


Abbildung 2: [Image source](#): Wickham et al. (2023) (all rights reserved)

Warum ‘tidy’ Daten?

“Glückliche Familien sind alle gleich; jede unglückliche Familie ist auf ihre eigene Art unglücklich.”

— Leo Tolstoy

“‘Tidy’ Datensätze sind alle gleich, aber jeder ‘untidy’ Datensatz ist auf seine eigene Weise unordentlich.”

— Hadley Wickham

Die Arbeit mit aufgeräumten Daten hat zwei wesentliche Vorteile:

1. Die Arbeit mit einer konsistenten Datenstruktur ermöglicht es uns, Konventionen zu übernehmen.

- Aufgeräumte Daten sind die allgemein vereinbarte Datenstruktur
- Konventionen/Werkzeuge basieren auf der Annahme dieser Struktur

2. Die vektorisierte Natur von R kann glänzen

- die meisten eingebauten R-Funktionen arbeiten mit *Vektorwerten* (und Spalten sind im Wesentlichen Vektoren)
- Alle Pakete im **tidyverse** sind darauf ausgelegt, mit aufgeräumten Daten zu arbeiten (z.B. **ggplot2** und **dplyr**)

Daten bereinigen (tidying)

- Umwandlung breiter Daten in lange Daten und langer Daten in breite Daten (neben anderen Schritten)
 - Ergebnis: aufgeräumte Daten (normalerweise)

‘Tidying’ Daten mit dem **tidyverse**

- Das Paket **tidyr** (aus **tidyverse**) hat zwei nützliche Funktionen zum Transponieren unserer Daten:
 - **pivot_longer()**: macht breite Daten länger
 - **pivot_wider()**: lange Daten breiter machen

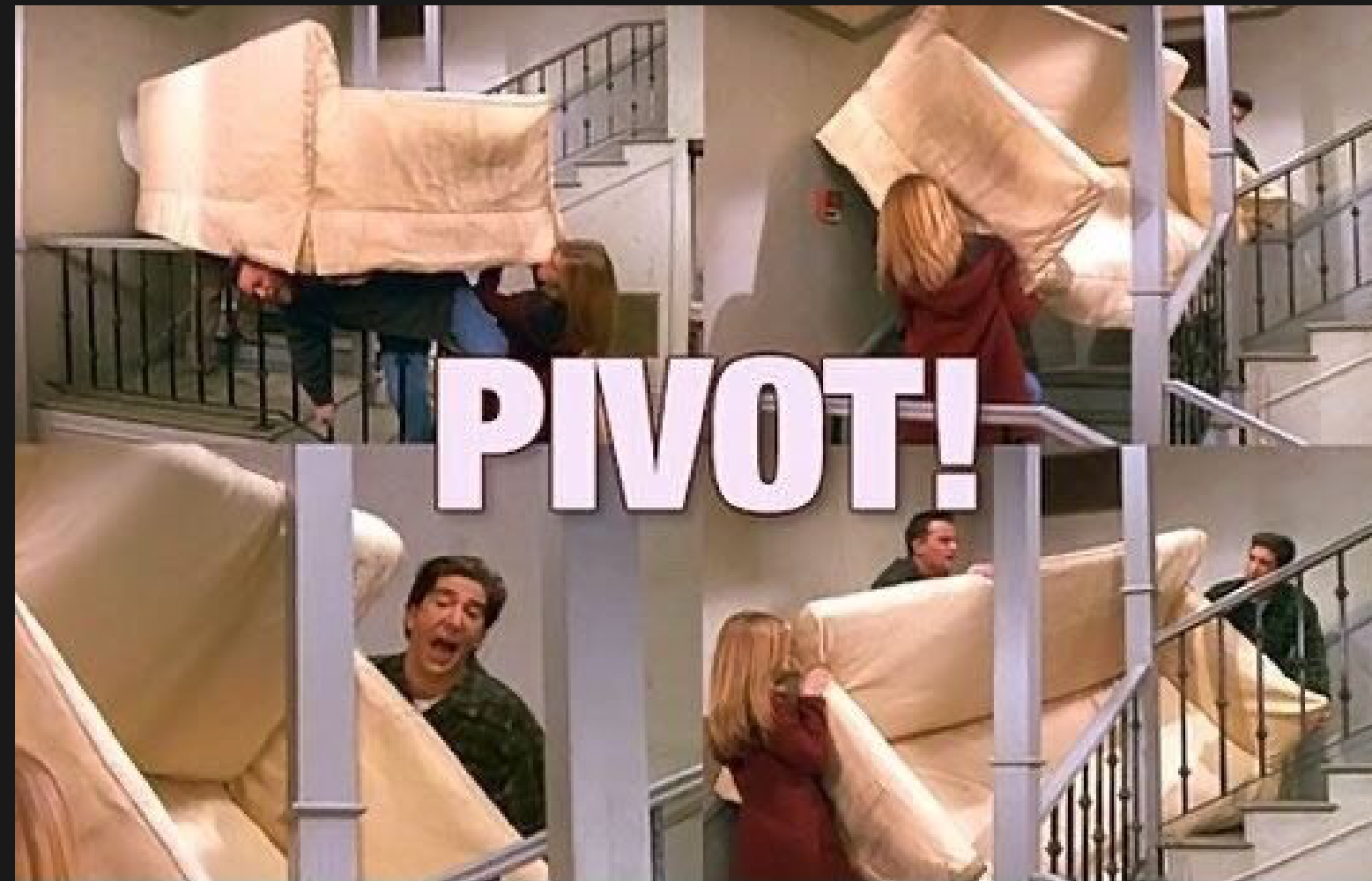


Abbildung 3: die berühmteste Verwendung des Wortes Pivot (zumindest für Millenials) ([Friends](#))

Breite versus lange Daten

- Wir müssen oft zwischen breiten und langen Datenformaten konvertieren, um verschiedene Arten von Zusammenfassungen oder Visualisierungen zu erstellen
- breite Daten: alle Beobachtungen zu einer Sache befinden sich in einer einzigen Zeile
 - ist *normalerweise* nicht aufgeräumt
- lange Daten: jede Beobachtung befindet sich in einer separaten Zeile
 - ist *normalerweise* aufgeräumt
- Beginnen wir mit dem typischsten Fall: Umwandlung breiter Daten in lange Daten

pivot_longer()

- im Datensatz `languageR_english.csv` (`df_eng`)
 - haben wir 4568 Beobachtungen (Zeilen)
 - Wir haben 5 Variablen (Spalten)
 - die Spalte `age_subject` gibt an, ob eine Beobachtung von einem Teilnehmer der Altersgruppe `old` oder `young` stammt
 - die Spalten `word` und `word_category` beschreiben Eigenschaften des Stimulus für eine bestimmte Beobachtung (d. h. das Wort)
 - die Spalte `rt_lexdec` enthält die Reaktionszeit für eine lexikalische Entscheidungsaufgabe
 - die Spalte `rt_naming` enthält die Antwortzeit für eine Wortbenennungsaufgabe

head(df_eng)

Tabelle 4: df_eng

| age_subject | word | word_category | rt_lexdec | rt_naming |
|-------------|------|---------------|-----------|-----------|
| young | ace | N | 623.61 | 456.3 |
| old | ace | N | 775.67 | 607.8 |
| young | act | V | 617.10 | 445.8 |
| old | act | V | 715.52 | 639.7 |
| young | add | V | 575.70 | 467.8 |
| old | add | V | 742.19 | 605.4 |

- Sind diese Daten in [Tabelle 4](#) aufgeräumt?
- Sind diese Daten zu breit oder zu lang?
- Wie können wir diese Daten länger machen?

Our goal

- wir wollen [Abbildung 4](#) produzieren

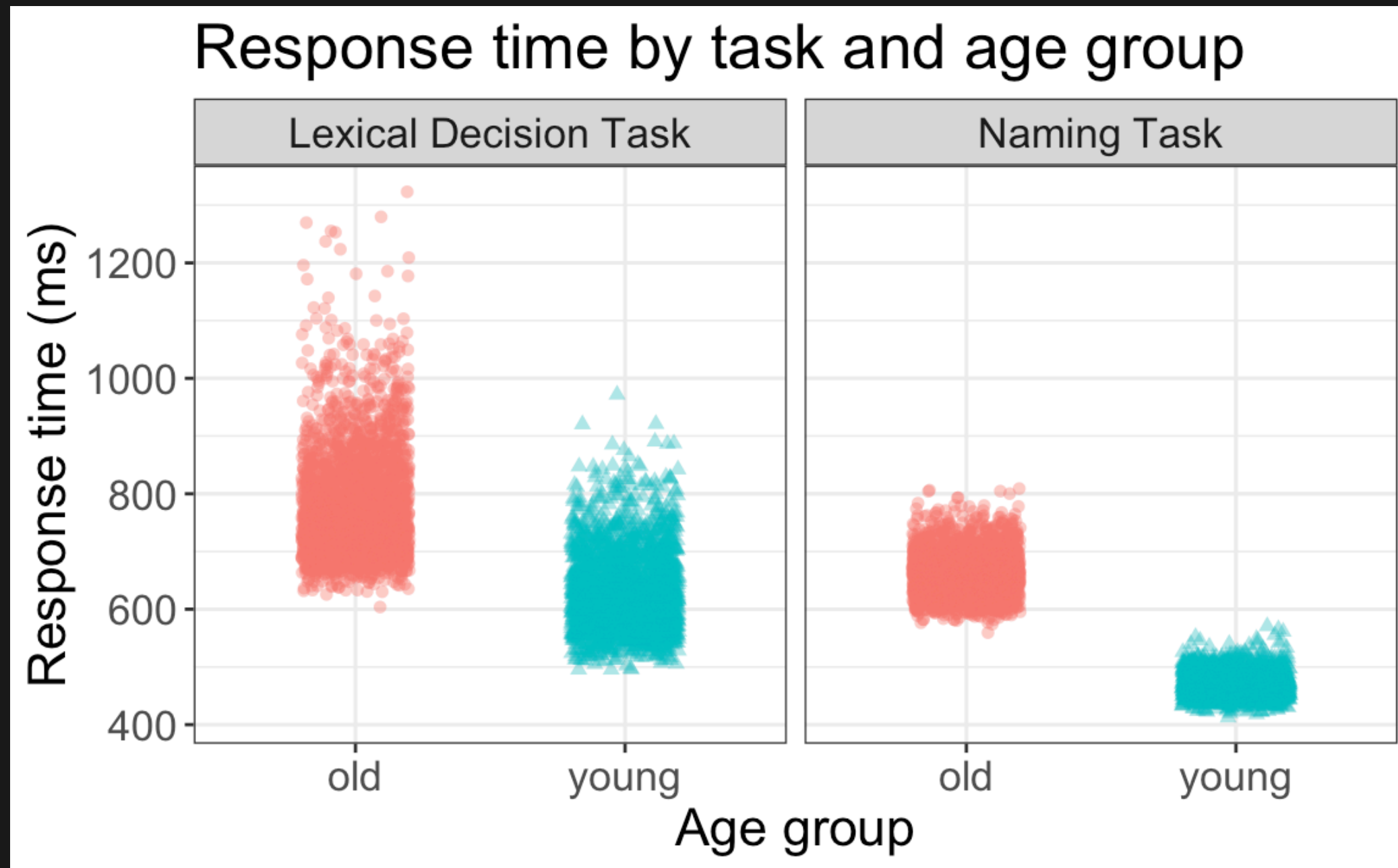


Abbildung 4: Our plot to be reproduced

- die beiden kontinuierlichen Variablen `rt_lexdec` und `rt_naming` erscheinen in Facetten
 - `facet_wrap()` nimmt eine *kategorische* Variable als Argument und erzeugt eine Facette für jede Kategorie
- wir brauchen also eine kategorische Variable, die die Ebenen `lexdec` und `naming` enthält
 - und eine *kontinuierliche* Variable, die die entsprechende Antwortzeit enthält

pivot_longer()

- Die Funktion `pivot_longer()` (von `tidyr`) konvertiert eine breite Datentabelle in ein längeres Format
 - wandelt die Namen der angegebenen Spalten in die Werte einer neuen kategorischen Spalte um
 - und kombiniert die Werte dieser Spalten in einer neuen Spalte

```
1 df_eng_long <-  
2   df_eng %>%  
3   pivot_longer(  
4     cols = starts_with("rt_"),  
5     names_to = "response",  
6     values_to = "time"  
7   )
```

```
1 df_eng_long <-  
2   df_eng %>%  
3   pivot_longer(  
4     cols = starts_with("rt_"),  
5     names_to = "response",  
6     values_to = "time"  
7   )
```

①

②

③

④

⑤

⑥

- ① Erstellen Sie ein neues Objekt namens **df_eng_long**, das...
- ② **df_eng**, und dann
- ③ mache es länger
- ④ indem du Spalten (**col =**) nimmst, die mit **rt_** beginnen
- ⑤ und eine Variable namens **response** erstellen, die die Namen aus **cols** enthält (**names_to =**)
- ⑥ und eine Variable namens **time** erstellen, die die Werte aus **cols** enthält (**values_to =**)

```
1 df_eng_long |> head()
```

```
# A tibble: 6 × 5
```

| | age_subject | word | word_category | response | time |
|---|-------------|-------|---------------|-----------|-------|
| | <chr> | <chr> | <chr> | <chr> | <dbl> |
| 1 | young | ace | N | rt_lexdec | 624. |
| 2 | young | ace | N | rt_naming | 456. |
| 3 | old | ace | N | rt_lexdec | 776. |
| 4 | old | ace | N | rt_naming | 608. |
| 5 | young | act | V | rt_lexdec | 617. |
| 6 | young | act | V | rt_naming | 446. |

- Vergleichen wir die Beobachtungen für die Wörter **ace** und **act** in
 - **df_eng** (Tabelle 5)
 - **df_eng_longer** (Tabelle 6)

Tabelle 5: **df_eng**

| age_subject | word | rt_lexdec | rt_naming |
|-------------|------|-----------|-----------|
| young | ace | 623.61 | 456.3 |
| old | ace | 775.67 | 607.8 |
| young | act | 617.10 | 445.8 |
| old | act | 715.52 | 639.7 |

Tabelle 6: **df_eng |> pivot_longer(...)**

| age_subject | word | response | time |
|-------------|------|-----------|--------|
| young | ace | rt_lexdec | 623.61 |
| young | ace | rt_naming | 456.30 |
| old | ace | rt_lexdec | 775.67 |
| old | ace | rt_naming | 607.80 |
| young | act | rt_lexdec | 617.10 |
| young | act | rt_naming | 445.80 |
| old | act | rt_lexdec | 715.52 |
| old | act | rt_naming | 639.70 |

- die beiden Tabellen enthalten genau die gleichen Informationen
 - 8 Werte für die Antwortzeit:
 - 4 für `rt_lexdec`
 - 4 für `rt_naming`
- Dies ist eine wichtige Erkenntnis: Wir haben keine Daten oder Beobachtungswerte geändert, sondern lediglich die Organisation der Datenpunkte neu strukturiert.

Plotten unserer 'tidy' Daten

- Versuchen wir nun, unser Diagramm zu erstellen:
 - `age_subject` auf der x-Achse
 - `time` auf der y-Achse
 - Kategorien `response` in Facetten

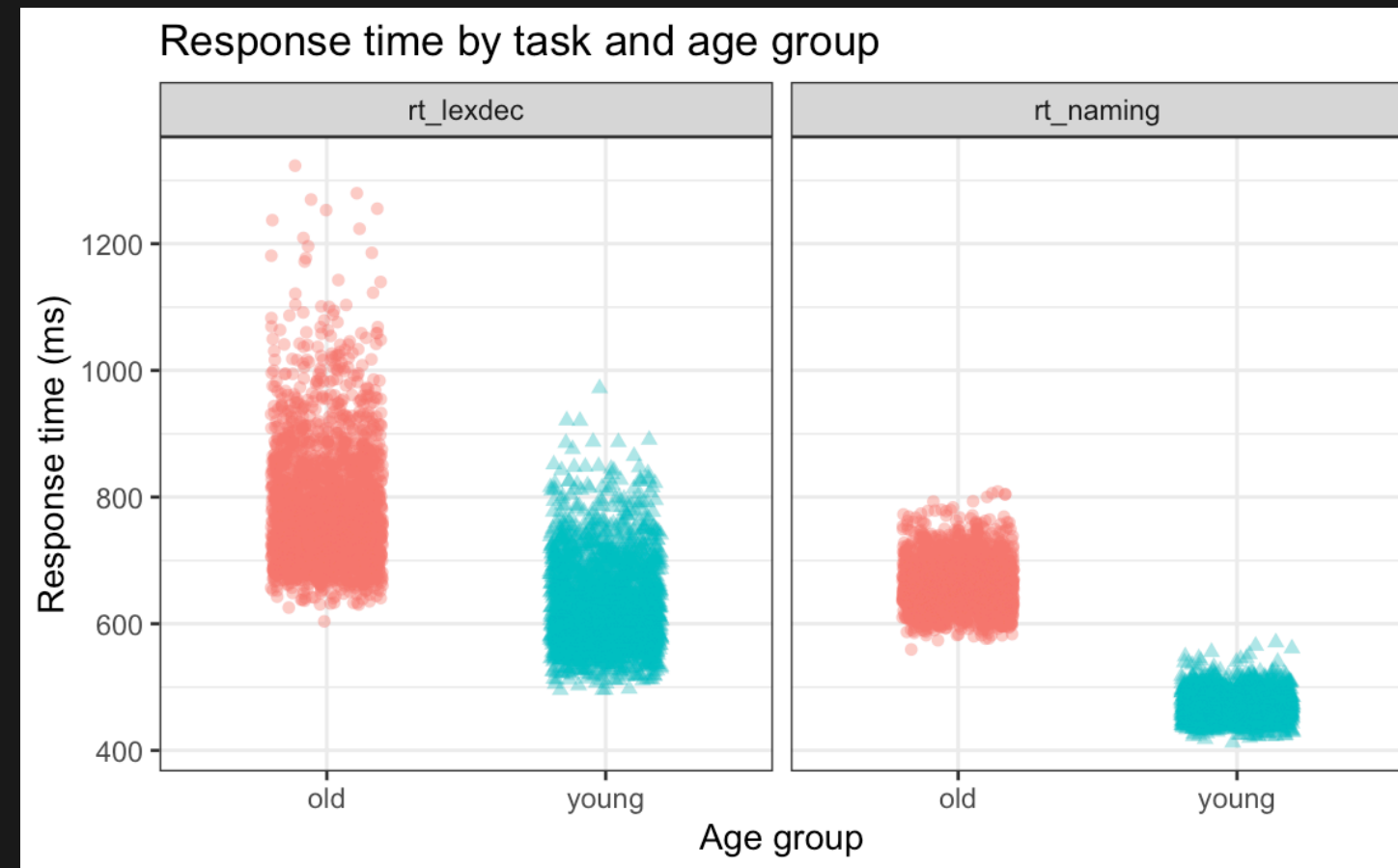


Abbildung 5: Response times per age group for the lexical decision task vs. naming task

Aufgabe 1: Tidy data

Beispiel 1

Abbildung 5 neu erstellen.

pivot_wider()

- Es kommt häufiger vor, dass man seine Daten verlängern will (man nimmt Spalten und macht aus deren Werten neue Zeilen)
 - aber manchmal möchte man seine Daten auch verbreitern (man nimmt Zeilen und verwandelt ihre Werte in neue Spalten)
- Die `tidyr`-Funktion `pivot_wider()` macht Datensätze *breiter*, indem sie Spalten vergrößert und Zeilen reduziert.
 - Dies ist hilfreich, wenn eine Beobachtung über mehrere Zeilen verteilt ist.
- Lassen Sie uns versuchen, `df_eng` breiter zu machen
 - Wir könnten zum Beispiel eine einzige Zeile pro *Wort* haben
 - mit einer einzigen Variablen für die Antwort des `young` Probanden und die Antwort des `old` Probanden

`pivot_wider()`

- `pivot_wider` nimmt ähnliche Argumente wie `pivot_longer()`, mit einigen leichten Unterschieden:
 - `id_cols` (optional): identifizierende Spalten (welche Spalten identifizieren jede Beobachtung eindeutig?)
 - `names_from`: wie soll die neue Spalte heißen, die die vorherigen Spaltennamen enthält (muss eine kategoriale Variable sein)?
 - `names_prefix` (optional): Präfix für die neuen Spaltennamen (optional)
 - `values_von`: neue Spaltenwerte

- lassen Sie uns zwei neue Variablen erstellen, die ihre Namen von **age_subject** und ihre Werte von **rt_lexdec** übernehmen

```
1 df_eng_wide <-  
2   df_eng %>%  
3   select(-rt_naming) |>  
4   pivot_wider(  
5     names_from = age_subject,  
6     values_from = rt_lexdec,  
7     names_prefix = "lexdec_"  
8   )
```

①

②

③

- ① neue Spaltennamen unter Verwendung der Werte in **age_subject** erstellen
- ② Erstelle neue Beobachtungswerte aus **rt_lexdec**
- ③ Hinzufügen von **lexdec_** am Anfang der neuen Spaltennamen

- Vergleichen wir die Beobachtungen für die Wörter **ace** und **act** in
 - **df_eng** (Tabelle 5)
 - **df_eng_longer** (Tabelle 6)

Tabelle 7: **df_eng**

| age_subject | word | word_category | rt_lexdec |
|-------------|------|---------------|-----------|
| young | ace | N | 623.61 |
| old | ace | N | 775.67 |
| young | act | V | 617.10 |
| old | act | V | 715.52 |

Tabelle 8: **df_eng_wide**

| word | word_category | lexdec_young | lexdec_old |
|------|---------------|--------------|------------|
| ace | N | 623.61 | 775.67 |
| act | V | 617.10 | 715.52 |

- Auch hier haben wir keine Daten oder Beobachtungswerte geändert, sondern lediglich die Anordnung der Datenpunkte neu strukturiert.

Eindeutige Werte




- Wir haben `rt_naming` entfernt, weil es auch einen eindeutigen Wert pro Wort pro Altersgruppe hat
- wir ändern nur die Breite und führen `NA`-Werte für `lexdec_young` für alte Themen und `NA`-Werte für `lexdec_old` für junge Themen ein
- Hätten wir sie nicht entfernt, sähen unsere ersten 6 Zeilen wie [Tabelle 9](#) aus
 - Vergleichen Sie dies mit der Ausgabe in [Tabelle 8](#), sehen Sie den Unterschied?

Tabelle 9: Wider data with missing values

| word | word_category | rt_naming | lexdec_young | lexdec_old |
|------|---------------|-----------|--------------|------------|
| ace | N | 456.3 | 623.61 | NA |
| ace | N | 607.8 | NA | 775.67 |
| act | V | 445.8 | 617.10 | NA |
| act | V | 639.7 | NA | 715.52 |

Lernziele

Heute haben wir gelernt...

- über breite und lange Daten 
- wie man breite Daten länger macht 
- wie man lange Daten breiter macht 

Hausaufgaben

Anhang 2 auf der Website des Kurses.

Session Info

Hergestellt mit R version 4.4.0 (2024-04-24) (Puppy Cup) und RStudioversion 2023.9.0.463 (Desert Sunflower).

```
1 sessionInfo()
```

```
R version 4.4.0 (2024-04-24)
```

```
Platform: aarch64-apple-darwin20
```

```
Running under: macOS Ventura 13.2.1
```

```
Matrix products: default
```

```
BLAS: /Library/Frameworks/R.framework/Versions/4.4-arm64/Resources/lib/libRblas.0.dylib
```

```
LAPACK: /Library/Frameworks/R.framework/Versions/4.4-arm64/Resources/lib/libRlapack.dylib; LAPACK version 3.12.0
```

```
locale:
```

```
[1] en_US.UTF-8/en_US.UTF-8/en_US.UTF-8/C/en_US.UTF-8/en_US.UTF-8
```

```
time zone: Europe/Berlin
```

```
tzcode source: internal
```

```
attached base packages:
```

```
stats, graphics, grDevices, utils, datasets, MASS, R6, Rcpp, RcppEigen, RcppArmadillo, RcppParallel, RcppGSL, RcppGL, RcppHDF, RcppMath, RcppRNG, RcppS4, RcppShiny, RcppTOML, RcppUnit, RcppV8, RcppX11, RcppX11R, RcppX11R2, RcppX11R3, RcppX11R4, RcppX11R5, RcppX11R6, RcppX11R7, RcppX11R8, RcppX11R9, RcppX11R10, RcppX11R11, RcppX11R12, RcppX11R13, RcppX11R14, RcppX11R15, RcppX11R16, RcppX11R17, RcppX11R18, RcppX11R19, RcppX11R20, RcppX11R21, RcppX11R22, RcppX11R23, RcppX11R24, RcppX11R25, RcppX11R26, RcppX11R27, RcppX11R28, RcppX11R29, RcppX11R30, RcppX11R31, RcppX11R32, RcppX11R33, RcppX11R34, RcppX11R35, RcppX11R36, RcppX11R37, RcppX11R38, RcppX11R39, RcppX11R40, RcppX11R41, RcppX11R42, RcppX11R43, RcppX11R44, RcppX11R45, RcppX11R46, RcppX11R47, RcppX11R48, RcppX11R49, RcppX11R50, RcppX11R51, RcppX11R52, RcppX11R53, RcppX11R54, RcppX11R55, RcppX11R56, RcppX11R57, RcppX11R58, RcppX11R59, RcppX11R60, RcppX11R61, RcppX11R62, RcppX11R63, RcppX11R64, RcppX11R65, RcppX11R66, RcppX11R67, RcppX11R68, RcppX11R69, RcppX11R70, RcppX11R71, RcppX11R72, RcppX11R73, RcppX11R74, RcppX11R75, RcppX11R76, RcppX11R77, RcppX11R78, RcppX11R79, RcppX11R80, RcppX11R81, RcppX11R82, RcppX11R83, RcppX11R84, RcppX11R85, RcppX11R86, RcppX11R87, RcppX11R88, RcppX11R89, RcppX11R90, RcppX11R91, RcppX11R92, RcppX11R93, RcppX11R94, RcppX11R95, RcppX11R96, RcppX11R97, RcppX11R98, RcppX11R99, RcppX11R100
```

Literaturverzeichnis

Nordmann, E., & DeBruine, L. (2022). *Applied Data Skills*. Zenodo. <https://doi.org/10.5281/zenodo.6365078>

Wickham, H., Çetinkaya-Rundel, M., & Grolemund, G. (2023). *R for Data Science* (2. Aufl.).