

# Mini-Workshop (Structural) Topic Models

Marko Bachl

Sommersemester 2020 | IJK Hannover



# Contents



# Chapter 1

## Überblick

### 1.1 Inhalt des virtuellen Mini-Workshops

- In diesem Mini-Workshop erläutere ich das praktische Vorgehen einer Datenanalyse mit *Structural Topic Models*. Wir behandeln die folgenden Schritte im Analyseprozess:
  - Modellspezifikation
  - Modellvergleich zur Auswahl eines geeigneten Modells
  - Interpretation der Topics im finalen Modell
  - Darstellung der Ergebnisse
  - Weitere Analysen
    - \* Identifikation verwandter Themen
    - \* Zusammenhänge der Themenprävalenz mit Kovariaten.
- Wir verwenden das Paket `{stm}` (?) zum Schätzen von Topic Models. Für die Variante der *Structural Topic Models* und die Implementation in diesem Paket sprechen *für mich* die folgenden Gründe
  - Gute Integration mit *R* und Paketen, die ich für die Arbeit mit Text-Daten verwende (insbesondere `{quanteda}` und `{tidytext}`)
  - Gute ergänzende Pakete zur Arbeit mit den Modellen (insbesondere `{stminsights}`)
  - Vergleichsweise schnelle Modellschätzung auch mit großen Datensätzen
  - Direktes Schätzen von Zusammenhängen von Topics mit Kovariaten
  - Initialisieren der Modellschätzung mit dem Spectral Algorithmus
  - Recht weit verbreitet in einem Feld, in dem ich viel lese (Politische Kommunikation nach einem weitem Verständnis)
- Die Darstellung basiert auf einer Analyse, die ich gemeinsam mit Elena Link durchgeführt habe. Wir untersuchten, wie das Thema Impfen in Online-Foren für Eltern diskutiert wurde. Wir verwenden aber nur einen *nicht repräsentativen* Ausschnitt aus dem Material, um die notwendige

Rechenleistung und -zeit zu verringern.

- Einen Preprint zur Analyse könnt ihr hier lesen: Vaccine-related Discussions in Online Communities for Parents. A Quantitative Overview.
- Die Dokumentation zur Studie ist hier verfügbar: [https://bachl.github.io/vaccine\\_discussions/](https://bachl.github.io/vaccine_discussions/). Daten und Analyse-Skripts gibt es im OSF. Dort werden auch die Datenerhebung mittels Web-Scraping und die Datenaufbereitung erläutert. Diese Inhalte sind *nicht* Teil dieses Workshops. Wenn ihr Fragen dazu habt, dürft ihr sie natürlich stellen.

## 1.2 Welche Inhalte wir *nicht* behandeln

- Auch wenn das im direkten Vergleich mit dem Parallel-Angebot zu Panel Data Analysis (meine Ausführlichkeit dort sind ein Grund für die spätere Lieferung dieser Materialien) enttäuschend sein mag: Die Inhalte in diesem Mini-Workshop entsprechen in ihrem Umfang wirklich nur dem, was ich zu Beginn des Digital-Semesters geplant und angekündigt hatte. Der Mini-Workshop ersetzt keine tiefer gehende Einarbeitung in die Methode, sondern ist als ein Einstieg zu verstehen.
- Wir behandeln hier keine theoretischen, statistischen oder auf die Software-Implementierung der Modellschätzung bezogenen Fragen. Die Grundlagen dazu können aus den Texten im LMS entnommen werden (??).
- Es gibt neben `{stm}` viele andere Implementationen in *R* und ihn anderer Software. Gefühlt gibt es alle 6 Monate eine neue Variante von Topic Models, alle 3 Monate eine neue Implementierung und jeden Monat ein Paket mit zusätzlichen Tools für die Arbeit mit Topic Models. Meine Entscheidung für `{stm}` ist keine informierte Entscheidung gegen andere Varianten, Implementierungen und Tools. Dieser Workshop ist keine Aufforderung, ausschließlich `{stm}` zu nutzen. Informiert euch gegebenenfalls selbst über Software-Lösungen, die für eure Bedürfnisse geeignet sind.
- Dieser Mini-Workshop ist kein *R*-Tutorial. Wenn ihr Interesse habt, *R*-Kenntnisse zu erwerben und zu vertiefen, empfehle ich R4DS.
- Dieser Mini-Workshop ist keine allgemeine Einführung in die computergestützte Inhaltsanalyse. Wenn ihr allgemein mit *R* arbeiten möchtet, empfehle ich zu diesem Thema die Einführung von Cornelius Puschmann.

## 1.3 Aufbau des Workshops

- Inhaltlicher Aufbau: Siehe Kapitel-Gliederung

## Material

- Dieses Dokument + R Skripte: (Hoffentlich) mehr oder weniger selbsterklärendes Material
  - Kuratierte Form ist dieses HTML-Dokument
  - Es gibt auch ein PDF, das ich aber nicht formatiert habe
- Daten: Ein Ausschnitt auf den Daten der oben genannten Beispielstudie. Eine genauere Beschreibung folgt im nächsten Abschnitt.
- Screencast: Zu einigen Analyseschritten stelle ich Screencasts zur Verfügung. Diese sind größtenteils ergänzend gedacht. Bis auf wenige Ausnahmen sollte das schriftliche Material selbsterklärend sein.
- Übungen: Zu einigen Analysen gibt es Übungsaufgaben.
  - XXX

## Pakete

Wir verwenden die folgenden Pakete

```
if (!require("pacman")) install.packages("pacman")
pacman::p_load(tidyverse, stm, stmights, tidytext, quanteda, lubridate, knitr,
               tictoc, furrr)
theme_set(theme_bw()) # ggplot theme

tibble(package = c("R", sort(pacman::p_loaded())) %>% mutate(version = map_chr(package,
~as.character(pacman::p_version(package = .x))) %>% knitr::kable()
```

package	version
R	3.6.2
dplyr	0.8.4
forcats	0.4.0
furrr	0.1.0
future	1.16.0
ggplot2	3.3.1
knitr	1.28
lubridate	1.7.4
pacman	0.5.1
purrr	0.3.3
quanteda	2.0.0
readr	1.3.1
stm	1.3.5
stminsights	0.4.0
stringr	1.4.0
tibble	2.1.3
tictoc	1.0
tidyr	1.0.2
tidytext	0.2.3
tidyverse	1.3.0



## Chapter 2

# Beispiel-Daten und Aufbereitung

### 2.1 Laden der Daten und Übersicht

- Wir verwenden einen Ausschnitt der Daten aus der Beispielstudie. Konkret handelt es sich um Posts mit dem Suchwort *impf*, die zwischen dem 1. Mai 2016 und dem 8. Juli 2019 im Elternforum Urbia veröffentlicht wurden. Ausgeschlossen wurden unter anderem
  - sehr kurze Posts (weniger als 19 Wörter)
  - Posts mit dem Wort *schimpf*
  - Posts zur Impfung von Haustieren (nach einem kurzen Diktionär)
- Die Dokumentation zur Studie gibt weitere Informationen zur Erhebung und Bereinigung der Rohdaten.
- Diese Daten können aus Copyright- und Privacy-Gründen nicht auf GitHub veröffentlicht werden. Ich habe Sie daher im LMS hochgeladen. Bitte ladet die ZIP-Datei herunter.
  - Wenn ihr sie mit dem Code aus dem Repository integrieren wollt, müsst ihr sie in den Ordner “data” unter “R” entpacken.

```
# Laden der Daten
```

```
d = read_rds("R/data/example_data.rds")
```

```
d %>%
```

```
print(n = 5)
```

```
## # A tibble: 12,369 x 5
```

```
##   post                                author   postdate      wc thread_title
```

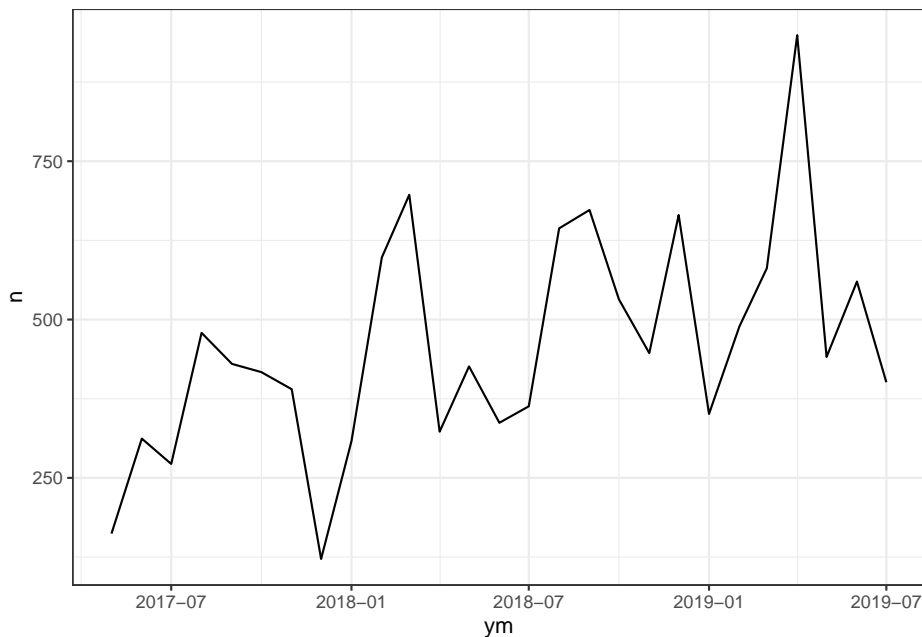
```
##   <chr>                                <chr>    <date>      <int> <chr>
```

```
## 1 Wenn Impfungen zu Todesfäll~ zwerg-b~ 2018-04-06    26 HPV-Impfung
```

```
## 2 Hallo Moni Danke für deine ~ Inaktiv  2017-06-03    21 Warum so oft Scheidenp~
```

```
## 3 Hallo ja sind glaube ich dr~ danerl 2017-06-05 42 Warum so oft Scheidenp~
## 4 Guten Morgen, gibt es hier ~ butterf~ 2017-05-14 133 Impfung Deutschland/Üs~
## 5 In Österreich wird im 3., 5~ butterf~ 2017-05-15 68 Impfung Deutschland/Üs~
## # ... with 1.236e+04 more rows
```

```
d %>%
  mutate(ym = round_date(postdate, "month")) %>%
  count(ym) %>%
  ggplot(aes(ym, n)) + geom_line()
```



```
d %>%
  pull("wc") %>%
  summary()
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##      20      37      60      85    101    2493
```

- Der Datensatz besteht aus 12,369 Posts.
  - Die Variable **post** enthält den vollen Text des Posts.
  - Die Variable **author** enthält den Accountnamen, von dem der Post abgegeben wurde.
  - Die Variable **date** enthält den Tag der Veröffentlichung.
  - Die Variable **wc** enthält die Zahl der Wörter des Posts.
  - Die Variable **thread\_title** enthält den Titel des Diskussions-Threads.
- Pro Monat sind zwischen ca. 120 und 1.000 Posts in unserer Stichprobe.
- Typische Posts haben einen Umfang von zwischen 40 und 100 Wörtern

(Zur Erinnerung: Sehr kurze Post wurden bereits ausgeschlossen).

## 2.2 Aufbereitung für das Schätzen der Topic Models

- Grundsätzlich gilt: Die verschiedenen Schritte bei der Aufbereitung des Text-Korpus kann die Ergebnisse wesentlich beeinflussen (??). Aber ist es häufig sehr schwierig, theoretisch informierte Entscheidungen zu treffen, da
  - unsere Theorien fast immer zu vage sind, um etwas über konkrete, manifeste Eigenschaften der Texte auszusagen
  - es schwer ist, die Folge einer Entscheidung für das technische Schätzen der Modelle und für die substanzielle Interpretation der Ergebnisse vorherzusagen,
  - Entscheidungen *post hoc* auf Basis der Ergebnisse wissenschaftstheoretisch und -praktisch problematisch sein können (*overfitting*, *harking* bzw. *hindsight bias*, etc.).
- In der zugrunde liegenden Studie habe ich versucht, diese Entscheidungen *a priori* zu treffen. Die Entscheidungen basieren aber zugegebenermaßen mehr auf vagen Vermutungen und für mich plausiblen und pragmatischen Überlegungen als auf einer konsistenten Theorie.
  - Entfernen von Stoppwörtern: Stoppwörter sind Wörter, die in einer Sprache häufig vorkommen und nicht wesentlich zur Bedeutung eines Texts beitragen. Hier habe ich auf Basis der deutschen Liste im Paket `{stopwords}` und der Worthäufigkeiten im Korpus eine Liste erstellt. Durch das *Pruning* der Dokument-Feature-Matrix (siehe unten) ist die Auswahl der Stoppwörter aber weniger entscheidend, da Wörter, die in sehr vielen Texten des Korpus vorkommen, ohnehin entfernt werden.
  - Zusätzliche Berücksichtigung von Bi- und Tri-Grammen: Ich habe die Kombinationen von zwei oder drei Wörtern, die häufig im Korpus vorkamen, daraufhin gesichtet, ob sie für das Thema Impfen und gesundheitsrelevante Diskussionen zusätzliche Informationen enthalten, die jedes einzelne Wort alleine nicht enthält. Diese Kombinationen wurden als zusätzliche Features aufgenommen.
  - Der Argumentation und den empirischen Ergebnissen von ? (deren Aufsatz übrigens einen großartigen Titel hat, großer NLP Nerd Humor) folgend habe ich auf Stemming oder Lemmatisierung verzichtet. In der Tat zeigt sich, dass Wörter mit dem gleichen Wortstamm, wie von ? beschrieben, häufig im selben Topic landen.
  - Üblichen Standards (z.B. ?) folgend habe ich alle Wörter in Kleinschreibung umgewandelt, Satzzeichen entfernt und URL entfernt. Zahlen habe ich beibehalten, da sie (wie die Ergebnisse auch zeigen) typische Merkmale bestimmter Perspektiven auf das Thema Impfen

sind.

- Da wir auch an der Veränderung der Topic-Häufigkeiten über die Zeit interessiert sind, wird die Variable mit dem Erscheinungstags des Posts in eine numerische Variable umgewandelt. Sie ist so skaliert, dass der aktuellste Post den Wert 0 hat. Diese Variable können wir dann als Prädiktor beim Schätzen des *Structural Topic Model* berücksichtigen.
- Unter *Pruning* versteht man das Entfernen von Features, die entweder in sehr weniger oder in sehr vielen Dokumenten vorkommen. Dadurch können die Größe des Datensatzes und in der Folge die zum Schätzen der Modelle nötigen Ressourcen wesentlich reduziert werden. Inhaltlich sollte das Entfernen dieser Features wenig ändern: Features, die in sehr vielen Dokumenten vorkommen, tragen nicht zur Differenzierung zwischen den Dokumenten bei. Features, die nur in sehr wenigen Dokumenten vorkommen, tragen nicht zur Definition von Topics bei, da diese durch das regelmäßige *gemeinsame* Vorkommen in Dokumenten identifiziert werden. Siehe ausführlich ?.
- Die Vorbereitung des Korpus und der Dokument-Feature-Matrix erfolgte mit Funktionen aus `{quanteda}`.
  - Mit der Funktion `corpus()` wird der Datensatz in einen Text-Korpus umgewandelt. In diesem Zuge wird auch die numerische Datums-Variable erstellt. Die Variable mit dem Text des Posts duplizieren wir, damit sie zusätzlich als Meta-Datum für jeden Text gespeichert wird. Das wird später hilfreich sein, wenn wir die Ergebnisse einer Modellschätzung explorieren.
  - `custom_stopwords` und `relevant_ngrams` zeigen die Stoppwörter und Wortkombinationen, die ausgeschlossen bzw. einbezogen werden. Letztere werden mit der Funktion `dictionary()` aus `{quanteda}` erstellt.
  - Mit der Funktion `dfm()` wird der Korpus in eine Dokument-Feature-Matrix umgewandelt. Dabei werden die Standard-Schritte der Textaufbereitung durchgeführt. Sie besteht aus 12,369 Posts in den Zeilen und 41,385 Features in den Spalten. In jeder Zelle ist angegeben, wie häufig ein Feature in einem Dokument vorkommt.
  - Mit der Funktion `dfm_trim()` wird das Pruning durchgeführt. Dabei werden alle Features, die in weniger als 0.5% oder mehr als 99% der Posts vorkommen, entfernt. Nach dem Pruning enthält die Matrix nur noch 1,150 Features.
  - Zuletzt muss die Matrix in das von `stm()` benötigte Format konvertiert werden. Dabei werden zwei Posts gelöscht, die nach der Bereinigung kein einziges Feature mehr enthalten. Wichtig für den Bericht der Fallzahl in einer Publikation!
- Am Ende seht ihr eine einfache Beschreibung der häufigsten Features im Korpus als Tabelle und Wordcloud.



```

dfm(stem = FALSE, tolower = TRUE, remove_punct = TRUE,
    remove = custom_stopwords,
    remove_url = TRUE, verbose = TRUE,
    thesaurus = relevant_ngrams)

## Creating a dfm from a corpus input...

## ... lowercasing

## ... found 12,369 documents, 41,651 features

## ... applying a dictionary consisting of 20 keys
## ... removed 286 features

impf_dfm

## Document-feature matrix of: 12,369 documents, 41,385 features (99.9% sparse) and 6 d
##          features
## docs    TROTZ_IMPFUNG GRIPPE_IMPFEN MMR_IMPFUNG HEPATITIS_B GUT_VERTRAGEN
## text1           0           0           0           0           0
## text2           0           0           0           0           0
## text3           0           0           0           0           0
## text4           1           0           0           0           0
## text5           0           0           0           0           0
## text6           0           0           0           0           0
##          features
## docs    6FACH_IMPFUNG 6_FACH 6_FACH_IMPFUNG MENINGOKOKKEN_B GUTE_BESSERUNG
## text1           0      0           0           0           0
## text2           0      0           0           0           0
## text3           0      0           0           0           0
## text4           1      0           0           0           0
## text5           0      0           0           0           0
## text6           0      0           0           0           0
## [ reached max_ndoc ... 12,363 more documents, reached max_nfeat ... 41,375 more feat

# Pruning
impf_dfm = impf_dfm %>%
  dfm_trim(max_docfreq = 0.99, min_docfreq = 0.005, docfreq_type = "prop")
impf_dfm

## Document-feature matrix of: 12,369 documents, 1,150 features (98.2% sparse) and 6 d
##          features
## docs    TROTZ_IMPFUNG GRIPPE_IMPFEN MMR_IMPFUNG GUT_VERTRAGEN 6_FACH
## text1           0           0           0           0           0
## text2           0           0           0           0           0
## text3           0           0           0           0           0
## text4           1           0           0           0           0
## text5           0           0           0           0           0

```

```
## text6          0          0          0          0          0
##          features
## docs  MENINGOKOKKEN_B GUTE_BESSERUNG ERHÖHTE_TEMPERATUR KEIN_FIEBER
## text1          0          0          0          0          0
## text2          0          0          0          0          0
## text3          0          0          0          0          0
## text4          0          0          0          0          0
## text5          0          0          0          0          0
## text6          0          0          0          0          0
##          features
## docs  KEIN_PROBLEM
## text1          0
## text2          0
## text3          0
## text4          0
## text5          0
## text6          0
## [ reached max_ndoc ... 12,363 more documents, reached max_nfeat ... 1,140 more features ]

# Überblick: Die häufigsten Features im Korpus
impf_dfm %>%
  colSums() %>%
  enframe() %>%
  arrange(desc(value)) %>%
  slice(1:20) %>%
  kable()
```

name	value
impfung	4519
impfen	4356
kind	3690
lassen	3257
immer	2705
mehr	2594
kinder	2503
gibt	2184
geimpft	2176
impfungen	2174
gut	2115
hallo	1898
einfach	1767
lg	1720
2	1704
bekommen	1585
ganz	1584
erst	1558
geht	1492
arzt	1411

```
# als (beliebte, wenn auch nur mittel informative) Wordcloud
impf_dfm %>%
  textplot_wordcloud()
```

```
## Warning in graphics::strwidth(word[i], cex = size[i]): Konvertierungsfehler für
## ' ' in 'mbcsToSbcs': Punkt ersetzt <f0>

## Warning in graphics::strwidth(word[i], cex = size[i]): Konvertierungsfehler für
## ' ' in 'mbcsToSbcs': Punkt ersetzt <9f>

## Warning in graphics::strwidth(word[i], cex = size[i]): Konvertierungsfehler für
## ' ' in 'mbcsToSbcs': Punkt ersetzt <98>

## Warning in graphics::strwidth(word[i], cex = size[i]): Konvertierungsfehler für
## ' ' in 'mbcsToSbcs': Punkt ersetzt <8a>

## Warning in text.default(x1, y1, word[i], cex = (1 + adjust) * size[i], offset =
## 0, : Konvertierungsfehler für ' ' in 'mbcsToSbcs': Punkt ersetzt <f0>

## Warning in text.default(x1, y1, word[i], cex = (1 + adjust) * size[i], offset =
## 0, : Konvertierungsfehler für ' ' in 'mbcsToSbcs': Punkt ersetzt <9f>

## Warning in text.default(x1, y1, word[i], cex = (1 + adjust) * size[i], offset =
## 0, : Konvertierungsfehler für ' ' in 'mbcsToSbcs': Punkt ersetzt <98>

## Warning in text.default(x1, y1, word[i], cex = (1 + adjust) * size[i], offset =
## 0, : Konvertierungsfehler für ' ' in 'mbcsToSbcs': Punkt ersetzt <8a>
```



```

## Warning in text.default(x1, y1, word[i], cex = (1 + adjust) * size[i], offset =
## 0, : Fontmetrik ist für das Unicode-Zeichen U+1f60a unbekannt

## Warning in graphics::strwidth(word[i], cex = size[i]): Konvertierungsfehler für
## ' ' in 'mbcsToSbcs': Punkt ersetzt <f0>

## Warning in graphics::strwidth(word[i], cex = size[i]): Konvertierungsfehler für
## ' ' in 'mbcsToSbcs': Punkt ersetzt <9f>

## Warning in graphics::strwidth(word[i], cex = size[i]): Konvertierungsfehler für
## ' ' in 'mbcsToSbcs': Punkt ersetzt <98>

## Warning in graphics::strwidth(word[i], cex = size[i]): Konvertierungsfehler für
## ' ' in 'mbcsToSbcs': Punkt ersetzt <82>

## Warning in text.default(x1, y1, word[i], cex = (1 + adjust) * size[i], offset =
## 0, : Konvertierungsfehler für ' ' in 'mbcsToSbcs': Punkt ersetzt <f0>

## Warning in text.default(x1, y1, word[i], cex = (1 + adjust) * size[i], offset =
## 0, : Konvertierungsfehler für ' ' in 'mbcsToSbcs': Punkt ersetzt <9f>

## Warning in text.default(x1, y1, word[i], cex = (1 + adjust) * size[i], offset =
## 0, : Konvertierungsfehler für ' ' in 'mbcsToSbcs': Punkt ersetzt <98>

## Warning in text.default(x1, y1, word[i], cex = (1 + adjust) * size[i], offset =
## 0, : Konvertierungsfehler für ' ' in 'mbcsToSbcs': Punkt ersetzt <82>

## Warning in text.default(x1, y1, word[i], cex = (1 + adjust) * size[i], offset =
## 0, : Fontmetrik ist für das Unicode-Zeichen U+1f602 unbekannt

## Warning in graphics::strwidth(word[i], cex = size[i]): Konvertierungsfehler für
## ' ' in 'mbcsToSbcs': Punkt ersetzt <f0>

## Warning in graphics::strwidth(word[i], cex = size[i]): Konvertierungsfehler für
## ' ' in 'mbcsToSbcs': Punkt ersetzt <9f>

## Warning in graphics::strwidth(word[i], cex = size[i]): Konvertierungsfehler für
## ' ' in 'mbcsToSbcs': Punkt ersetzt <98>

## Warning in graphics::strwidth(word[i], cex = size[i]): Konvertierungsfehler für
## ' ' in 'mbcsToSbcs': Punkt ersetzt <89>

## Warning in text.default(x1, y1, word[i], cex = (1 + adjust) * size[i], offset =
## 0, : Konvertierungsfehler für ' ' in 'mbcsToSbcs': Punkt ersetzt <f0>

## Warning in text.default(x1, y1, word[i], cex = (1 + adjust) * size[i], offset =
## 0, : Konvertierungsfehler für ' ' in 'mbcsToSbcs': Punkt ersetzt <9f>

## Warning in text.default(x1, y1, word[i], cex = (1 + adjust) * size[i], offset =
## 0, : Konvertierungsfehler für ' ' in 'mbcsToSbcs': Punkt ersetzt <98>

## Warning in text.default(x1, y1, word[i], cex = (1 + adjust) * size[i], offset =
## 0, : Konvertierungsfehler für ' ' in 'mbcsToSbcs': Punkt ersetzt <89>

```

```

## Warning in text.default(x1, y1, word[i], cex = (1 + adjust) * size[i], offset =
## 0, : Fontmetrik ist für das Unicode-Zeichen U+1f609 unbekannt

## Warning in graphics::strwidth(word[i], cex = size[i]): Konvertierungsfehler für
## ' ' in 'mbcsToSbcs': Punkt ersetzt <f0>

## Warning in graphics::strwidth(word[i], cex = size[i]): Konvertierungsfehler für
## ' ' in 'mbcsToSbcs': Punkt ersetzt <9f>

## Warning in graphics::strwidth(word[i], cex = size[i]): Konvertierungsfehler für
## ' ' in 'mbcsToSbcs': Punkt ersetzt <99>

## Warning in graphics::strwidth(word[i], cex = size[i]): Konvertierungsfehler für
## ' ' in 'mbcsToSbcs': Punkt ersetzt <88>

## Warning in text.default(x1, y1, word[i], cex = (1 + adjust) * size[i], offset =
## 0, : Konvertierungsfehler für ' ' in 'mbcsToSbcs': Punkt ersetzt <f0>

## Warning in text.default(x1, y1, word[i], cex = (1 + adjust) * size[i], offset =
## 0, : Konvertierungsfehler für ' ' in 'mbcsToSbcs': Punkt ersetzt <9f>

## Warning in text.default(x1, y1, word[i], cex = (1 + adjust) * size[i], offset =
## 0, : Konvertierungsfehler für ' ' in 'mbcsToSbcs': Punkt ersetzt <99>

## Warning in text.default(x1, y1, word[i], cex = (1 + adjust) * size[i], offset =
## 0, : Konvertierungsfehler für ' ' in 'mbcsToSbcs': Punkt ersetzt <88>

## Warning in text.default(x1, y1, word[i], cex = (1 + adjust) * size[i], offset =
## 0, : Fontmetrik ist für das Unicode-Zeichen U+1f648 unbekannt

## Warning in graphics::strwidth(word[i], cex = size[i]): Konvertierungsfehler für
## ' ' in 'mbcsToSbcs': Punkt ersetzt <f0>

## Warning in graphics::strwidth(word[i], cex = size[i]): Konvertierungsfehler für
## ' ' in 'mbcsToSbcs': Punkt ersetzt <9f>

## Warning in graphics::strwidth(word[i], cex = size[i]): Konvertierungsfehler für
## ' ' in 'mbcsToSbcs': Punkt ersetzt <98>

## Warning in graphics::strwidth(word[i], cex = size[i]): Konvertierungsfehler für
## ' ' in 'mbcsToSbcs': Punkt ersetzt <85>

## Warning in text.default(x1, y1, word[i], cex = (1 + adjust) * size[i], offset =
## 0, : Konvertierungsfehler für ' ' in 'mbcsToSbcs': Punkt ersetzt <f0>

## Warning in text.default(x1, y1, word[i], cex = (1 + adjust) * size[i], offset =
## 0, : Konvertierungsfehler für ' ' in 'mbcsToSbcs': Punkt ersetzt <9f>

## Warning in text.default(x1, y1, word[i], cex = (1 + adjust) * size[i], offset =
## 0, : Konvertierungsfehler für ' ' in 'mbcsToSbcs': Punkt ersetzt <98>

## Warning in text.default(x1, y1, word[i], cex = (1 + adjust) * size[i], offset =
## 0, : Konvertierungsfehler für ' ' in 'mbcsToSbcs': Punkt ersetzt <85>

```





## Chapter 3

# Modellspezifikation, Modellvergleich und Modellauswahl

### 3.1 Modellspezifikation

- Die Funktion `stm()` aus dem gleichnamigen Paket bietet zahlreiche Möglichkeiten, Details der Modellspezifikation und -schätzung anzupassen. Wir beschränken uns im Folgenden auf drei wesentliche Einstellungen:
  - `K`: Die Zahl der Topics.
  - `prevalence`: Eine Formel zur Vorhersage der Topic-Prävalenzen
  - `init.type`: Wie soll der Startpunkt für die Modellschätzung gewählt werden?
- Zu weiteren Details siehe für einen Überblick `?stm` und `?` für eine ausführliche Erläuterung.
- Eine Modell-Spezifikation könnte z.B. so aussehen:

```
modelfit = stm(documents = impf_stm$documents,  
               vocab = impf_stm$vocab,  
               data = impf_stm$meta,  
               K = 10,  
               prevalence = ~s(date_num),  
               init.type = "Spectral")
```

- Mit den ersten drei Inputs übergeben wir die Daten aus der im letzten Abschnitt erstellten Dokument-Feature-Matrix.

- Mit `K` geben wir an, wie viele Themen es geben soll. Mit dieser Syntax würde ein Modell mit  $k = 10$  Themen geschätzt. Wie wir bei der Wahl eines geeigneten  $k$  vorgehen können, ist Thema des folgenden Unterabschnitts.
- Mit der Formel zu **prevalence** geben wir an, welche Dokument-Variablen mit dem Auftreten der Topics zusammenhängen.
  - In der Formel wird die abhängige Variable vor der Tilde ( $\sim$ ) freigelassen. Es wird immer der Zusammenhang mit dem Auftreten von allen  $k$  Topics geschätzt. In diesem Beispiel schätzen wir, wie sich das Auftreten der Topics über den Untersuchungszeitraum hinweg verändert. Details zum Schätzen von Zusammenhängen mit Kovariaten folgen später in diesem Workshop.
  - Mit `init.type` wird angegeben, wie `stm()` die Ausgangswerte für die Modellschätzung bestimmen soll. Die Default-Einstellung ist “Spectral”. Ich empfehle diese Einstellung aus folgenden Gründen:
    - \* Sie ist deterministisch, d.h., sie führt gegeben derselben Daten und desselben Modells immer zu derselben Lösung. So wird die Reproduzierbarkeit sichergestellt.
    - \* Sie ist effizient, d.h., dass von diesem Startpunkt aus relativ schnell die finale Lösung gefunden wird.
    - \* Wenn eine andere Einstellung für die Ausgangswerte gewählt wird, müssen mehrere Schätzungen für eine Spezifikation durchgeführt werden. Nur so kann geprüft werden, ob die Ausgangswerte das Ergebnis beeinflussen.
- Allgemein muss beachtet werden, dass die Schätzung eines Structural Topic Model mit `stm()` trotz der Effizienz der Implementierung sehr rechenintensiv ist. Die Schätzung des oben beschriebenen Modells dauert auf meinem recht leistungsfähigen Notebook bereits ca. eine Minute. Es empfiehlt sich daher, die Modelle immer in neue Objekte zu speichern und diese ggf. direkt auf der Festplatte zu sichern. Um die Berechnungszeiten in diesem Workshop kurz zu halten, stelle ich die Ergebnisse der Modellschätzungen über das LMS zur Verfügung. Wenn ihr diese herunterladet und in den Ordner “data” kopiert, muss das Modell nicht neu geschätzt werden.

## 3.2 Modellvergleich

### 3.2.1 Allgemeines Vorgehen

- Eine zentrale Frage ist die Wahl eines geeigneten  $k$ , also der Zahl von Topics, die in den Dokumenten identifiziert werden sollen. Wichtig ist zuerst die Feststellung, dass es in der angewandten Analyse kein *per se* richtiges oder falsches  $k$  gibt.
- Wie viele Topics nützlich sind, hängt von Umfang von Zusammensetzung des Materials und vom substantiellen Forschungsinteresse ab.

- Um ein geeignetes  $k$  zu finden, gehen wir in der Regel modellvergleichend vor. Wir schätzen Modelle mit unterschiedlich vielen Topics und prüfen dann, welche Modelle besser zu den Daten und zum Forschungsinteresse passen.
- Hinweise für einen allgemeinen Ausgangspunkt, in welchem Bereich nützliche  $k$  zu finden sein könnten, liefert die Paket-Hilfe:

The most important user input in parametric topic models is the number of topics. There is no right answer to the appropriate number of topics. More topics will give more fine-grained representations of the data at the potential cost of being less precisely estimated. [...] For short corpora focused on very specific subject matter (such as survey experiments) 3-10 topics is a useful starting range. For small corpora (a few hundred to a few thousand) 5-50 topics is a good place to start. Beyond these rough guidelines it is application specific. Previous applications in political science with medium sized corpora (10k to 100k documents) have found 60-100 topics to work well. For larger corpora 100 topics is a useful default size. Of course, your mileage may vary. — ?stm

- Hier werden zwei wichtige Kriterien, die unser Nachdenken über die Spannweite von zu Berücksichtigten  $k$  leiten können, deutlich:
  - Quantität des Materials: Je mehr Dokumente, desto mehr Topics.
  - Varianz im Inhalt: Je mehr inhaltliche Varianz, desto mehr Topics (an einem Beispiel: für 10k Nachrichtenbeiträge aus dem Wirtschaftsteil brauchen wir weniger Topics als für 10k Nachrichtenbeiträge, die aus allen Ressorts kommen).
- Im vorliegenden Fall haben wir einen kleinen bis mittleren Korpus (ca. 13k Dokumente, die größtenteils recht kurz sind). Wir können von einer mittleren inhaltlichen Varianz ausgehen. Einerseits haben wir Posts bewusst danach ausgewählt, dass sie sich mit dem Thema Impfen beschäftigen, was die Varianz einschränkt. Andererseits wissen wir, dass in Online-Foren die verschiedensten Perspektiven auf dieses Thema vorkommen können, was für Varianz sorgt.
- Wir gehen im Folgenden in mehreren Schritten vor:
  - 1) Um eine allgemeine Orientierung zu erhalten, in welcher Range Modelle zu finden sind, die gut zu den Daten passen, schätzen wir 10 Modelle von  $k = 10$  bis  $k = 100$  mit einem Abstand von jeweils 10 Topics. Diese Modelle vergleichen wir anhand von einigen statistischen Maßen, um die Zahl der Kandidatenmodelle einzuschränken.
  - 2) Wir interpretieren die besten Modelle substantiell und entscheiden, welche Topic-Anzahl für das Forschungsinteresse hilfreicher scheint.
  - 3) Wir schätzen weitere Modelle basierend auf den Ergebnissen aus 2) mit kleineren Abständen zwischen den  $k$ . Wir prüfen, wie sich die Topics verändern. Zudem achten wir darauf, ob bei Modellen mit größeren  $k$  interessante Topics hinzukommen oder ob sich mehr Ambivalenzen zeigen.

- 4) Wir entscheiden uns für ein Modell.
- Da das Schätzen der Modelle recht lange dauert, parallelisieren wir die Berechnung. Dazu nutze ich das Paket `furrr`. Es sei an dieser Stelle darauf hingewiesen, dass das Paket vor allem unter Windows mit RStudio für Probleme sorgen kann. Es ist daher empfehlenswert, das Skript zum Schätzen und Speichern der Modelle in der Konsole oder im Terminal auszuführen. Noch schneller geht es für diesen Workshop, die bereits geschätzten Modelle aus dem LMS zu laden.

### 3.2.2 Quantitativer Vergleich der ersten Modelle

```
# 1) Modelle mit K = 10, ..., K = 100
# Modelle schätzen bzw. laden
if (file.exists("R/data/models10_100.rds")) {
  # Schneller: Modelle aus LMS laden
  many_models = read_rds("R/data/models10_100.rds")
} else {
  # Vorsicht: Schätzen dauert auf meinem MacBook Pro 2020 i9 32 GB RAM 10 Minuten
  Ks = seq(10, 100, by = 10)
  tic()
  plan(multiprocess(workers = 10))
  many_models = tibble(K = Ks) %>%
    mutate(topic_model = future_map(K, ~stm(documents = impf_stm$documents,
                                             vocab = impf_stm$vocab,
                                             data = impf_stm$meta,
                                             init.type = "Spectral",
                                             K = ., verbose = FALSE),
                                             .progress = TRUE))

  plan(sequential)
  toc()
  saveRDS(many_models, "R/data/models10_100.rds")
}

# Quantitative Indikatoren der Modellqualität berechnen
# Inspired by https://juliasilge.com/blog/evaluating-stm/
if (file.exists("R/data/eval10_100.rds")) {
  # Schneller: Modellevaluation aus LMS laden
  model_eval = read_rds("R/data/eval10_100.rds")
} else {
  # Evaluation dauert auf meinem MacBook Pro 2020 i9 32 GB RAM 24 Sekunden
  tic()
  heldout = make.heldout(documents = impf_stm$documents, vocab = impf_stm$vocab)
  plan(multiprocess(workers = 10))
  model_eval = many_models %>%
    mutate(exclusivity = future_map(topic_model, exclusivity),
```

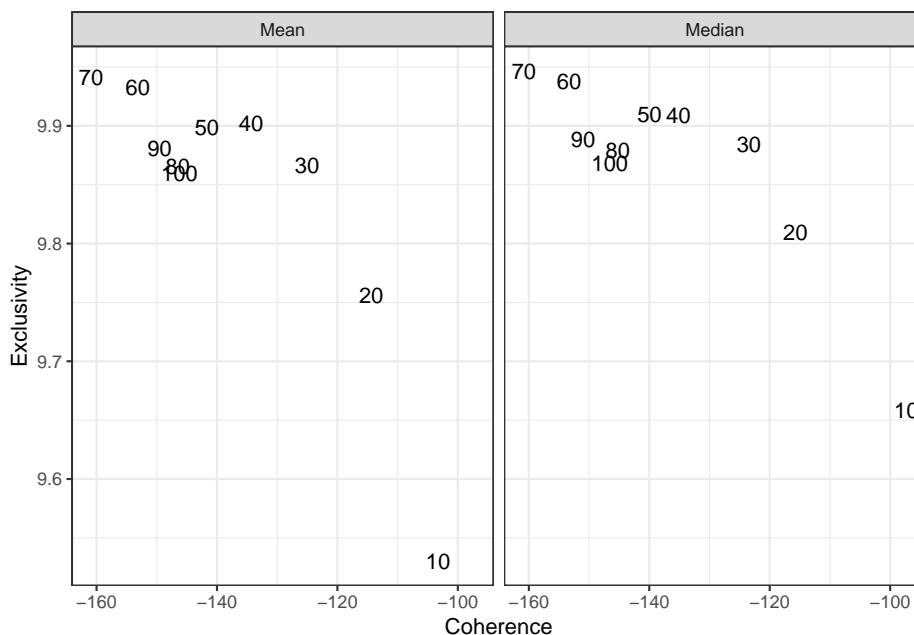


```

    semantic_coherence = future_map(topic_model, semanticCoherence, impf_stm$documents),
    eval_heldout = future_map(topic_model, eval.heldout, heldout$missing),
    residual = future_map(topic_model, checkResiduals, impf_stm$documents),
    residuals = future_map_dbl(residual, "dispersion"),
    held_out_likelihood = future_map_dbl(eval_heldout, "expected.heldout")) %>%
  select(-topic_model)
plan(sequential)
toc()
saveRDS(model_eval, "R/data/eval10_100.rds")
}

# Exklusivität und Semantische Kohärenz (Mean, Median)
model_eval %>%
  select(K, Exclusivity = exclusivity, Coherence = semantic_coherence) %>%
  mutate_at(-1, .funs = list(Mean = ~map_dbl(.x, mean), Median = ~map_dbl(.x, median))) %>%
  select_if(negate(is.list)) %>%
  gather(metric, value, -K) %>%
  separate(metric, c("measure", "metric"), "_") %>%
  spread(measure, value) %>%
  ggplot(aes(Coherence, Exclusivity, label = K)) + geom_text() + facet_wrap("metric")

```

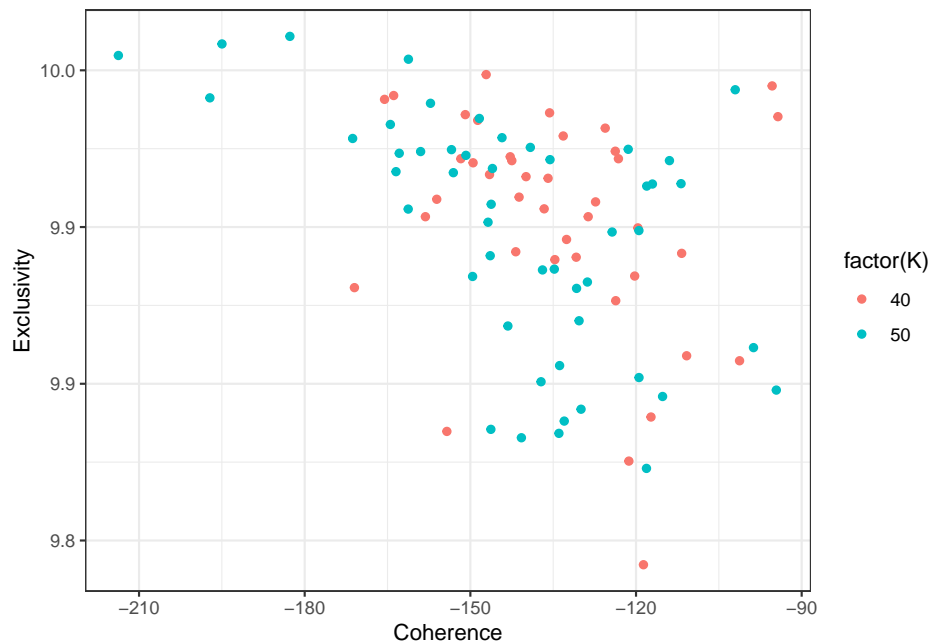


```

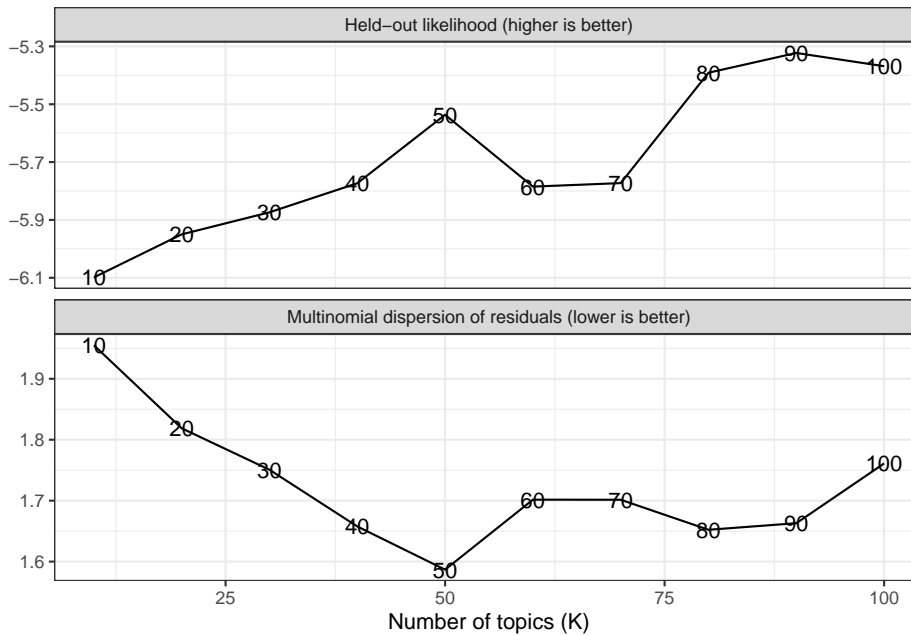
# Exklusivität und Semantische Kohärenz von k = 40, 50
model_eval %>%
  filter(K %in% c(40, 50)) %>%
  select(K, Exclusivity = exclusivity, Coherence = semantic_coherence) %>%

```

```
unnest(c(Exclusivity, Coherence)) %>%
  ggplot(aes(Coherence, Exclusivity, color = factor(K))) + geom_point()
```



```
# Held-out-likelihood und multinomiale Residuen
model_eval %>%
  select(K, `Held-out likelihood (higher is better)` = held_out_likelihood, `Multinomial Residuals` = multinomial_residuals) %>%
  gather(measure, value, -K) %>%
  ggplot(aes(K, value, label = K)) + geom_line() + geom_text() + facet_wrap("measure",
```



```
# Speichern der Modelle mit 30 und 40 Topics für qualitative Analyse
# Benennung und Datenstruktur für stm insights
out = impf_stm
m30 = many_models$topic_model[[3]]
m40 = many_models$topic_model[[4]]
m60 = many_models$topic_model[[6]]
# save(out, m30, m40, m60, file = "R/data/models30_40_60.rdata")
```

- Wir betrachten zuerst die *semantische Kohärenz* und die *Exklusivität* der Topics in den Modellen. Beide Metriken sind Eigenschaften der einzelnen Themen. In der ersten Grafik sind daher der Mittelwerte und die Mediane aller Themen in einem Modell dargestellt. Die absoluten Werte beider Metriken haben keine substantielle Bedeutung. Von Interesse ist der Vergleich der Modelle. Je höher die Exklusivität, desto geringer ist die Wahrscheinlichkeit, dass die typischsten Terme eines Topics in anderen Topics vorkommen (?). Je höher die semantische Kohärenz, desto wahrscheinlicher kommen die typischsten Wörter eines Topics gemeinsam in einem Dokument mit diesem Topic vor (?). Zwischen den beiden Metriken besteht in der Regel ein negativer Zusammenhang. Daher ist es notwendig, eine Balance zwischen beiden zu finden.
  - Im vorliegenden Beispiel können wir zuerst die Modelle mit  $k \geq 80$  ausschließen. Obwohl sie mehr Topics benötigen, sind ihre Topics im Mittel weniger exklusiv und weniger kohärent als die Topics der Modelle  $k = 60$  oder  $k = 70$ . Eine Erklärung dafür kann sein, irgendwann zwischen dem 70. und dem 80. Topic keine substantiell neuen

Aspekte mehr im Korpus zu finden sind. Die neuen Topics sind dann redundant zu schon bestehenden Topics.

- Ebenso können wir für die meisten Zwecke die Modelle mit  $k \leq 20$  vernachlässigen, da die mittlere Exklusivität ihrer Topics wesentlich geringer ist als die der übrigen Modelle. Das Modell mit  $k = 20$  könnte vielleicht infrage kommen, wenn wir Sparsamkeit sehr hoch gewichten, also den Korpus durch möglichst wenige Topics beschreiben wollen.
- Schließlich ist das Modell mit  $k = 50$  nicht besonders attraktiv. Die Topics sind im Mittel nur weniger kohärent, aber nicht exklusiver, als die Topics des Modells mit  $k = 40$ . Der Detailvergleich der beiden Modelle in der nächsten Abbildung, in der jedes Topic mit einem Punkt dargestellt ist, zeigt, dass dies vor allem auf vier Topics zurück geht, die weniger kohärent sind, ohne eine besonders gute Exklusivität aufzuweisen.
- Gegeben der Metriken *semantische Kohärenz* und *Exklusivität* spricht einiges dafür, entweder Modelle mit ca. 60-70 Topics oder Modelle mit 30-40 Topics weiter zu verfolgen.
- Die *held-out likelihood* und die *Dispersion der multinomialen Residuen* sind zwei Metriken, die die Abweichung des Modells von den Daten quantifizieren. Auch sie sind wieder im relativen Modellvergleich zu interpretieren. Die *held-out likelihood* gibt Auskunft darüber, wie gut das Vorkommen von Wörtern in einem Dokument, das nicht zum Schätzen des Modells genutzt wurde, vorhergesagt werden kann. Mit der *Dispersion der multinomialen Residuen* wird die Abweichung der durch das Modell vorhergesagten von den beobachteten Wörtern in den Dokumenten auf Basis des gesamten Datensatzes quantifiziert. Ein Wert von 1 wäre ideal, wird in angewandten Beispielen mit echten Texten aber kaum erreicht. Beide Metriken sind eng verwandt und legen in der Regel ähnliche Entscheidungen nahe.
  - Das Modell mit  $k = 50$  ist nach beiden Metriken gut geeignet.
  - Nach der *held-out likelihood* sind die Modelle mit  $k \geq 80$  noch etwas besser - diese haben sich aber in der semantischen Kohärenz und Exklusivität nicht sonderlich gut bewährt.
  - Die nach semantischer Kohärenz und Exklusivität besten Modelle liegen nach diesen beiden Metriken etwa gleich auf.
- Die Befunde sind auch für die didaktischen Zwecke dieses Workshops gut geeignet. Es wird klar, dass die Modellwahl sich nicht einfach automatisieren lässt. Die quantitativen Kriterien helfen und lediglich, den Raum möglicher Modelle einzuschränken.
  - Für mich kommen auf Basis der berichteten Metriken Modelle mit zwischen 30 und 70 Topics infrage.
  - Wenn man an einer besonders sparsamen Lösung interessiert ist, könnte man auch  $k = 20$  in Erwägung ziehen.
  - Wenn man besonders detailliertere Lösungen sucht, sind auch die Modelle mit  $k \geq 80$  nicht ausgeschlossen. Hier müsste man dann die

Differenzierungen zwischen den Topics in der Tiefe untersuchen und klarstellen.

- Im nächsten Schritt vergleichen wir die Modelle mit 30 und mit 60 Topics, um zu verstehen, welche Konsequenzen es für die inhaltliche Interpretation hat, ein eher sparsames oder ein eher detailliertes Modell zu wählen. Nach dieser Entscheidung können wir dann weiter überlegen, welches der eher sparsamen bzw. detaillierten Modelle für unser Forschungsinteresse besser geeignet ist.

### 3.2.3 Qualitativer Vergleich der Modelle mit 30 und 60 Topics

- Um die Ergebnisse eines Topic Model interpretieren zu können, müssen wir zuerst verstehen, wie das Ergebnis eines Topic Models aussieht. Es enthält im wesentlichen zwei zweidimensionale Vektoren von Koeffizienten.
  - beta: Für jedes Topic die Wahrscheinlichkeit, dass ein Dokument, in dem das Feature vorkommt, das Topic hat. Jedes Topic erhält für jedes Feature einen beta-Koeffizienten zwischen 0 und 1, die zusammen 1 ergeben.
  - theta bzw. gamma (uneinheitlich benannt): Für jedes Dokument die Wahrscheinlichkeit, dass das Dokument das Topic enthält. Jedes Dokument erhält für jedes Topic einen theta- bzw. gamma-Koeffizienten zwischen 0 und 1, die zusammen 1 ergeben.
- Um die Bedeutung der Topics zu interpretieren, betrachten wir daher zwei Modell-Outputs.
  - Die Features, die am typischsten für Dokumente mit einem Topic sind, also die höchsten beta-Koeffizienten (ggf. nach Korrekturen für die Verteilung der Features im gesamten Korpus).
  - Die Dokumente, die mit der größten Wahrscheinlichkeit ein Thema enthalten (manchmal auch interpretiert als zum größten Anteil aus einem Thema bestehen), also die Dokumente mit den höchsten gamma- bzw. theta-Koeffizienten.
- Zur qualitativen Interpretation der Modelle anhand dieser Outputs kann ich das Paket `{stminsights}` empfehlen — vor allem denjenigen, die lieber mit einer intuitiven grafischen Benutzeroberfläche als direkt mit *R* arbeiten. Aber auch ich selbst schätze das Tool für einen schnellen Überblick über mehrere Modelle.
- Leider ist es zurzeit nicht einfach, `{stminsights}` direkt zum Laufen zu bringen, da einige Pakete, auf denen es aufbaut, Bugs bzw. Kompatibilitätsprobleme haben. Eine Anleitung, wie das Paket zurzeit installiert werden kann, gibt es hier:

Important note: The shiny app for the CRAN release of `stminsights` does currently not work properly due to bugs introduced by recent changes in the Shiny package [...]. Please use the Github version of `stminsights` for now. This will require the

development version of Shiny which can be installed by running  
`devtools::install_github('rstudio/shiny')`.

You can download and install the latest development version of  
`stminsights` by running `devtools::install_github('cschwem2er/stminsights')`  
 — <https://github.com/cschwem2er/stminsights>

- Zur Vorbereitung der Analyse mit `{stminsights}` müssen die Objekte der geschätzten Modelle (hier: `m30`, `m40` [brauchen wir etwas später] und `m60`) und die Daten, die zur Modellschätzung verwendet wurden (hier: `out = impf_stm`), als `.rdata` Datei gespeichert werden. Dabei muss das Daten-Objekt *unbedingt* den Namen `out` haben. Damit der Text der Dokumente in `{stminsights}` angezeigt werden kann, muss dieser zusätzlich als Variable in den Meta-Daten enthalten sein (siehe Abschnitt zur Datenaufbereitung).
- Durch das ausführen der Funktion `run_stminsights()` wird eine grafische Benutzeroberfläche im Browser gestartet, mit dem die qualitative Analyse durchgeführt wird. Eine Beschreibung der Oberfläche findet ihr als Video im LMS [kommt nach Fertigstellen des Textmaterials].

#### `run_stminsights()`

- Wer die Interpretation der Topics lieber in R durchführt, kann z.B. den folgenden Code verwenden und anpassen.
  - Wir sammeln zuerst die 20 typischsten Texte für jedes Topic in einem Datensatz. Dabei bereiten wir die Texte auch für eine Ausgabe in der Konsole vor. Der Parameter “gamma” oder “theta” (uneinheitlich benannt, aber derselbe Parameter) kann mit `tidytext::tidy()` extrahiert werden. Dann werden die Texte aus den Ursprungsdaten abgespielt.
  - Mit `stm::labelTopics()` können wir die typischsten Features für ein Topic anzeigen. `n` bestimmt die Zahl der Features, `topics` das Topic, das wir gerade beschreiben möchten. Es werden die typischsten Features nach vier verschiedenen Kriterien ausgegeben (Diese werden auch in `{stminsights}` angezeigt):
    - \* *Highest Prob* zeigt die Features, die mit der größten Wahrscheinlichkeit in Dokumenten mit einem Topic vorkommen (= Features mit den höchsten beta-Koeffizienten für das Topic). Dabei wird die Gesamthäufigkeit der Features im Korpus *nicht* berücksichtigt. Wörter, die allgemein sehr häufig vorkommen, sind nach dieser Metrik typisch für verschiedene Topics. Die anderen drei Metriken versuchen, diese Schwäche auf verschiedene Art und Weise zu korrigieren.
    - \* *FREX* steht für *most frequent and exclusive* Features. Hier werden die Features aufgelistet, die möglichst typisch für ein Dokument mit einem Topic, aber möglichst nicht sehr typisch für Dokumente mit anderen Topics sind. Siehe `?calcfrex` für

technische Details.

- \* Die *Lift*-Metrik setzt die Wahrscheinlichkeit, dass ein Feature in einem Dokument mit einem Topic vorkommt, zur Wahrscheinlichkeit, dass ein Feature in einem beliebigen Dokument vorkommt, ins Verhältnis. Siehe `?calclift` für technische Details.
- \* *Score* gewichtet für die Wahrscheinlichkeit, mit der ein Feature in Dokumenten mit einem anderen Topics vorkommt. Siehe `?calcscore` für technische Details.
- Mit den Auszügen aus dem Datensatz typischer Dokumente (gefiltert nach Topic) können wir diese Features zudem im Kontext der gesamten Texte sehen.
- Im abschließenden Datensatz `topic_labels` halten wir für jedes Topic ein aussagekräftiges, möglichst kurzes Label fest, das wir jetzt zur eigenen Übersicht und später auch zur Ergebnisdarstellung verwenden. Zusätzlich empfehle ich, eine kurze Zusammenfassung für jedes Topic auf einem Medium eigener Wahl (präferierte analoger oder digitaler Notizzettel) festzuhalten.

```
# Laden der Modelle
load("R/data/models30_40_60.rdata")

# Beispiel für das Modell mit k = 30
# Erstellen eines Datensatzes mit den typischsten Dokumenten
top_docs = tidy(m30, "gamma") %>%
  arrange(desc(gamma)) %>%
  group_by(topic) %>%
  slice(1:20) %>%
  mutate(rank = 1:20) %>%
  left_join(mutate(select(out$meta, txt), document = 1:n())) %>%
  mutate(out = paste0(round(gamma, 2), ": ", txt))

# Ausgabe der typischen Feature für ein Topic (hier Topic 1)
m30 %>%
  labelTopics(n = 10, topics = 1)

## Topic 1 Top Words:
## Highest Prob: 2, 3, 1, monate, erst, monaten, 4, 6, alt, 5
## FREX: 3, 2, monate, monaten, 6, 4, 1, monat, +, 5
## Lift: österreich, zecken, monat, fach, 3, 2, monate, 6, +, monaten
## Score: österreich, 2, 3, monate, 1, monaten, 4, alt, 6, +

# Ausgabe der typischen Texte für ein Topic (hier Topic 1)
top_docs %>%
  filter(topic == 1) %>%
  filter(rank %in% 1:10) %>%
  .$out %>%
```

```

str_squish() %>%
cat(sep = "\n\n")

## 0.62: da fängt der 12 lebensmonat an, ein Jahr alt ist dein Baby dann erst mit voll
##
## 0.59: Man spricht erst sobald der Monat vollendet ist vom z.B. 12. Monat. Deine Toch
##
## 0.59: In Österreich wird im 3., 5. und 12. Monat geimpft. Ich bin auf einer anderen
##
## 0.53: so mein ich ja :) sie befindet sich ab dem Tag ihrer Geburt ja im erst Lebens
##
## 0.51: Wir mussten 2 Monate auf den Termin beim Kardiologen warten. Hatten den Termin
##
## 0.5: Korrektur: Bei der 2+1 Impfung sind ca 84% nach der 2 . Impfung immun. Bei der
##
## 0.47: hallo rebecca es gibt auch "Mittelwege", zB das 2+1-Schema. Man beginnt erst m
##
## 0.47: Unter 2 Jahren muss 2x geimpft werden. Mein Sohn ist 6 Monate. Wird am Freitag
##
## 0.44: Dann wär er aber 1 Monat vor Geburtstag 12 Monate : Wenn du die Wochen zählst
##
## 0.43: Darf ich über die Fragestellung hinaus fragen, wie genau der Ablauf dann genau
# Datensatz für Topic Labels
# Bis Topic 30 weiterführen
m30_topic_labels = tribble(
  ~topic, ~label,
  1,      "Alter von Babies",
  2,      "label B",
  3,      "label C"
)

```

- Die Labels aus meiner vorläufigen Interpretation finden sich in R/topic\_labels.R (hier nicht dargestellt).
  - Die meisten Topics des Modells mit  $k = 30$  erwiesen sich als gut interpretierbar. An der einen oder anderen Stelle schien es aber so, als seien mehrere Aspekte der Impfdiskussionen in einem Topic vermischt. Dies könnte auf eine etwas größere Topic-Anzahl hindeuten.
  - Die Interpretation der Topics des Modells mit  $k = 60$  fiel mir wesentlich schwerer. Die typischen Features vieler Topics waren zwar exklusiv für das jeweilige Topic, waren aber nicht direkt als ein Aspekt Diskussionen zum Thema impfen zu interpretieren. Häufig scheinen hier auch Terme, die z.B. aus Gründen der Formulierung häufiger zusammen auftauchen, ein Topic zu bilden.
  - Der qualitative Vergleich der beiden Modelle legt für mich nahe, dass ein nützliches Modell eher in der Region etwas über  $k = 30$  zu finden

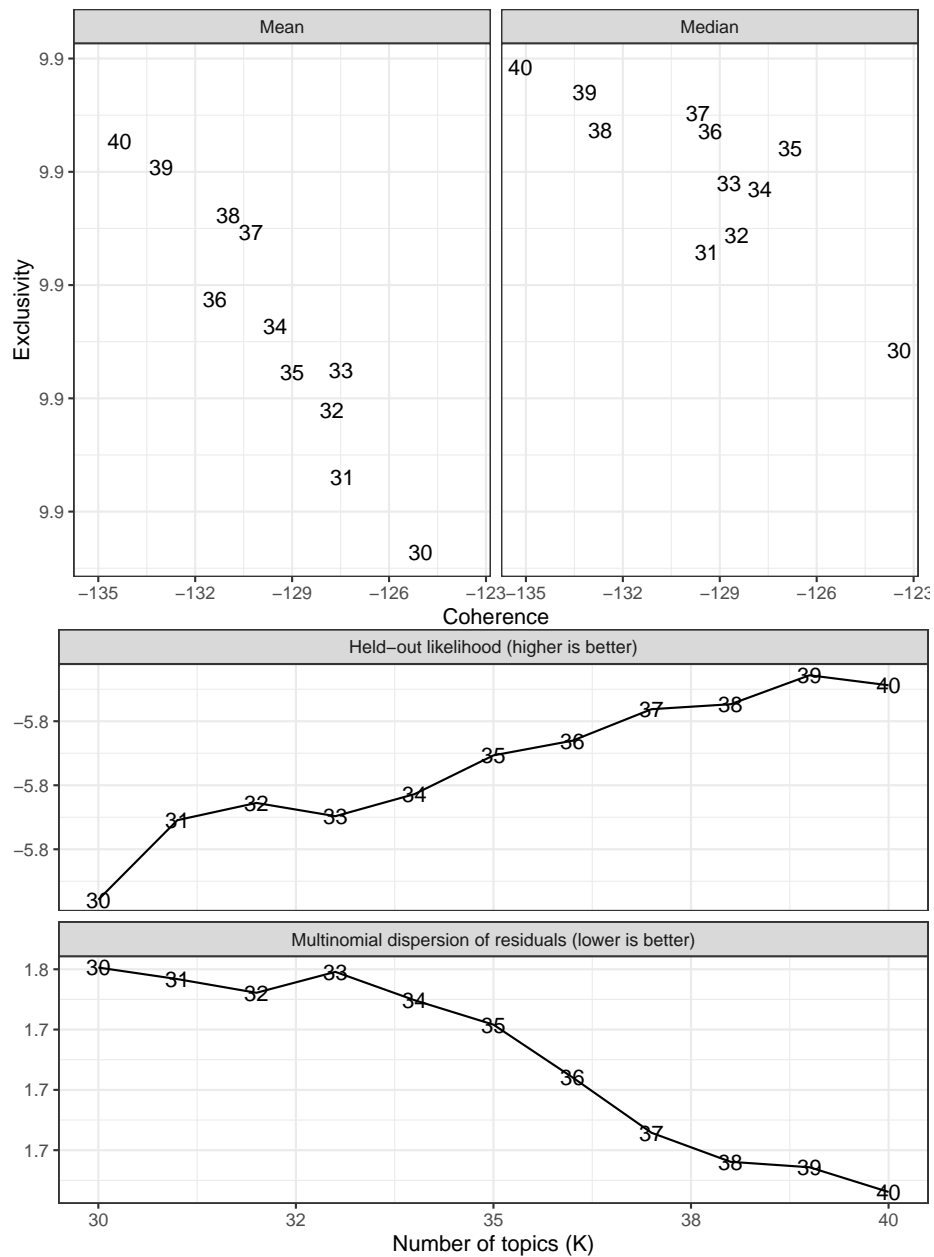


ist als bei  $k = 60$ . Der Grund dafür ist, dass ich meist Modelle mit möglichst vielen interpretierbare Topics bevorzuge. Wenn — wie in diesem Beispiel — die Dokumente bereits so ausgewählt wurden, dass die meisten Dokumente prinzipiell für das Forschungsinteresse relevant sind, halte ich es für sinnvoll, über einen möglichst großen Anteil der inhaltlichen Varianz in diesen Dokumenten etwas auszusagen. Topics, die aus der Analyse ausgeschlossen werden müssen, helfen nicht dabei, die Inhalte der Dokumente zu beschreiben und zu verstehen.

- Beachtet aber, dass das nicht die einzige sinnvolle Vorgehensweise ist. Topic Models enthalten fast immer einige Topics, die sich nicht mit Blick auf das Forschungsinteresse interpretieren lassen (?). Diese werden dann häufig von den weiteren Analysen ausgeschlossen. Die Analysen der interpretierbaren Topics sind davon nicht weiter beeinträchtigt, solange die nicht interpretierbaren Topics keine relevanten Informationen aus den interpretierbaren Topic entfernen. Wir könnten nach dieser Logik also auch das Modell mit  $k = 60$  zum Ausgangspunkt der Suche nach dem nützlichsten Modell machen und viele Topics ausschließen. Die Argumentation wäre dann, dass die verbliebenen Topics besonders klar herausstechen, da die nicht relevanten Charakteristika der Dokumente in den nicht relevanten Topics stecken.
- Aus diesem Grund habe ich als nächsten Schritt das Modell mit  $k = 40$  zum Vergleich herangezogen. Dabei bestätigt sich nach meiner Wahrnehmung die Vermutung, dass (etwas) mehr als 30 Topics eine geeignete Wahl sein könnten.

### 3.2.4 Schätzen und Vergleich von weiteren Modellen zwischen $k = 30$ und $k = 40$

- Die Spezifikation und der quantitative Vergleich erfolgt nach demselben Prinzip wie zuvor. Der Code ist daher hier nicht dargestellt, kann aber in `R/model_compare2.R` gefunden werden. Wir vergleichen alle Modelle zwischen  $k = 30$  und  $k = 40$ .



- Ein Blick auf die Skalierung der Achsen zeigt, dass die Modelle sich untereinander nur wenig unterscheiden. Die quantitativen Metriken helfen uns bei der Entscheidung zwischen diesen Modellen nicht weiter.
- Es folgt ein recht aufwändiger qualitativer Prozess, in dem wir entscheiden, ob die zusätzlichen Themen die Beschreibung der Dokumentinhalte substantiell verbessern.