



# Text as Data

---

**SICSS, May 13, 2025**

Kristen Scotti



# What is text analysis?

---

- Turning unstructured text into structured data for analysis
- Uses computational methods to extract patterns, meaning, or trends from text
- Combines linguistics, statistics, and programming

Examples:

- Analyzing social media sentiment
- Tracking topics in newspaper archives
- Identifying key themes in interview transcripts



# Learning objectives

---

- Understand the basic steps of a text analysis pipeline
- Load, clean, and prepare textual data
- Apply simple natural language processing (NLP) techniques
- Visualize word patterns and frequency
- Think critically about interpretation and bias in text data

A decorative plaid pattern with intersecting lines in red, green, yellow, and blue on a dark background, located on the left side of the slide.

# Text preprocessing

---

- Preprocessing = cleaning and preparing text before analysis
- Typical tasks:
  - Remove whitespace, punctuation, and HTML
  - Normalize text (e.g., lowercase)
  - Tokenize words
  - Remove stopwords



# Why does it matter?

```
duke_web_scrape <- "Duke Experts: A Trusted Source for Policy Makers\n\n\t\t\t\t\t"
```

# Why does it matter?

```
duke_web_scrape <- "Duke Experts: A Trusted Source for Policy Makers\n\n\t\t\t\t\t"
```

```
gsub("\\t", "", duke_web_scrape)
```

```
duke_web_scrape <- "Duke Experts: A Trusted Source for Policy Makers\n\n\t\t\t\t\t"
```

```
gsub("\\t", "", duke_web_scrape)
```

```
gsub("\\t|\\n", "", duke_web_scrape)
```



# What is GREP?

---

- GREP = Globally search a Regular Expression and Print
- Used for pattern matching
- Functions: `grep()`, `grepl()`
  - `grep()`
    - returns the positions (indices) of matches in an input vector
  - `grepl()`
    - returns a logical vector (True/False)



# GREP Example

---

- What do you get when you run this code?

```
{r}  
grep1("Experts", duke_web_scrape)
```

- Change "grep1" to "grep" – what do you get now?

# GREP Example

---

- What do you get when you run this code?

```
{r}  
some_text <- c("apple", "banana", "grape", "apricot")  
grep("ap", some_text)
```

- Change "grepl" to "grep" – what do you get now?

# GREP Example

---

- What do you get when you run this code?

```
## {r}  
some_text<-c("Friends","don't","let","friends","make","wordclouds")  
some_text[grepl("^[F]", some_text)]
```

## GREP Example - Special Characters

---

- What do you get when you run this code?

```
## {r}  
text_chunk<-c("[R is FUN!]")  
gsub("[", "", text_chunk)
```

```
## {r}  
text_chunk<-c("[R is FUN!]")  
gsub('\\[\\]', "", text_chunk)
```



# Regex Cheat Sheet

- <https://github.com/rstudio/cheatsheets/blob/main/regex.pdf>

### Basic Regular Expressions in R

Cheat Sheet

Character Classes	
<code>[[digit:]]</code> or <code>\d</code>	Digits; <code>[0-9]</code>
<code>\D</code>	Non-digits; <code>[^0-9]</code>
<code>[[lower:]]</code>	Lower-case letters; <code>[a-z]</code>
<code>[[upper:]]</code>	Upper-case letters; <code>[A-Z]</code>
<code>[[alpha:]]</code>	Alphabetic characters; <code>[a-zA-Z]</code>
<code>[[alnum:]]</code>	Alphanumeric characters; <code>[A-Za-z0-9_]</code>
<code>\w</code>	Word characters; <code>[A-Za-z0-9_]</code>
<code>\W</code>	Non-word characters
<code>[[xdigit:]]</code> or <code>\x</code>	Hexadic digits; <code>[0-9A-Fa-f]</code>
<code>[[blank:]]</code>	Space and tab
<code>[[space:]]</code> or <code>\s</code>	Space, tab, vertical tab, newline, form feed, carriage return
<code>\S</code>	Not space; <code>[^[:space:]]</code>
<code>[[punct:]]</code>	Punctuation characters; <code>!"#\$%&amp;'()*+,-./:;&lt;=&gt;?@[\^_`{ }~</code>
<code>[[graph:]]</code>	Graphical characters; <code>[[!alnum:]][!punct:]]</code>
<code>[[print:]]</code>	Printable characters; <code>[[!alnum:]][!punct:]]\s</code>
<code>[[cntrl:]]</code> or <code>\c</code>	Control characters; <code>\n</code> , <code>\r</code> etc.

Special Metacharacters	
<code>\n</code>	New line
<code>\r</code>	Carriage return
<code>\t</code>	Tab
<code>\v</code>	Vertical tab
<code>\f</code>	Form feed

Lookaheads and Conditionals*	
<code>(?=)</code>	Lookahead (requires <code>PERL = TRUE</code> ), e.g. <code>(?=xyz)</code> : position followed by "xyz"
<code>(?!)</code>	Negative lookahead (PERL = TRUE); position NOT followed by pattern
<code>(?&lt;=)</code>	Lookbehind (PERL = TRUE), e.g. <code>(?&lt;=xyz)</code> : position following "xyz"
<code>(?&lt;!)</code>	Negative lookbehind (PERL = TRUE); position NOT following pattern
<code>?(&lt;i&gt;ifthen</code>	If-then-condition (PERL = TRUE); use lookaheads, optional char, etc in if-clause
<code>?(&lt;i&gt;ifthenelse</code>	If-then-else-condition (PERL = TRUE)

\*see, e.g., <http://www.regular-expressions.info/lookaround.html>  
<http://www.regular-expressions.info/conditional.html>

### Functions for Pattern Matching

Detect pattern

Locate pattern

Extract pattern

Replace pattern

Extract Patterns

Locate Patterns

Replace Patterns

Split a String using a Pattern

Character Classes and Groups

Anchors

Quantifiers

Escaping Characters

Case Conversions

Greedy Matching

Note

A decorative plaid pattern with diagonal lines in red, green, and yellow on a dark blue background, located on the left side of the slide.

# What is tokenization?

---

- Tokenization = splitting text into meaningful units
  - Characters
  - Words (most common)
  - Sentences
  - N-grams: sequences of words of length  $n$ ; for, "this is a sentence"
    - $N = 1$  (unigrams): this, is, a, sentence
    - $N = 2$  (bigrams): This is, is a, a sentence
    - $N = 3$  (trigrams): This is a, is a sentence
    - N-grams can be useful when word-order is important; nuances in sequences of words

# Let's tokenize...

```
```{r}
sampleText <- "R is great for data analysis. It's also free!"
```
```

```
```{r}
data_frame(line = 1, text = sampleText) %>%
  unnest_tokens(word, text)
#unnest_tokens(character, text, token = "characters")
#unnest_tokens(sentence, text, token = "sentences")
#unnest_tokens(bigram, text, token = "ngrams", n = 2)
#unnest_tokens(trigram, text, token = "ngrams", n = 3)
```
```

| line  | word     |
|-------|----------|
| <dbl> | <chr>    |
| 1     | r        |
| 1     | is       |
| 1     | great    |
| 1     | for      |
| 1     | data     |
| 1     | analysis |
| 1     | it's     |
| 1     | also     |
| 1     | free     |

A decorative plaid pattern with red, green, and yellow lines on a dark blue background, located on the left side of the slide.

# Exercise

---

- Tokenize the sentence: "Learning R is fun and empowering!" into bigrams and trigrams.
- How does the meaning or context change based on how your unit of analysis?



# Exercise

---

```
```{r}
# Write your code here!
text <- "Learning R is fun and empowering!"
tibble(line = 1, text = text) %>%
  unnest_tokens(word, text)

# Bigrams
unnest_tokens(tibble(line = 1, text = text), bigram, text, token = "ngrams", n = 2)
# Trigrams
unnest_tokens(tibble(line = 1, text = text), trigram, text, token = "ngrams", n = 3)

```
```

```
data.frame(line = 1, text = "Learning R is fun and empowering!") %>%
  unnest_tokens(word, text)
```



## Exercise

---

"Learning R"

"R is"

"is fun"

"fun and"

"and empowering"

"Learning R is"

"R is fun"

"is fun and"

"fun and empowering"



## How does `unnest_tokens()` handle punctuation or emojis?

---

- Punctuation marks are removed automatically (unless you customize the tokenizer)
- All text is converted to lowercase
- Emojis are kept if they are treated as characters or part of the Unicode word class, depending on the tokenization method




## How does `unnest_tokens()` handle non-English characters?

---

- Tries to follow the Unicode standard, meaning:
- Accented characters (e.g., é, ñ, ü) and most Latin script characters are handled correctly
- Languages without spaces (like Chinese or Japanese) won't be tokenized correctly unless you use a custom tokenizer or a different package (e.g., jiebaR for Chinese or tokenizers.bpe for multilingual support).





# What is a Corpus?

---

- A corpus is a collection of documents
- Tweets, articles, reviews, transcripts, ...

## Load some data...

---

```
## {r}  
original_data <- read.csv("Corona_NLP_train.csv") %>%  
  select(created_at = TweetAt, text = OriginalTweet)  
head(original_data)
```

**After you load your data, tokenize it by word!**

# Tokenize

---

```
```{r}
# Write your code here!

tidy_covid <- original_data %>%
  unnest_tokens(word, text)

head(tidy_covid)
```
```

|   | <b>created_at</b><br><chr> | <b>word</b><br><chr> |
|---|----------------------------|----------------------|
| 1 | 3/16/2020                  | menyrbie             |
| 2 | 3/16/2020                  | phil_gahan           |
| 3 | 3/16/2020                  | chrisitv             |
| 4 | 3/16/2020                  | https                |
| 5 | 3/16/2020                  | t.co                 |
| 6 | 3/16/2020                  | ifz9fan2pa           |

# Counting

---

```
tidy_covid %>%  
  count(word) %>%  
  arrange(desc(n))
```

| word<br><chr> | n<br><int> |
|---------------|------------|
| the           | 29744      |
| to            | 25785      |
| and           | 16344      |
| https         | 15143      |
| t.co          | 15132      |
| of            | 14444      |
| coronavirus   | 12929      |
| a             | 12767      |
| in            | 12556      |
| for           | 9413       |



# Remove stop words

```
```{r}
data("stop_words")

tidy_covid <- tidy_covid %>%
  anti_join(stop_words)

head(tidy_covid)
```
```

```
```{r}
tidy_covid %>%
  count(word) %>%
  arrange(desc(n))
```
```

| word<br><chr> | n<br><int> |
|---------------|------------|
| https         | 15143      |
| t.co          | 15132      |
| coronavirus   | 12929      |
| 19            | 7948       |
| covid         | 7583       |
| supermarket   | 4990       |
| food          | 4963       |
| prices        | 4942       |
| store         | 4916       |
| grocery       | 4340       |

## Remove Noise (Links, Twitter Junk)

---

```
~~~{r}
tidy_covid <- tidy_covid %>%
  filter(!word %in% c("https", "rt", "t.co", "amp"))
~~~
```

```
~~~{r}
tidy_covid %>%
  count(word) %>%
  arrange(desc(n))
~~~
```

| word<br><chr> | n<br><int> |
|---------------|------------|
| coronavirus   | 12929      |
| 19            | 7948       |
| covid         | 7583       |
| supermarket   | 4990       |
| food          | 4963       |
| prices        | 4942       |
| store         | 4916       |
| grocery       | 4340       |
| people        | 4036       |
| covid19       | 3192       |

# Clean Punctuation, Numbers, and Whitespace

```
```{r}
tidy_covid$word <- gsub("[[:punct:]]", "", tidy_covid$word)
tidy_covid$word <- gsub("\\d+", "", tidy_covid$word)
tidy_covid <- tidy_covid %>% filter(word != "")
```
```

```
```{r}
tidy_covid %>%
  count(word) %>%
  arrange(desc(n))
```
```

| word<br><chr> | n<br><int> |
|---------------|------------|
| covid         | 14400      |
| coronavirus   | 12939      |
| supermarket   | 4990       |
| food          | 4963       |
| prices        | 4944       |
| store         | 4917       |
| grocery       | 4340       |
| people        | 4036       |
| consumer      | 2732       |
| shopping      | 2294       |



## Exercise: Clean a vector

---

```
c("Running!", "n95", "COVID-19", "great!!!", "1234")
```



# Exercise: Clean a vector

`c("Running!", "n95", "COVID-19", "great!!!", "1234")`

**word**

<chr>

running

n

covid

great

```
~~~{r}
# Write your code here!
words <- tibble(word = c("Running!", "n95", "COVID-19", "great!!!", "1234")) %>%
  mutate(word = gsub("[[:punct:]]", "", word),
         word = gsub("\\d+", "", word),
         word = tolower(word))
~~~
```

# Stemming

---

- Reduces words to their roots
- “typing” → “type”
- Use wordStem() from the SnowballC package

```
```{r}
tidy_covid <- tidy_covid %>%
  mutate(word = wordStem(word, language = "en"))
```
```

# Word frequency analysis

---

```
```{r}
top_words <- tidy_covid %>%
  count(word) %>%
  arrange(desc(n))
```
```

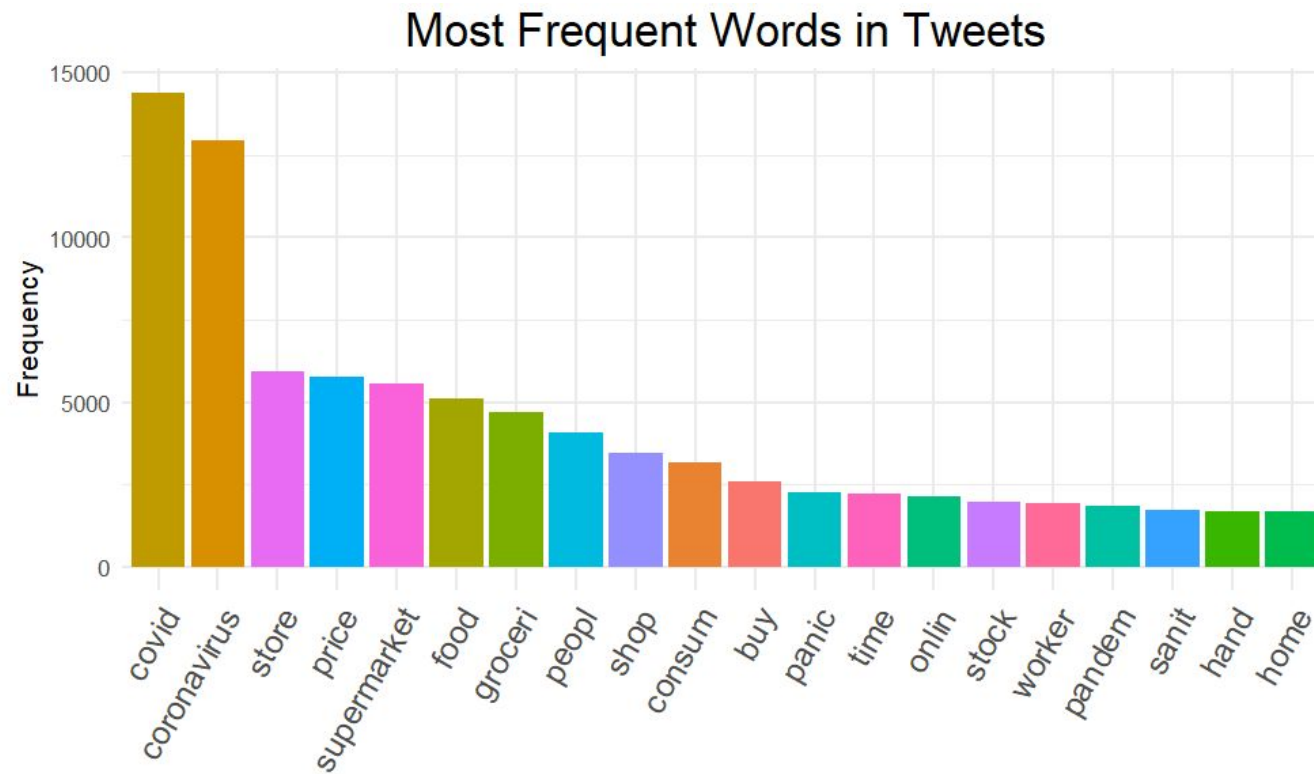
# Visualize...


---

```
...{r}
top_words %>%
  slice(1:20) %>%
  ggplot(aes(x=reorder(word, -n), y=n, fill=word))+
  geom_bar(stat="identity")+
  theme_minimal()+
  theme(axis.text.x =
    element_text(angle = 60, hjust = 1, size=13))+
  theme(plot.title =
    element_text(hjust = 0.5, size=18))+
  ylab("Frequency")+
  xlab("")+
  ggtitle("Most Frequent Words in Tweets")+
  guides(fill=FALSE)
...
```



# Visualize...





# TF-IDF

---

- TF = Term Frequency
- IDF = Inverse Document Frequency
- Highlights words that are important in specific documents

# TF-IDF

---

```
```{r}
tidy_covid_tfidf <- tidy_covid %>%
  count(created_at, word) %>%
  bind_tf_idf(word, created_at, n)
```
```

# TF-IDF

---

Now let's see what the most unusual words are (highest TF-IDF score)

```
```{r}
top_tfidf <- tidy_covid_tfidf %>%
  arrange(desc(tf_idf))
```

```
top_tfidf$word[1:10]
```

```
[1] "glasgow"          "stayathomeandstaysaf" "ajqfcbjb"
[4] "apmacanada"       "bigmart"              "bzuivvzu"
[7] "cmoldmak"         "cmqrsxnuv"            "cpykboekvt"
[10] "creativepapilo"
```



# TF-IDF

---

Least distinctive word

```
```{r}
top_tfidf <- tidy_covid_tfidf %>%
  arrange(tf_idf)

# View the top tf-idf word
top_tfidf$word[1]
```

# TF-IDF

Show top N most usual/unusual words with their TF-IDF score

```
```{r}
top_tfidf <- tidy_covid_tfidf %>%
  # arrange(desc(tf_idf)) %>%
  arrange(tf_idf) %>%
  select(word, tf_idf)

# View the top 10 words and their scores
head(top_tfidf, 10)
```
```

|    | word<br><chr> | tf_idf<br><dbl> |
|----|---------------|-----------------|
| 1  | basic         | 0               |
| 2  | call          | 0               |
| 3  | chines        | 0               |
| 4  | citi          | 0               |
| 5  | come          | 0               |
| 6  | communiti     | 0               |
| 7  | consum        | 0               |
| 8  | contin        | 0               |
| 9  | coronavirus   | 0               |
| 10 | covid         | 0               |



# Exercise

---

Compare top TF-IDF terms by date

# Exercise

---

Compare top TF-IDF terms by date

```
## {r}
data <- tibble(date = c("2020-01-01", "2020-01-01", "2020-01-02", "2020-01-02"),
               word = c("mask", "covid", "vaccine", "covid"))
data %>% count(date, word) %>%
  bind_tf_idf(word, date, n) %>%
  arrange(desc(tf_idf))
##
```



# Exercise

---

Compare top TF-IDF terms by date

| <b>date</b><br><chr> | <b>word</b><br><chr> | <b>n</b><br><int> | <b>tf</b><br><dbl> | <b>idf</b><br><dbl> | <b>tf_idf</b><br><dbl> |
|----------------------|----------------------|-------------------|--------------------|---------------------|------------------------|
| 2020-01-01           | mask                 | 1                 | 0.5                | 0.6931472           | 0.3465736              |
| 2020-01-02           | vaccine              | 1                 | 0.5                | 0.6931472           | 0.3465736              |
| 2020-01-01           | covid                | 1                 | 0.5                | 0.0000000           | 0.0000000              |
| 2020-01-02           | covid                | 1                 | 0.5                | 0.0000000           | 0.0000000              |

# Sentiment Analysis

- Use `get_sentiments("bing")`
- Join with tokenized text
- Count sentiment by date

```
{r}  
head(get_sentiments("bing"))
```

A tibble: 6 × 2

| <b>word</b><br><chr> | <b>sentiment</b><br><chr> |
|----------------------|---------------------------|
| 2-faces              | negative                  |
| abnormal             | negative                  |
| abolish              | negative                  |
| abominable           | negative                  |
| abominably           | negative                  |
| abominate            | negative                  |

6 rows

## Apply the Bing sentiment dictionary to our dataset

---

```
## {r}
covid_sentiment <- tidy_covid %>%
  inner_join(get_sentiments("bing")) %>%
  count(created_at, sentiment)

head(covid_sentiment)
```

```
## {r}
tidy_covid$date <- as.Date(tidy_covid$created_at, format = "%m/%d/%Y")
```



## Apply the Bing sentiment dictionary to our dataset

---

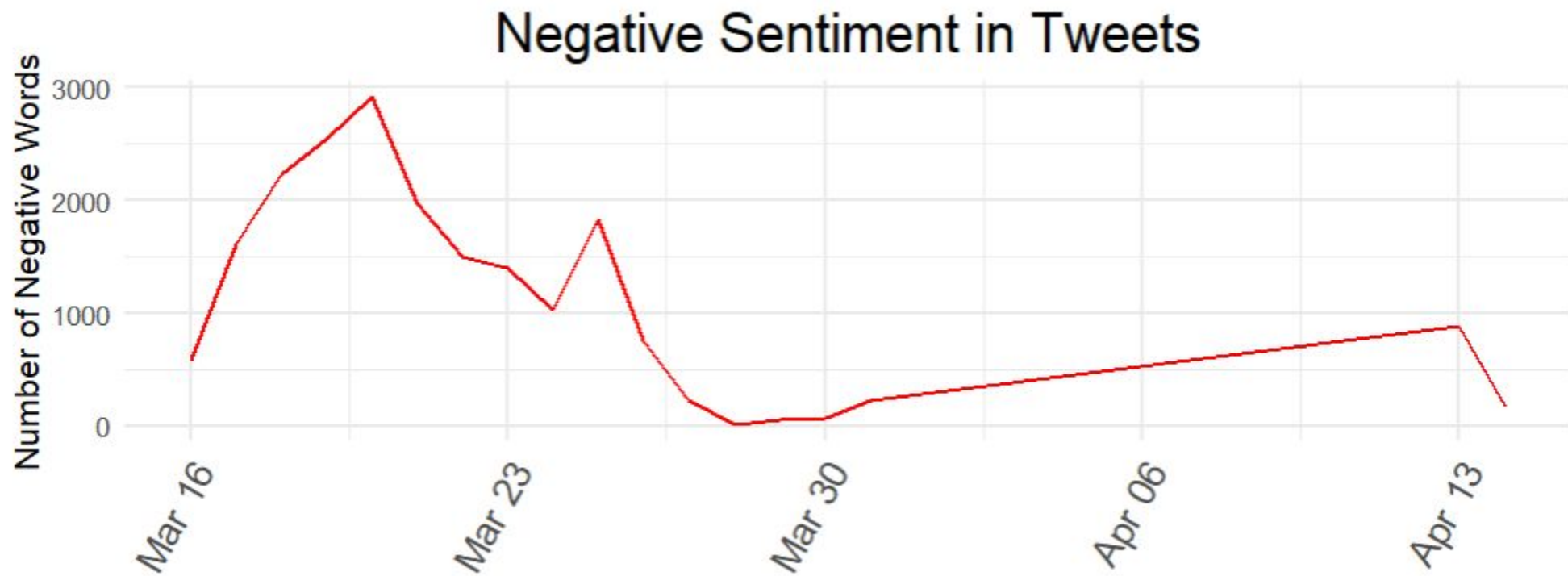
```
{r}  
covid_sentiment_plot <-  
  tidy_covid %>%  
    inner_join(get_sentiments("bing")) %>%  
    filter(sentiment=="negative") %>%  
    count(date, sentiment)  
...  
}
```



# Plot

```
~~~{r}
ggplot(covid_sentiment_plot, aes(x=date, y=n))+
  geom_line(color="red", size=.5)+
  theme_minimal()+
  theme(axis.text.x =
    element_text(angle = 60, hjust = 1, size=13))+
  theme(plot.title =
    element_text(hjust = 0.5, size=18))+
  ylab("Number of Negative Words")+
  xlab("")+
  ggtitle("Negative Sentiment in Tweets")+
  theme(aspect.ratio=1/4)
~~~
```

# Plot





# Questions...

---

- Can one tweet have both positive and negative sentiment?
- What if a word isn't in the lexicon?

# Dictionary-Based Analysis

---

```
##{r}  
economic_dictionary <- c("economy", "unemployment", "trade", "tariffs", "shopping")  
##
```



# Dictionary-Based Analysis

---

Now use `stringr::str_detect()` to filter tweets that mention any of these terms:

```
library(stringr)
economic_tweets <- original_data %>%
  filter(str_detect(text, paste(economic_dictionary, collapse = "|")))
head(economic_tweets)
```

# Dictionary-Based Analysis

```
```{r}
# Filter original tweets to only those that match your economic dictionary
economic_tweets <- original_data %>%
  filter(str_detect(text, paste(economic_dictionary, collapse = "|")))

# Then tokenize only these tweets
tidy_economic <- economic_tweets %>%
  unnest_tokens(word, text)

# Now join to Bing and analyze sentiment
economic_sentiment <- tidy_economic %>%
  inner_join(get_sentiments("bing")) %>%
  count(created_at, sentiment)

# Optionally, plot only negative sentiment
economic_sentiment_plot <- economic_sentiment %>%
  filter(sentiment == "negative") %>%
  mutate(date = as.Date(created_at, format = "%m/%d/%Y")) # fix format if needed

ggplot(economic_sentiment_plot, aes(x = date, y = n)) +
  geom_line(color = "red", linewidth = 0.5) +
  theme_minimal() +
  labs(x = "", y = "Negative Words", title = "Negative Sentiment in Economic Tweets") +
  theme(axis.text.x = element_text(angle = 60, hjust = 1))
```
```



## Negative Sentiment in Economic Tweets





## Exercise: Create your own dictionary

---

- e.g., healthcare, misinformation, etc.
- Filter or visualize sentiment over time
- Compare peaks in negative/positive sentiment





# Topic Modeling

---

- A method to discover hidden themes ("topics") in large collections of text
- Assumes:
  - Each document is a mixture of topics
  - Each topic is a mixture of words
- Helps explore, summarize, or classify text data
- Latent Dirichlet Allocation (LDA)
- Structural Topic Modeling (STM)

## LDA on Associated Press Data

---

- Built-in dataset: AssociatedPress → DocumentTermMatrix
  - Rows = documents
  - Columns = terms/words
  - Values = how many times each word appears in each document
- Fit model with 10 topics)

```
data("AssociatedPress")
```

```
# Fits a Latent Dirichlet Allocation (LDA) model with 10 topics on the AP dataset.  
AP_topic_model<-LDA(AssociatedPress, k=10, control = list(seed = 321))
```

## LDA on Associated Press Data

---

```
# Fits a Latent Dirichlet Allocation (LDA) model with 10 topics on the AP dataset.  
AP_topic_model<-LDA(AssociatedPress, k=10, control = list(seed = 321))
```

## Extract and Visualize Top Terms

---

```
AP_topics <- tidy(AP_topic_model, matrix = "beta")
```

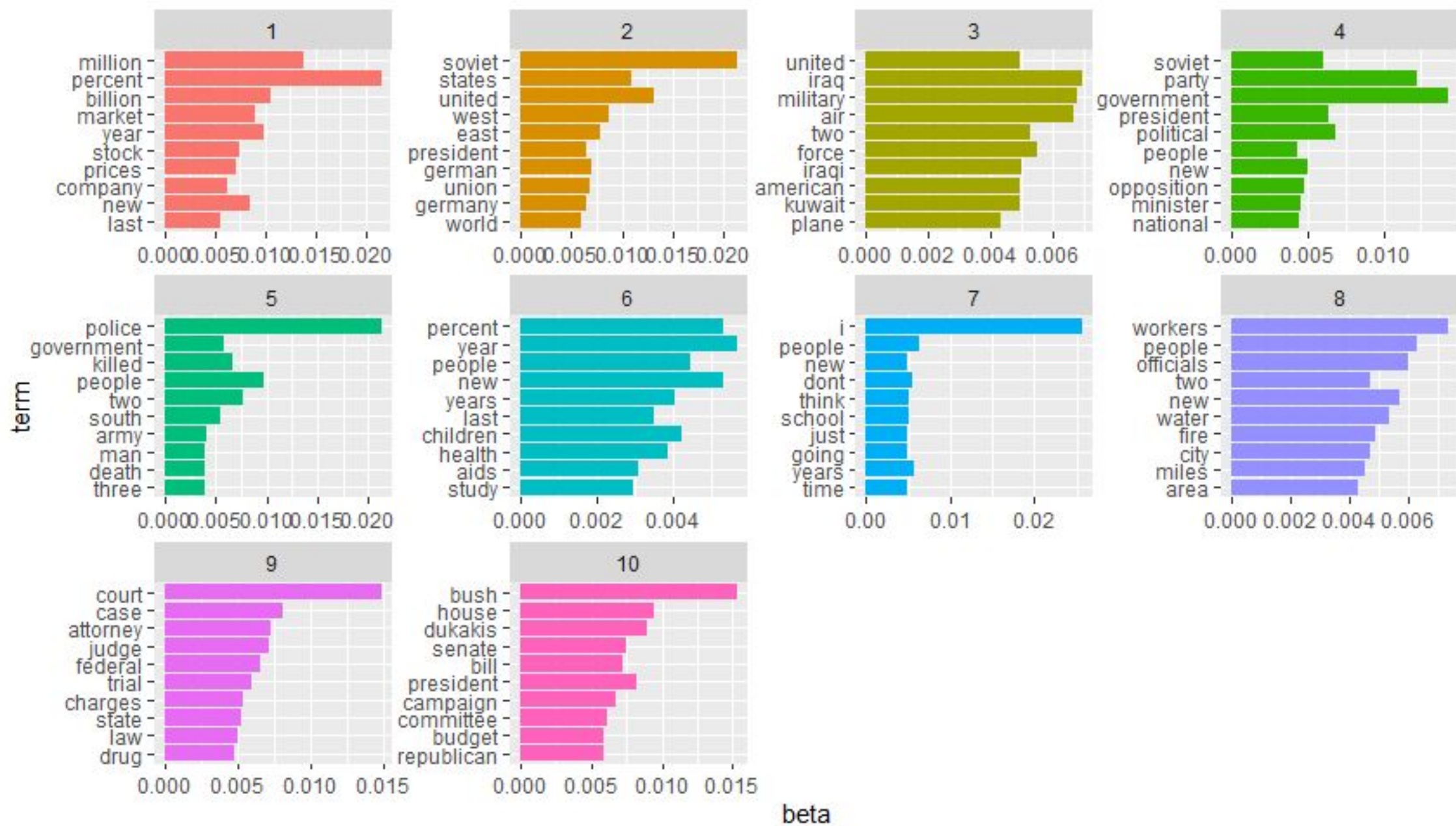
```
# Finds the top 10 terms for each topic, based on the highest beta values (topic-word probabilities).
ap_top_terms <- AP_topics %>%
  group_by(topic) %>%
  top_n(10, beta) %>%
  ungroup() %>%
  arrange(topic, -beta)
```



## Extract and Visualize Top Terms

---

```
ap_top_terms %>%  
  mutate(term = reorder(term, beta)) %>%  
  ggplot(aes(term, beta, fill = factor(topic))) +  
  geom_col(show.legend = FALSE) +  
  facet_wrap(~ topic, scales = "free") +  
  coord_flip()
```



# STM example on Poliblogs

---

- Dataset: 2008 political blog posts
- Preprocess with `textProcessor()`
- Format with `prepDocuments()`

```
google_doc_id <- "1LcX-JnpGB0lU1iDnXnxB6WFqBywUKpew" # google file ID
poliblogs<-read.csv(sprintf("https://docs.google.com/uc?id=%s&export=download",
                             google_doc_id), stringsAsFactors = FALSE)
```



# Preprocessing

---

```
# Preprocess and prepare documents
processed <- textProcessor(poliblogs$documents, metadata = poliblogs)
out <- prepDocuments(processed$documents, processed$vocab, processed$meta)
docs <- out$documents
vocab <- out$vocab
meta <- out$meta
```

Document: A single text entry, like a blog post.

Vocabulary: The list of all words that remain after cleaning.

Metadata: Additional info about each document, such as date or political affiliation.





# What is a covariate?

---

- A covariate is extra information about a document that might explain why certain topics appear more often
- A blog's political rating might influence which topics are more common.
- The date might explain how topics shift over time.



# What is K?

---

- K is the number of topics to extract
- Chosen by YOU — not learned from the data
- Too small → overly broad topics
- Too large → fragmented or incoherent topics
- Use `searchK()` to evaluate values for K

# Finding k

---

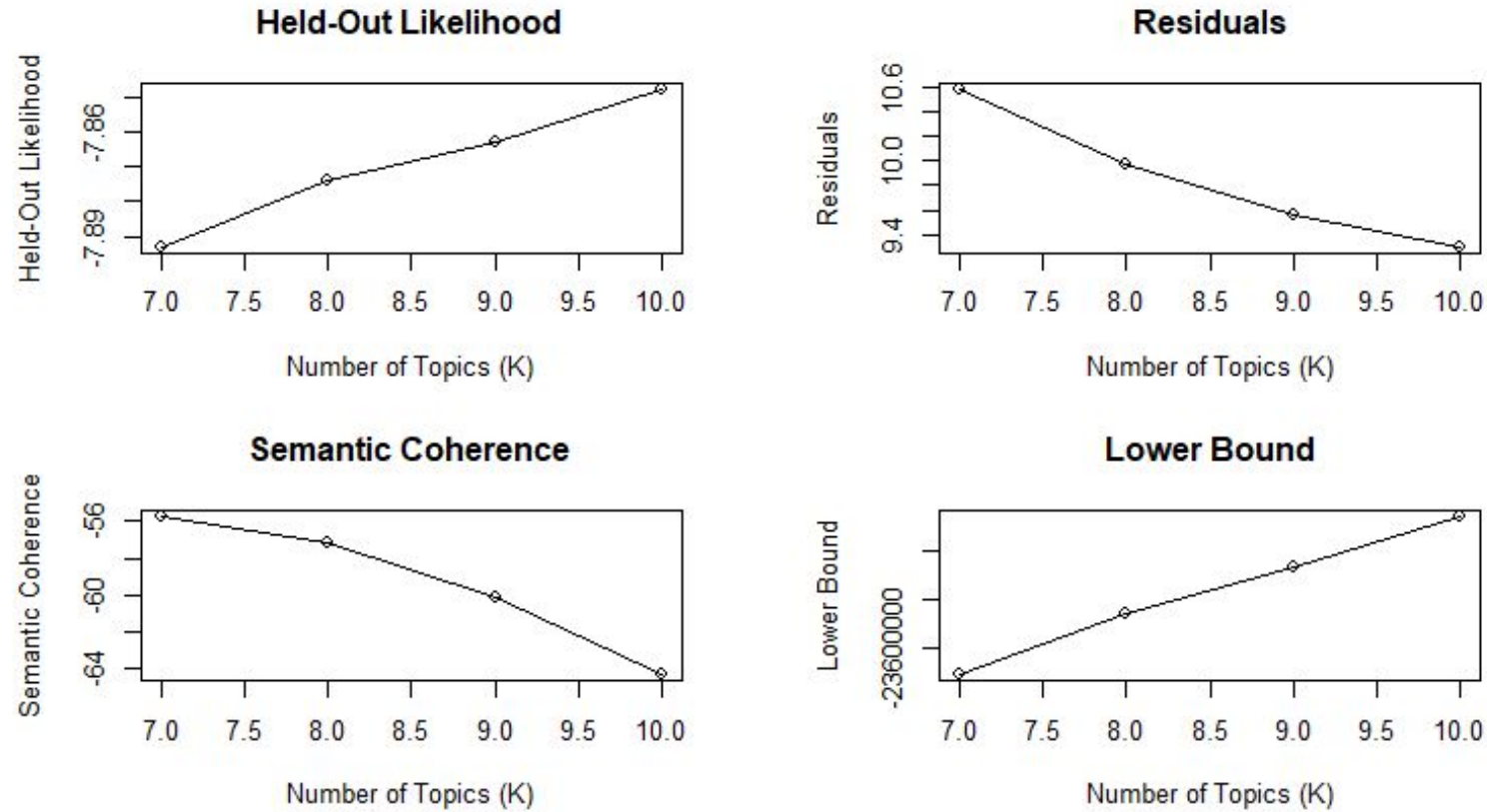
```
# -----  
# OPTION A: Run and Save searchK  
# -----  
# Warning: This process is computationally intensive,  
# Search for Optimal Number of Topics (K)  
findingk <- searchK(  
  documents = docs,  
  vocab = vocab,  
  K = 7:10,  
  prevalence = ~ rating + s(day),  
  data = meta,  
  verbose = FALSE  
)  
save(findingk, file = "findingk.Rda")
```

```
# -----  
# OPTION B: Load precomputed searchK  
# -----  
load("findingk.Rda")  
plot(findingk)
```

prevalence = ~ rating + s(day) includes  
covariates:

- rating: Political leaning
- s(day): A spline, which models non-linear change over time

## Diagnostic Values by Number of Topics

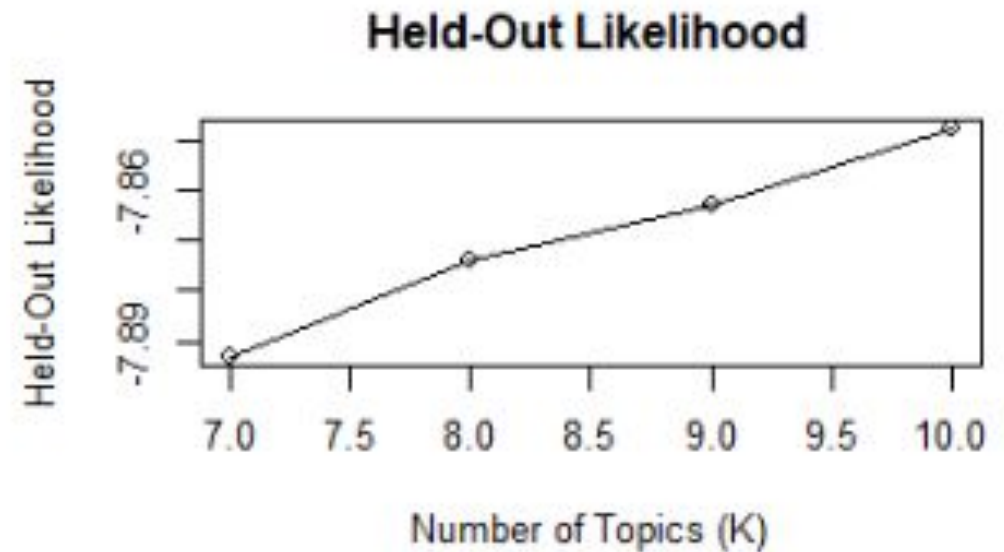




# Held-out likelihood

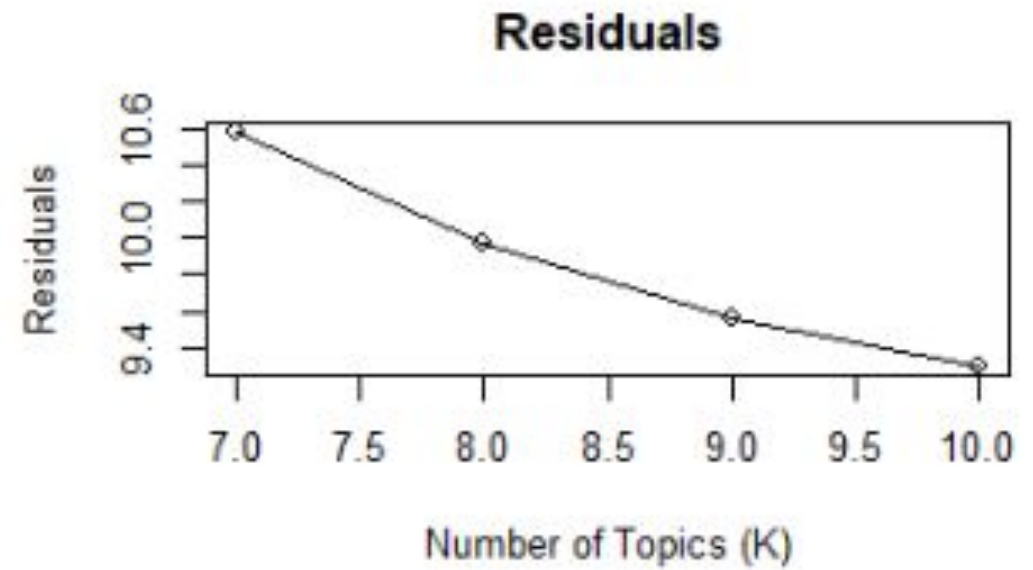
---

- Higher is better
- Increases steadily → 10 is best here



# Residuals

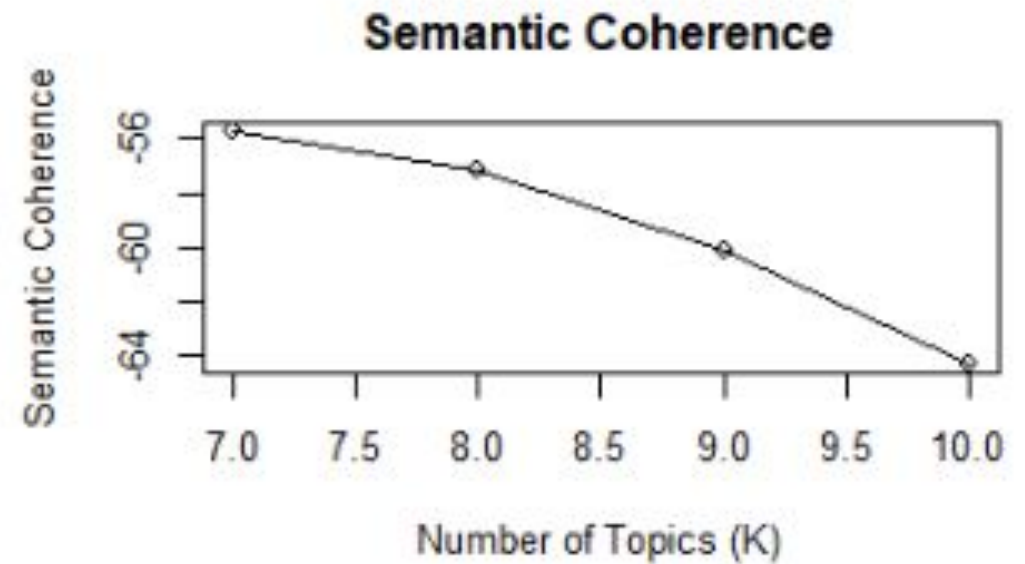
- Higher is better
- Increases steadily → 10 is best here



# Semantic Coherence

---

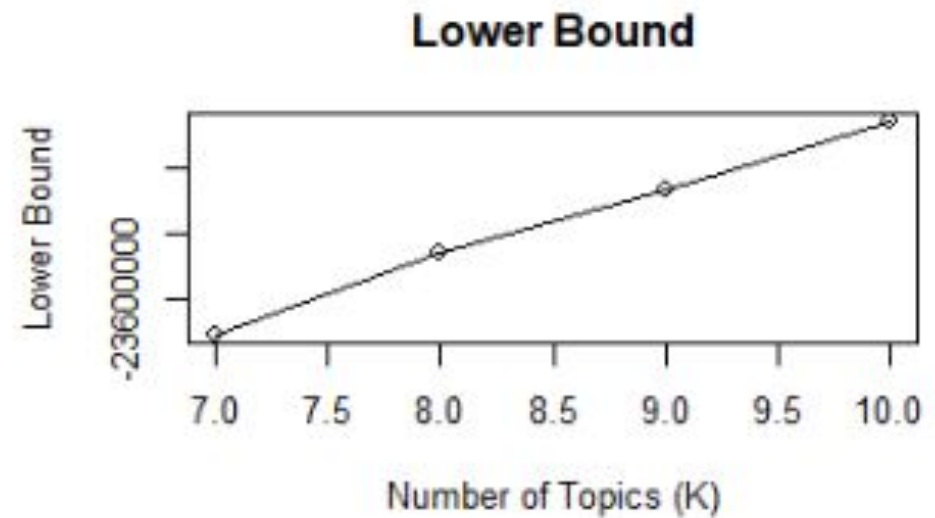
- Higher is better
- Decreases as K increases  
→ 7 is best here



# Lower bound

---

- Less negative is better
- Improves with  $K$ , but used mainly for convergence





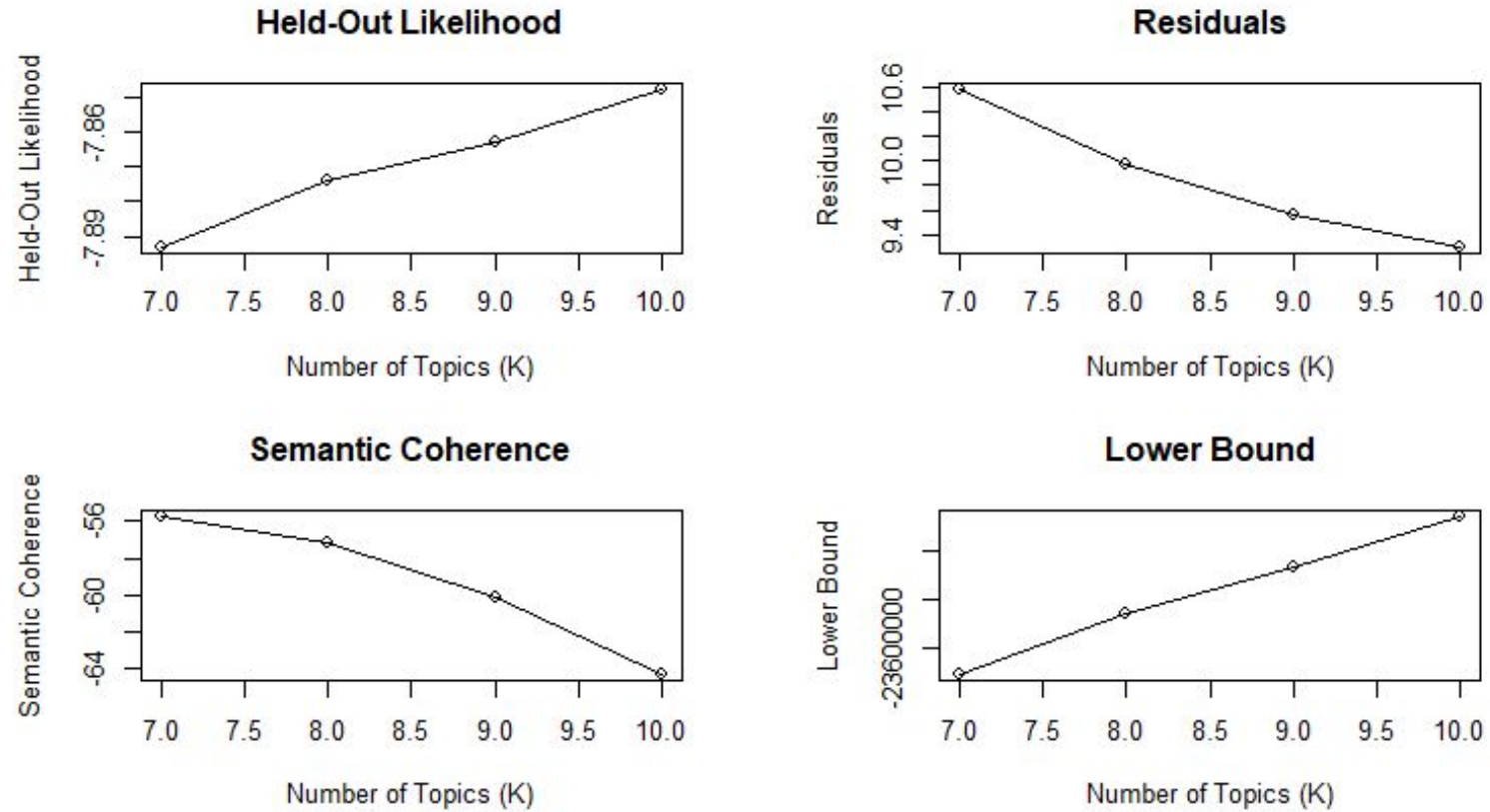


# Trade-offs

---

- Held-Out Likelihood and Residuals both point to  $K = 10$  as the most statistically powerful model—it generalizes best to new data and captures more structure.
- Semantic Coherence drops off significantly after  $K = 7$ , which suggests that topics become less interpretable—i.e., the top words in each topic co-occur less often in documents

## Diagnostic Values by Number of Topics



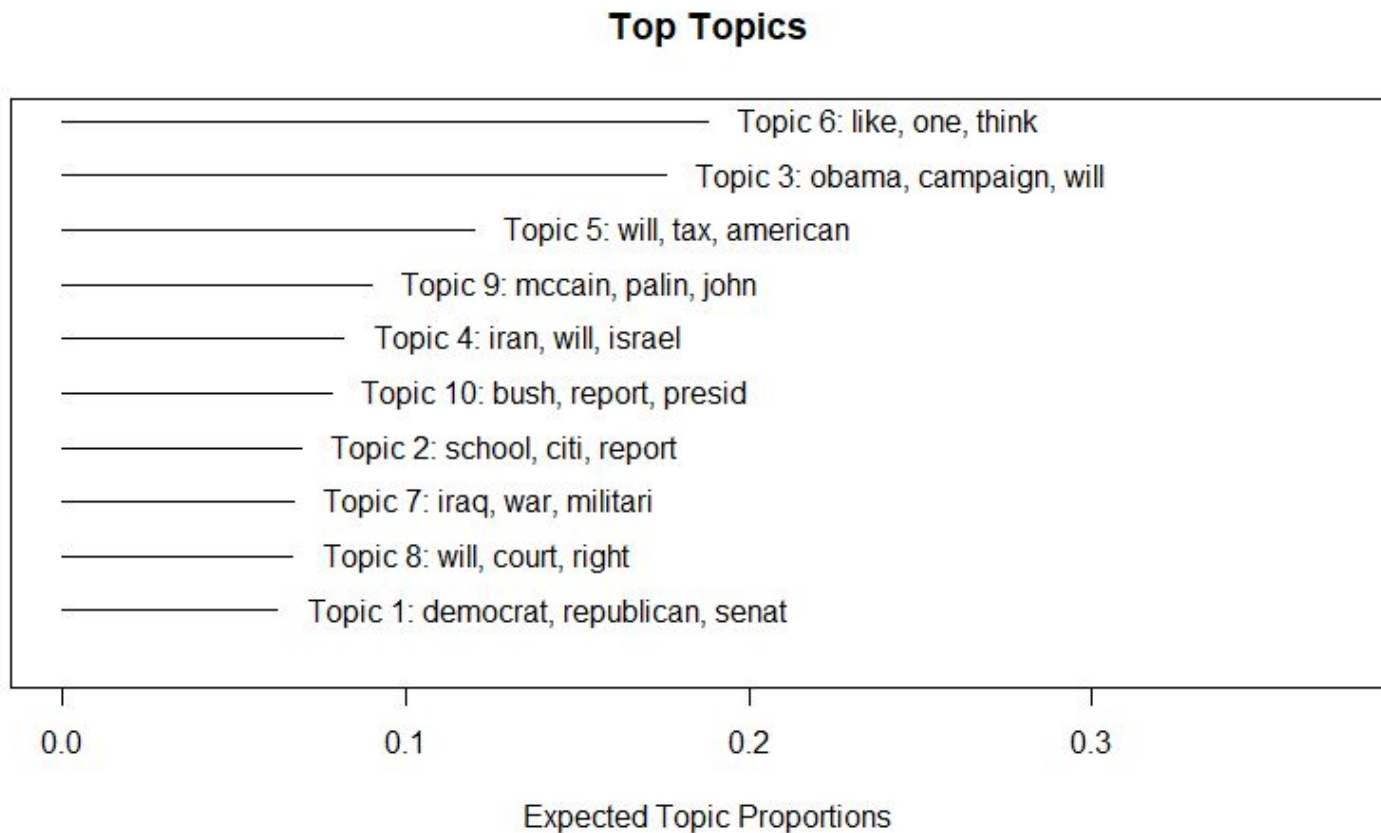
## STM – Fit

---

```
# -----  
# Fit STM with chosen K (e.g., 10)  
# -----  
First_STM <- stm(  
  documents = docs,  
  vocab = vocab,  
  K = 10, # use selected K from searchK  
  prevalence = ~ rating + s(day),  
  max.em.its = 75,  
  data = meta,  
  init.type = "spectral",  
  verbose = FALSE  
)  
  
plot(First_STM)
```

# Top topics

---





## Look at example texts for a topic

---

```
# -----  
# Examine example topics  
# -----  
findThoughts(First_STM, texts = poliblogs$documents,  
              n = 2, topics = 3)
```

# Estimate and plot topic effects

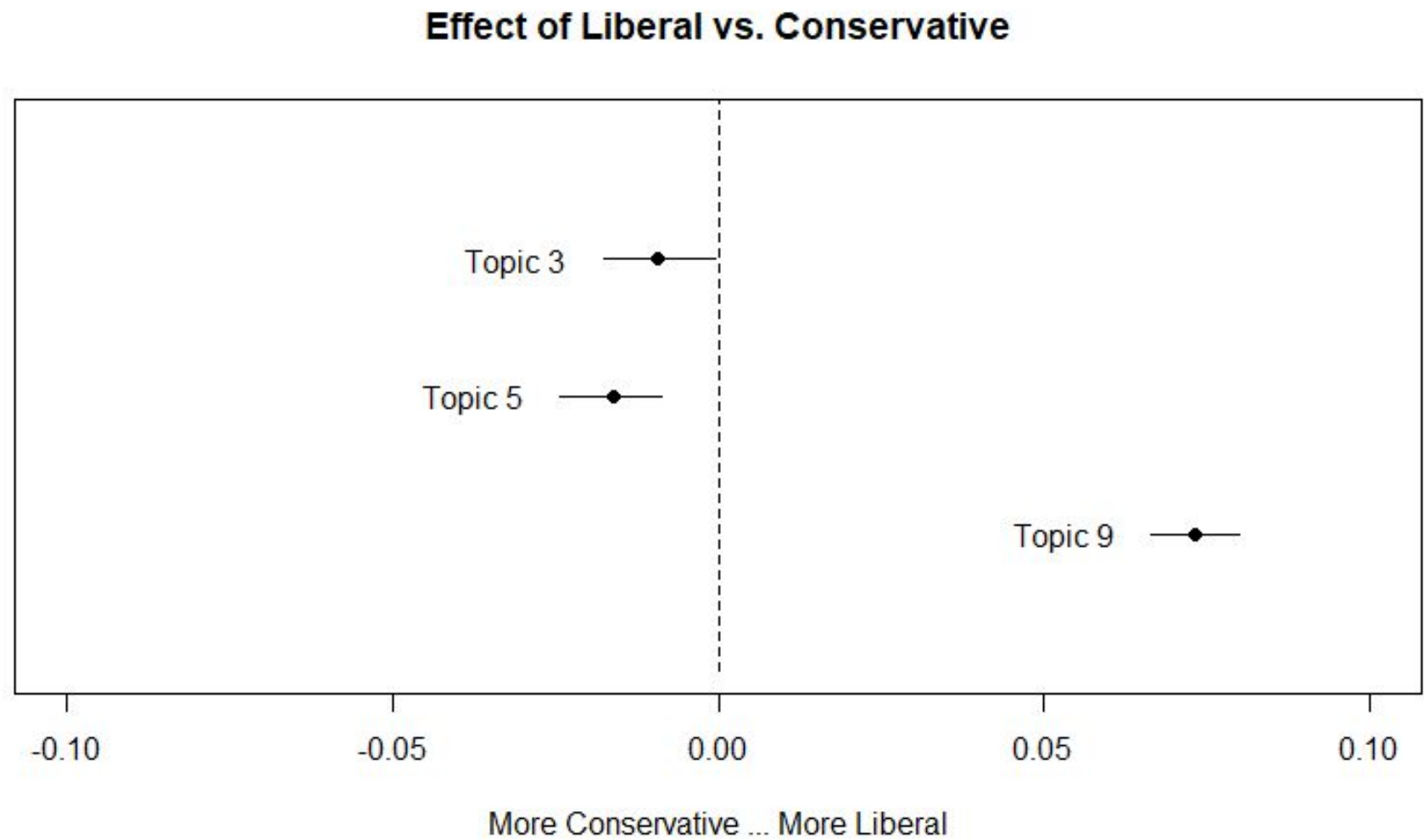
---

```
# -----  
# Estimate and plot topic effects by covariate  
# -----  
predict_topics<-estimateEffect(formula = 1:10 ~ rating + s(day), stmobj = First_STM,  
                               metadata = out$meta, uncertainty = "global")
```

# Visualize Covariate Effects

---

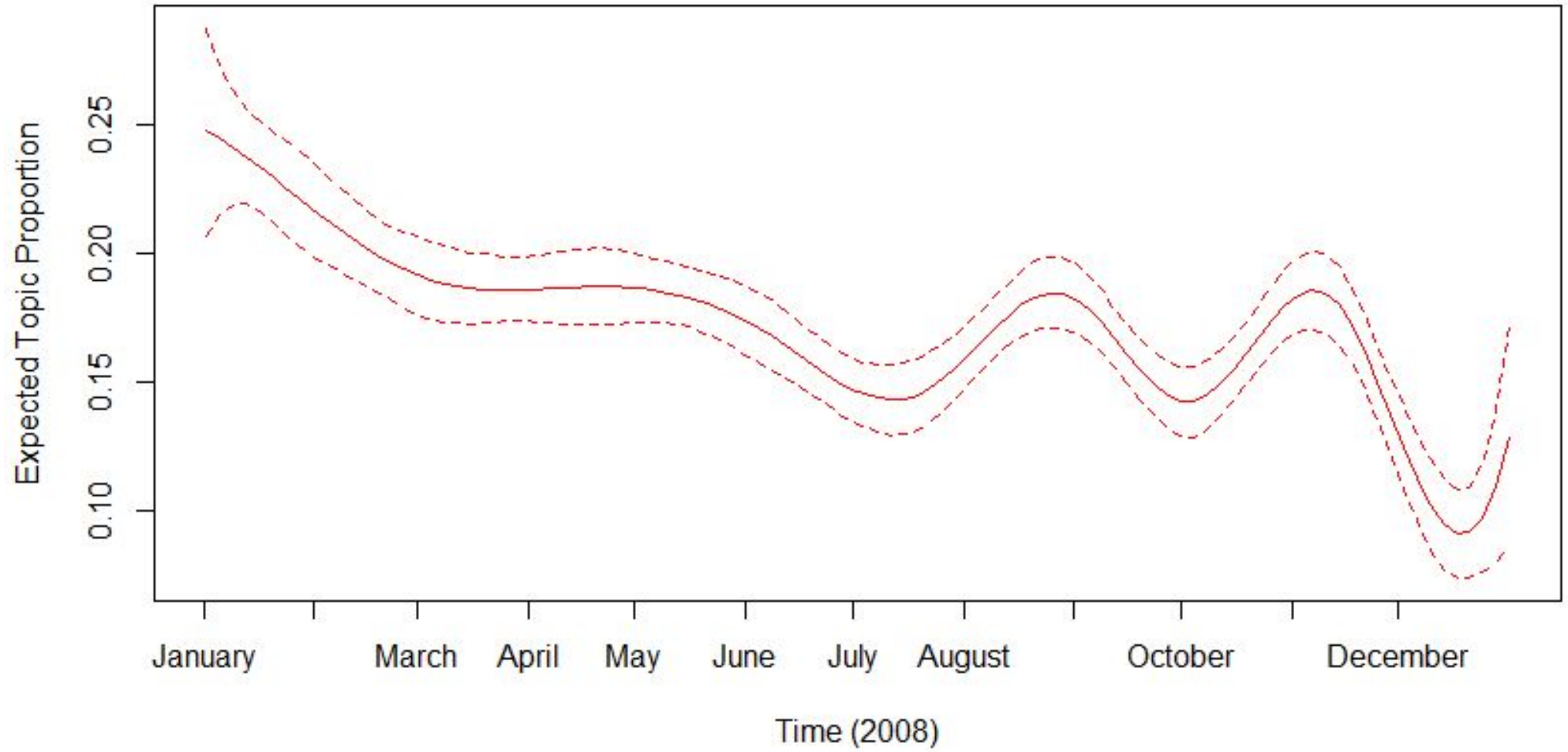
```
# Liberal vs. Conservative
plot(predict_topics, covariate = "rating", topics = c(3, 5, 9),
      model = First_STM, method = "difference",
      cov.value1 = "Liberal", cov.value2 = "Conservative",
      xlab = "More Conservative ... More Liberal",
      main = "Effect of Liberal vs. Conservative",
      xlim = c(-.1, .1), labeltype = "custom",
      custom.labels = c('Topic 3', 'Topic 5', 'Topic 9'))
```





```
# Change over time
plot(predict_topics, "day", method = "continuous", topics = 3,
      model = First_STM, printlegend = FALSE, xaxt = "n", xlab = "Time (2008)")

monthseq <- seq(from = as.Date("2008-01-01"), to = as.Date("2008-12-01"), by = "month")
monthnames <- months(monthseq)
axis(1, at = as.numeric(monthseq) - min(as.numeric(monthseq)), labels = monthnames)
```





# Brainstorming time...

---

- **Topic or Research Question**
  - What kind of text would you analyze? What do you want to know?
- **Approach/Tools**
  - Would you use sentiment analysis? TF-IDF? Topic modeling? Something else?
- **Potential Pitfalls**
  - What challenges might come up with the data, methods, or interpretation?
  - Are there any risks of bias or misinterpretation in your data, methods, or results? e.g., biased training data, missing context, or over-reliance on sentiment dictionaries