

Text Modeling



USING TIDY PRINCIPLES

Julia Silge | SDSS | 29 May 2019

Let's install some packages

```
install.packages(c("tidyverse",
                   "tidytext",
                   "gutenbergr",
                   "stm",
                   "glmnet",
                   "yardstick"))
```



Find us at...

 @juliasilge
 @juliasilge
 juliasilge.com

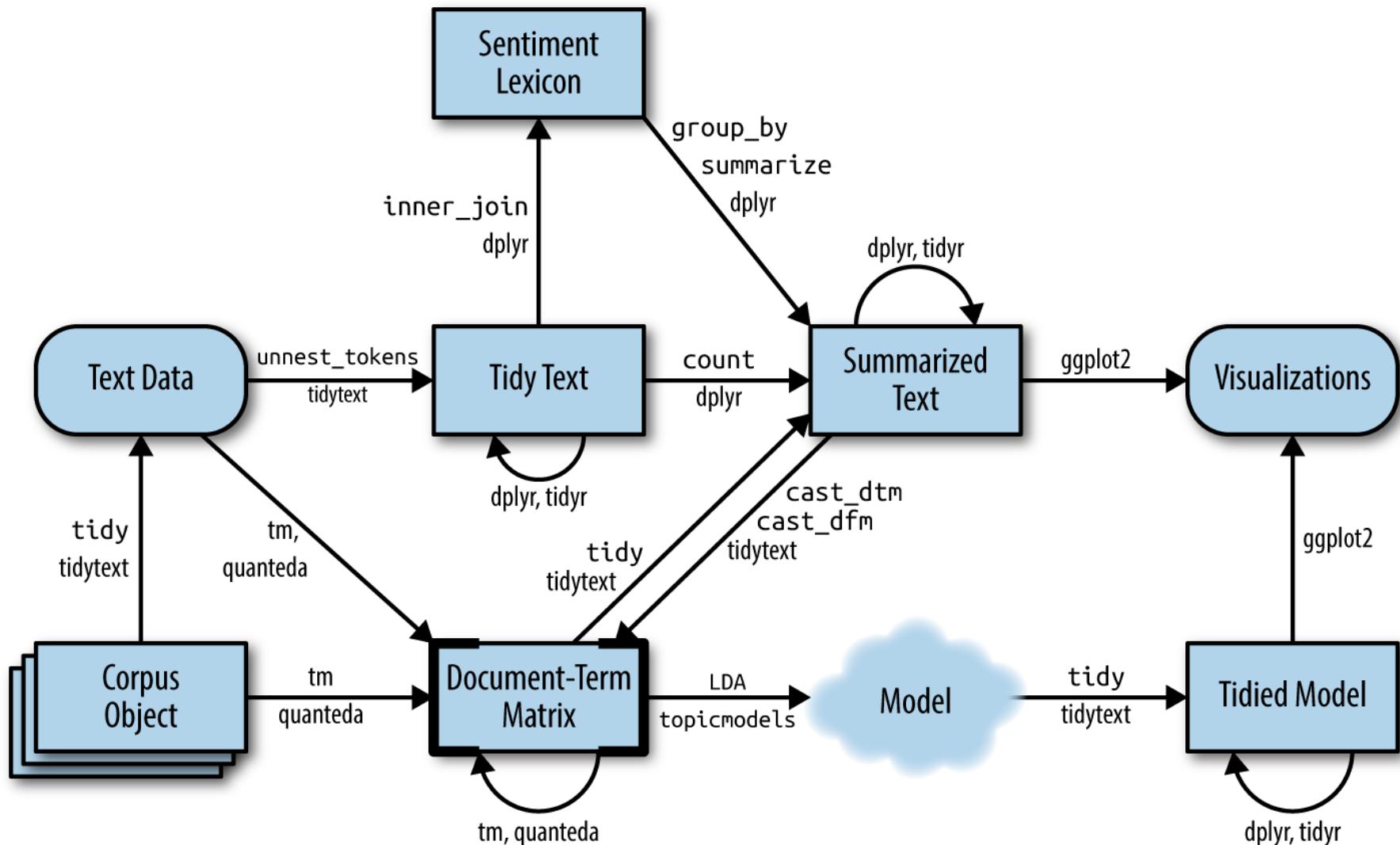


Find us at...

 @dataandme
 @batpigandme
 maraaverick.rbind.io

TIDYING AND CASTING





Two powerful NLP techniques

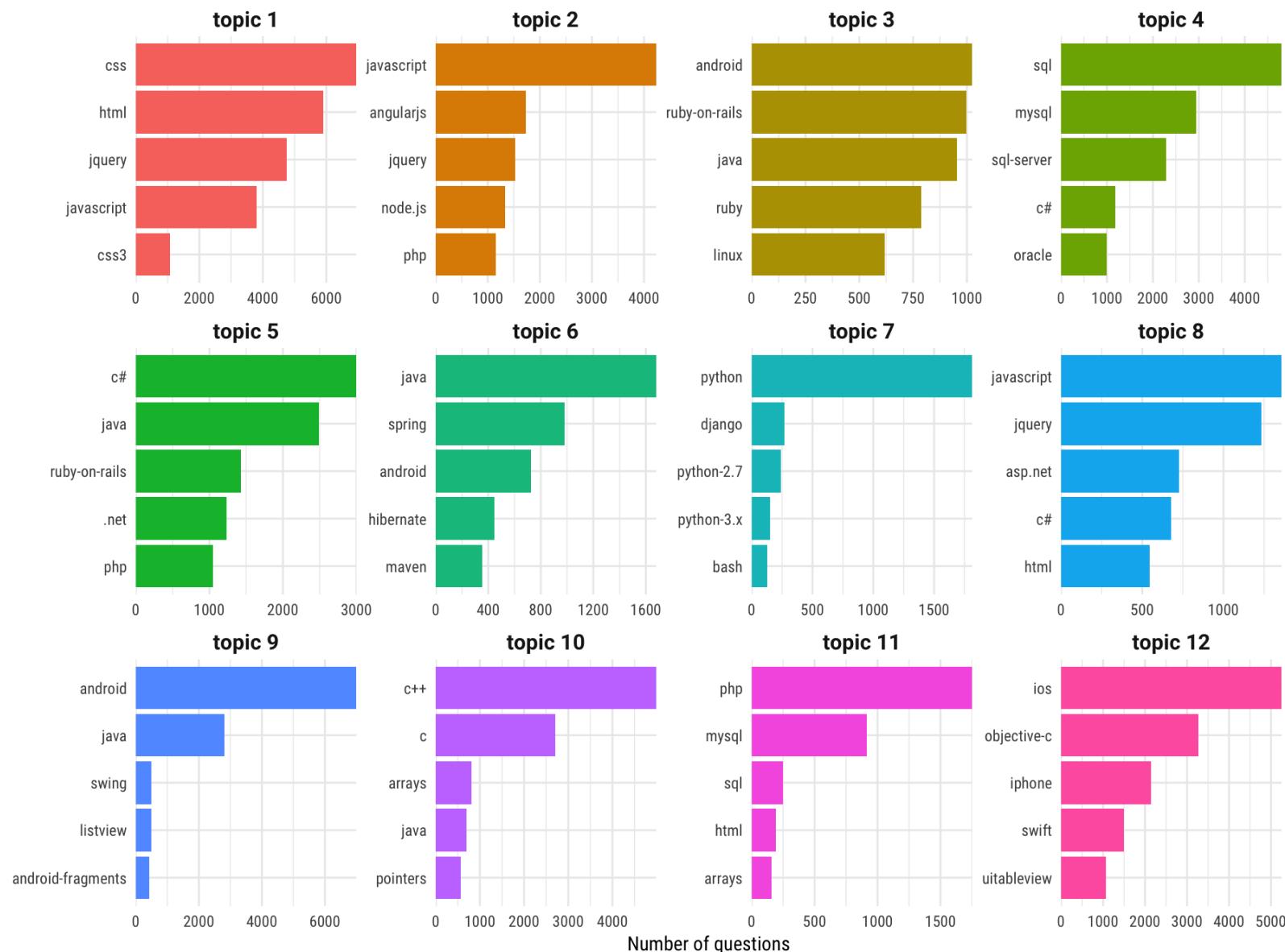
- Topic modeling
- Text classification

Topic modeling

- Each DOCUMENT = mixture of topics
- Each TOPIC = mixture of words

Top tags for each LDA topic

For questions with >80% probability for that topic



GREAT LIBRARY HEIST



Downloading your text data

```
library(tidyverse)
library(gutenbergr)

titles <- c("Twenty Thousand Leagues under the Sea",
           "The War of the Worlds",
           "Pride and Prejudice",
           "Great Expectations")

books <- gutenberg_works(title %in% titles) %>%
  gutenberg_download(meta_fields = "title")

books
```

```
## # A tibble: 51,663 x 3
##       gutenberg_id text
##   <int> <chr>
## 1          36 The War of the Worlds
## 2          36 ""
## 3          36 by H. G. Wells [1898]
## 4          36 ""
## 5          36 ""
## 6          36 "    But who shall dwell in these worlds ...
## 7          36 "    inhabited? . . . Are we or they L...
## 8          36 "    world? . . . And how are all thin...
## 9          36 "        KEPLER (quoted in The Anatomy o...
## 10         36 ""
## # ... with 51,653 more rows
```

		title
		<chr>
1	36	The War of the...
2	36	The War of the...
3	36	The War of the...
4	36	The War of the...
5	36	The War of the...
6	36	The War of the...
7	36	The War of the...
8	36	The War of the...
9	36	The War of the...
10	36	The War of the...

Someone has torn your books apart! 😭

```
by_chapter <- books %>%
  group_by(title) %>%
  mutate(chapter = cumsum(str_detect(text,
                                      regex("^\u03b7chapter ",
                                             ignore_case = TRUE)))) %>%
  ungroup() %>%
  filter(chapter > 0) %>%
  unite(document, title, chapter)

by_chapter

## # A tibble: 51,602 x 3
##   gutenberg_id text
##       <int> <chr>
## 1           36 CHAPTER ONE
## 2           36 ""
## 3           36 THE EVE OF THE WAR
## 4           36 ""
## 5           36 ""
## 6           36 No one would have believed in the last ye...
## 7           36 century that this world was being watched...
## 8           36 intelligences greater than man's and yet ...
## 9           36 men busied themselves about their various...
## 10          36 scrutinised and studied, perhaps almost a...
## # ... with 51,592 more rows
```

Can we put them back together?

```
library(tidytext)

word_counts <- by_chapter %>%
  unnest_tokens(word, text) %>%
  anti_join(get_stopwords(source = "smart")) %>%
  count(document, word, sort = TRUE)

word_counts

## # A tibble: 111,650 x 3
##   document          word     n
##   <chr>            <chr>  <int>
## 1 Great Expectations_57 joe    88
## 2 Great Expectations_7  joe    70
## 3 Pride and Prejudice_18 mr     66
## 4 Great Expectations_17 biddy  63
## 5 Great Expectations_27 joe    58
## 6 Great Expectations_38 estella 58
## 7 Great Expectations_2  joe    56
## 8 Great Expectations_23 pocket  53
## 9 Great Expectations_15 joe    50
## 10 Great Expectations_18 joe   50
## # ... with 111,640 more rows
```

Can we put them back together?

```
words_sparse <- word_counts %>%  
  cast_sparse(document, word, n)
```

```
class(words_sparse)
```

```
## [1] "dgCMatrix"  
## attr(,"package")  
## [1] "Matrix"
```

Train a topic model

Use a sparse matrix or a `quanteda::dfm` object as input

```
library(stm)

topic_model <- stm(words_sparse, K = 4,
                     verbose = FALSE, init.type = "Spectral")

summary(topic_model)

## A topic model with 4 topics, 193 documents and a 18360 word dictionary.

## Topic 1 Top Words:
## Highest Prob: mr, elizabeth, mrs, darcy, bennet, miss, jane
## FREX: elizabeth, darcy, bennet, bingley, wickham, collins, lydia
## Lift: wickham, nephew, phillips, brighton, meryton, bourgh, mend
## Score: elizabeth, darcy, bennet, bingley, wickham, jane, lydia
## Topic 2 Top Words:
## Highest Prob: captain, nautilus, sea, nemo, ned, conseil, land
## FREX: nautilus, nemo, ned, conseil, canadian, ocean, seas
## Lift: vanikoro, indian, d'urville, reefs, scotia, shark's, solidification
## Score: nautilus, nemo, ned, conseil, canadian, ocean, captain
## Topic 3 Top Words:
## Highest Prob: mr, joe, miss, time, pip, looked, herbert
## FREX: joe, pip, herbert, wemmick, havisham, estella, biddy
## Lift: towel, giv, whimple, meantersay, jew, rot, barnard's
## Score: joe, wemmick, pip, jagers, havisham, estella, herbert
## Topic 4 Top Words:
## Highest Prob: people, martians, man, time, black, men, night
## FREX: martians, martian, woking, mars, curate, pine, ulla
## Lift: martians, mars, curate, shepperton, henderson, hood, ripley
## Score: martians, martian, woking, cylinder, curate, ulla, pine
```

Exploring the output of topic modeling

Time for tidying!

```
chapter_topics <- tidy(topic_model, matrix = "beta")  
chapter_topics
```

```
## # A tibble: 73,440 x 3  
##   topic term      beta  
##   <int> <chr>    <dbl>  
## 1     1 joe  8.69e-104  
## 2     2 joe  3.03e-139  
## 3     3 joe  1.21e- 2  
## 4     4 joe  3.28e- 19  
## 5     1 mr   1.90e- 2  
## 6     2 mr   1.91e- 4  
## 7     3 mr   1.22e- 2  
## 8     4 mr   1.15e- 45  
## 9     1 biddy 3.21e- 80  
## 10    2 biddy 3.84e-149  
## # ... with 73,430 more rows
```

Exploring the output of topic modeling

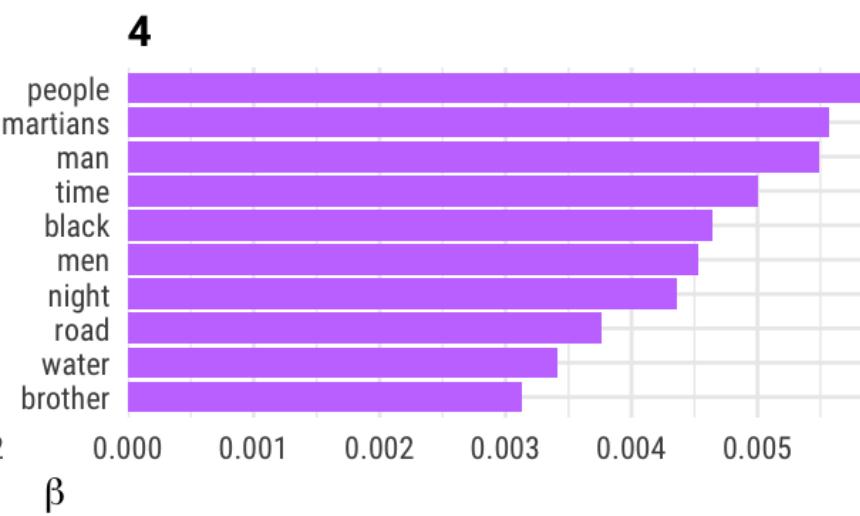
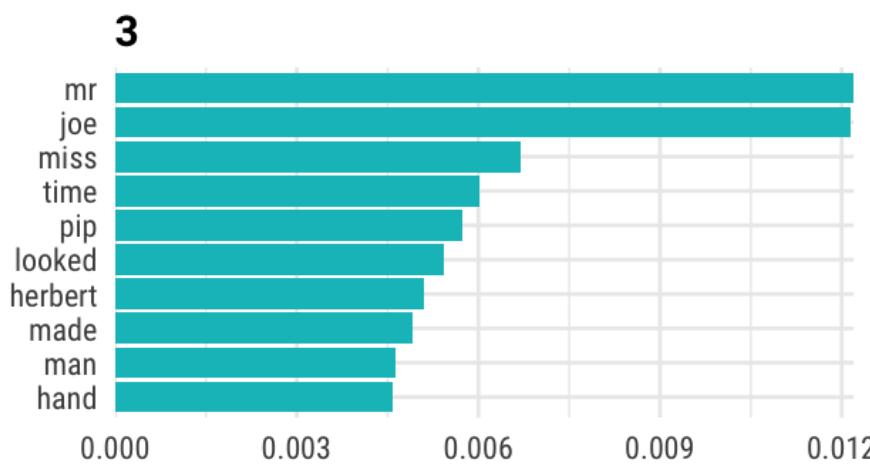
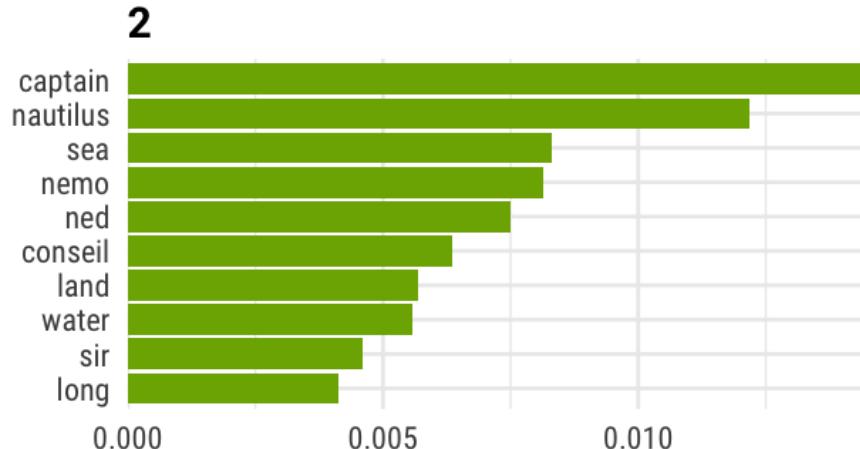
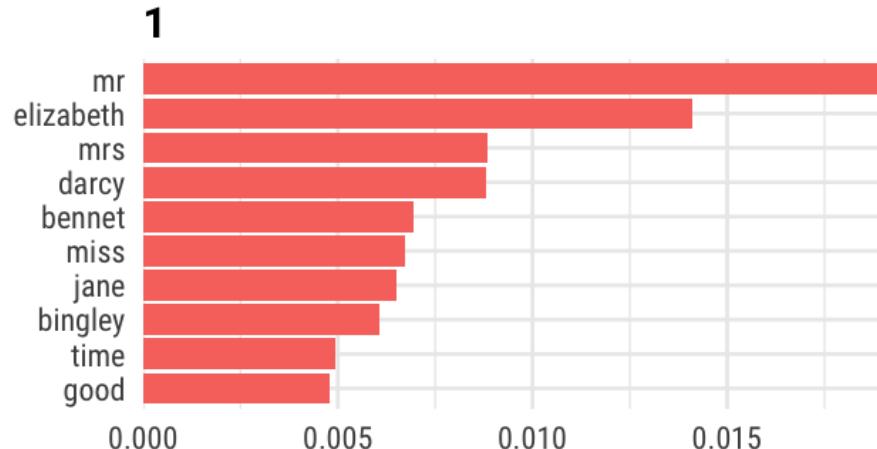
```
top_terms <- chapter_topics %>%
  group_by(topic) %>%
  top_n(10, beta) %>%
  ungroup() %>%
  arrange(topic, -beta)
```

```
top_terms
```

```
## # A tibble: 40 x 3
##   topic term      beta
##   <int> <chr>    <dbl>
## 1     1 mr      0.0190
## 2     1 elizabeth 0.0141
## 3     1 mrs      0.00886
## 4     1 darcy    0.00881
## 5     1 bennet   0.00694
## 6     1 miss     0.00674
## 7     1 jane     0.00652
## 8     1 bingley  0.00607
## 9     1 time     0.00493
## 10    1 good     0.00480
## # ... with 30 more rows
```

Exploring the output of topic modeling

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



β

How are documents classified?

```
chapters_gamma <- tidy(topic_model, matrix = "gamma",
                        document_names = rownames(words_sparse))
```

```
chapters_gamma
```

```
## # A tibble: 772 x 3
##   document          topic     gamma
##   <chr>            <int>    <dbl>
## 1 Great Expectations_57     1  0.000792
## 2 Great Expectations_7      1  0.00340
## 3 Pride and Prejudice_18    1  1.000
## 4 Great Expectations_17    1  0.0480
## 5 Great Expectations_27    1  0.000367
## 6 Great Expectations_38    1  0.00110
## 7 Great Expectations_2      1  0.000531
## 8 Great Expectations_23    1  0.432
## 9 Great Expectations_15    1  0.000565
## 10 Great Expectations_18    1  0.000277
## # ... with 762 more rows
```

How are documents classified?

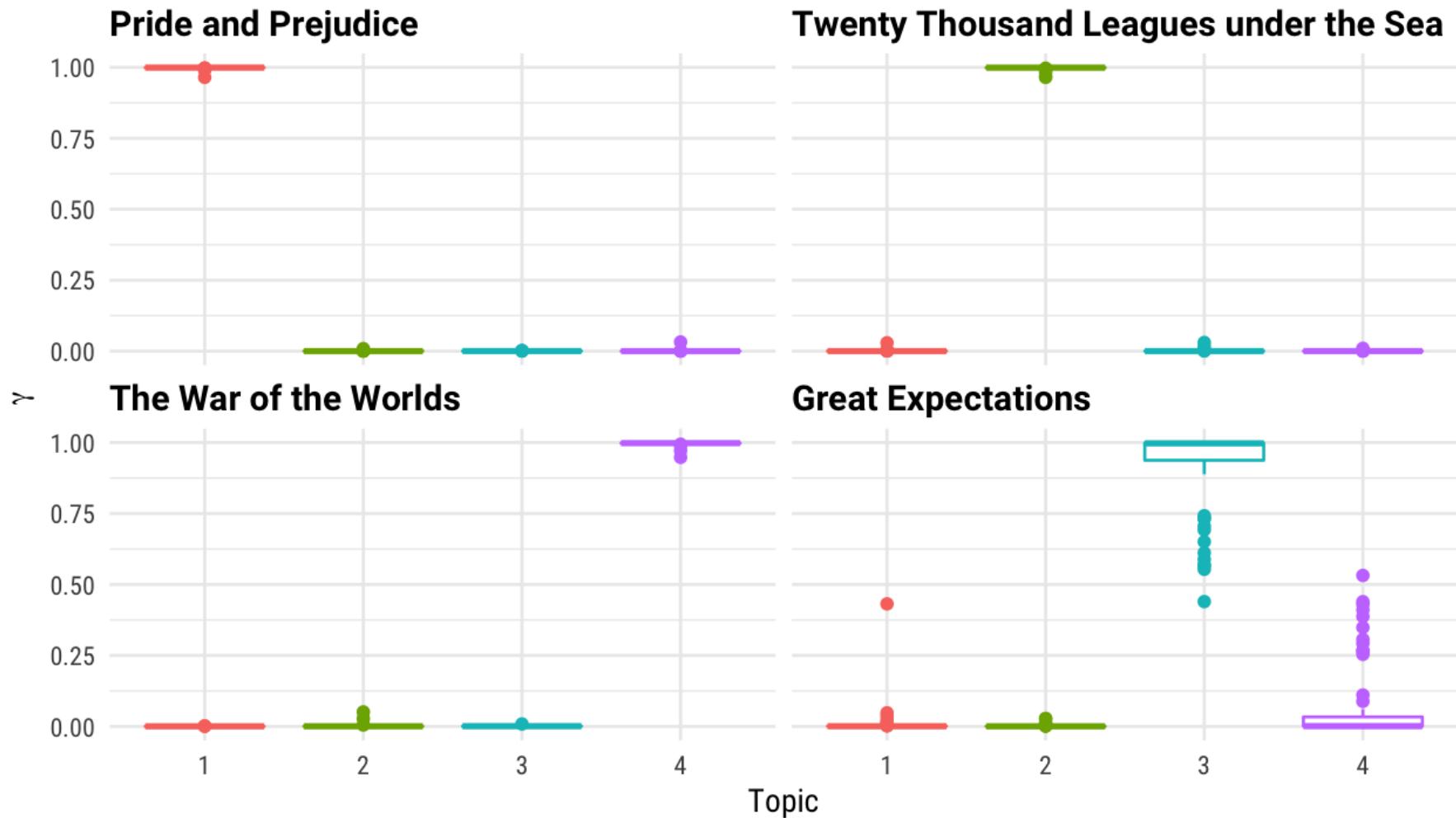
```
chapters_parsed <- chapters_gamma %>%
  separate(document, c("title", "chapter"),
          sep = "_", convert = TRUE)

chapters_parsed

## # A tibble: 772 x 4
##   title           chapter topic    gamma
##   <chr>          <int>  <int>    <dbl>
## 1 Great Expectations     57      1  0.000792
## 2 Great Expectations      7      1  0.00340
## 3 Pride and Prejudice     18      1  1.000
## 4 Great Expectations     17      1  0.0480
## 5 Great Expectations     27      1  0.000367
## 6 Great Expectations     38      1  0.00110
## 7 Great Expectations      2      1  0.000531
## 8 Great Expectations     23      1  0.432
## 9 Great Expectations     15      1  0.000565
## 10 Great Expectations     18      1  0.000277
## # ... with 762 more rows
```

How are documents classified?

```
chapters_parsed %>%
  mutate(title = fct_reorder(title, gamma * topic)) %>%
  ggplot(aes(factor(topic), gamma)) +
  geom_boxplot() +
  facet_wrap(~ title)
```



GOING FARTHER



Tidying model output

Which words in each document are assigned to which topics?

- `augment()`
- Add information to each observation in the original data

The screenshot shows the IBM Data Science Experience interface. On the left, an RStudio-like environment displays R code for topic modeling. The code includes loading tidytext and stm packages, reading a Gutenberg download, and preparing the data for analysis. In the center, a file browser shows a folder named 'topic-modeling' containing 'topic-modeling.Rproj' and 'game_in_a_few.Rmd'. On the right, a console window shows the execution of the R code, including the creation of a tibble and its contents.

Topic modeling with R and tidy data principles

1,971 views

1 like

0 dislikes

SHARE

...



Julia Silge

Published on Dec 18, 2017

SUBSCRIBE 74

Watch along as I demonstrate how to train a topic model in R using the tidytext and stm packages on a collection of Sherlock Holmes stories. In this video, I'm working in IBM Cloud's Data Science Experience environment.

SHOW MORE

Using stm

- Document-level covariates

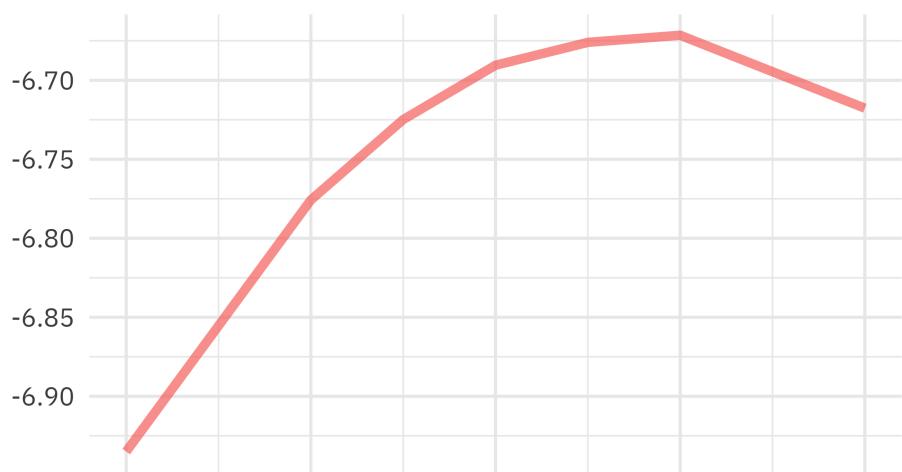
```
topic_model <- stm(words_sparse, K = 0, init.type = "Spectral",
                     prevalence = ~s(Year),
                     data = covariates,
                     verbose = FALSE)
```

- Use functions for `semanticCoherence()`,
`checkResiduals()`, `exclusivity()`, and more!
- Check out <http://www.structuraltopicmodel.com/>
- See my blog post for how to choose **K**, the number of topics

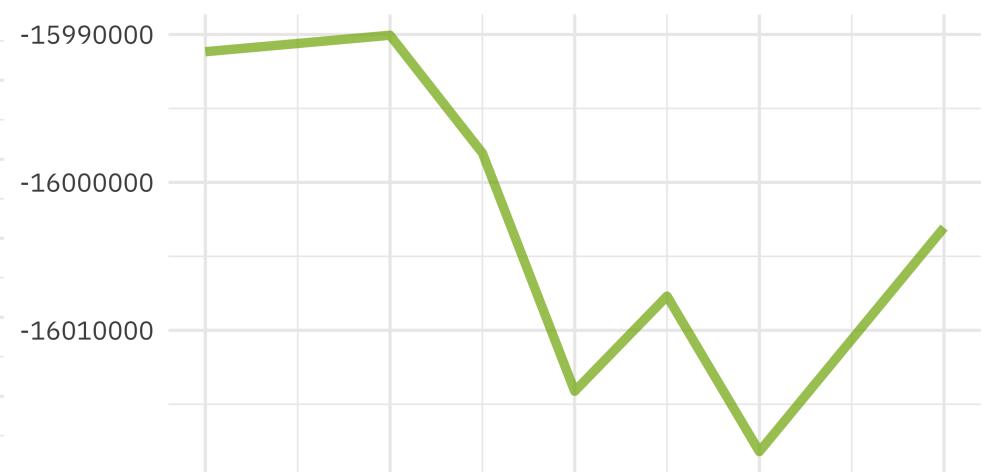
Model diagnostics by number of topics

These diagnostics indicate that a good number of topics would be around 60

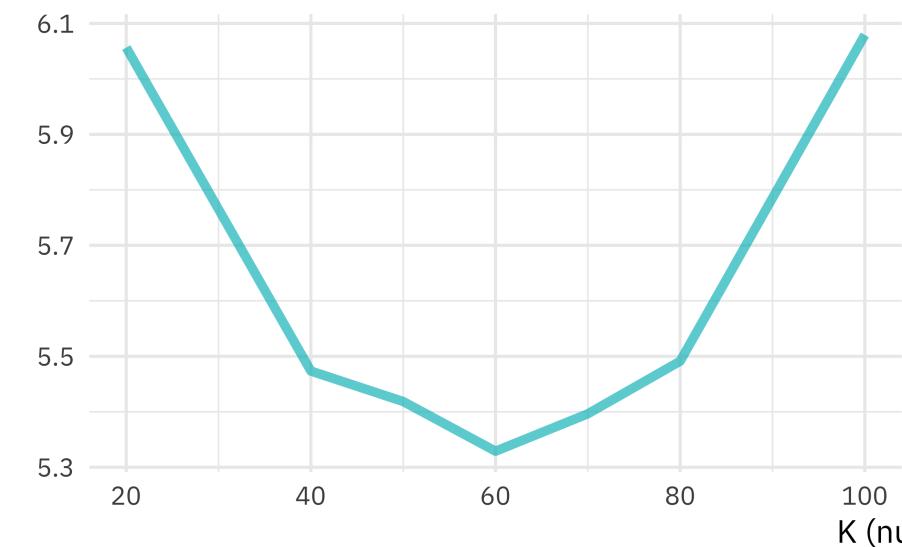
Held-out likelihood



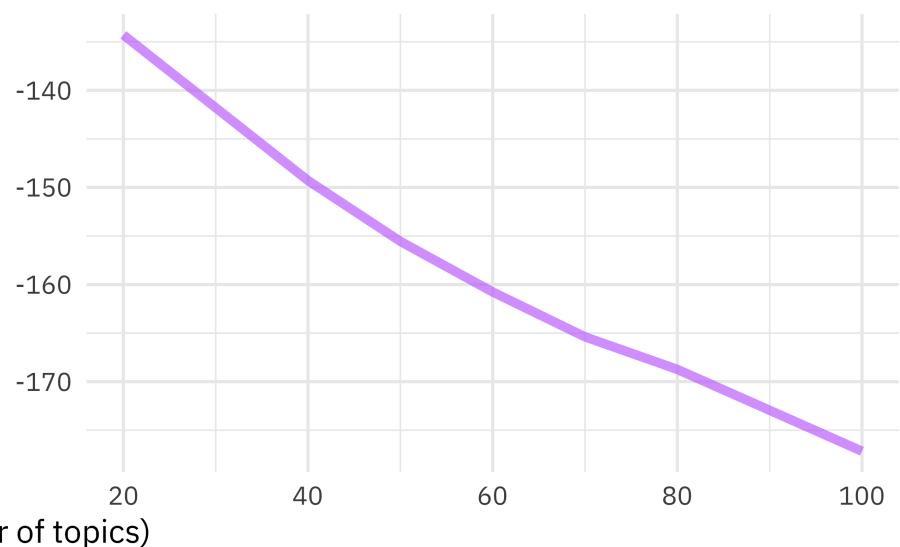
Lower bound



Residuals



Semantic coherence



Stemming?

Advice from Schofield & Mimno

"Comparing Apples to Apple: The Effects of Stemmers on Topic Models"

“

**Despite their frequent use in topic modeling,
we find that stemmers produce no
meaningful improvement in likelihood and
coherence and in fact can degrade topic
stability.**

”

TEXT CLASSIFICATION



Downloading your text data

```
library(tidyverse)
library(gutenbergr)

titles <- c("The war of the worlds",
           "Pride and Prejudice")

books <- gutenberg_works(title %in% titles) %>%
  gutenberg_download(meta_fields = "title") %>%
  mutate(document = row_number())

books
```

```
## # A tibble: 19,504 x 4
##   gutenberg_id text
##   <int> <chr>
## 1 36  The war of the worlds
## 2 36  ""
## 3 36  by H. G. Wells [1898]
## 4 36  ""
## 5 36  ""
## 6 36  "    But who shall dwell in these ...
## 7 36  "    inhabited? . . . Are we or...
## 8 36  "    world? . . . And how are a...
## 9 36  "        KEPLER (quoted in The An...
##10 36  ""

## # ... with 19,494 more rows
```

	title	document
	<chr>	<int>
1	The War of t...	1
2	The War of t...	2
3	The War of t...	3
4	The War of t...	4
5	The War of t...	5
6	The War of t...	6
7	The War of t...	7
8	The War of t...	8
9	The War of t...	9
10	The War of t...	10

Making a tidy dataset

Use this kind of data structure for EDA! 

```
library(tidytext)
```

```
tidy_books <- books %>%
  unnest_tokens(word, text) %>%
  group_by(word) %>%
  filter(n() > 10) %>%
  ungroup
```

tidy_books

Cast to a sparse matrix

And build a dataframe with a response variable

```
sparse_words <- tidy_books %>%
  count(document, word, sort = TRUE) %>%
  cast_sparse(document, word, n)

books_joined <- tibble(document = as.integer(rownames(sparse_words))) %>%
  left_join(books %>%
    select(document, title))
```

Train a glmnet model

```
library(glmnet)
library(doMC)
registerDoMC(cores = 8)

is_jane <- books_joined$title == "Pride and Prejudice"

model <- cv.glmnet(sparse_words, is_jane, family = "binomial",
                    parallel = TRUE, keep = TRUE)
```

Tidying our model

Tidy, then filter to choose some lambda from glmnet output

```
library(broom)

coefs <- model$glmnet.fit %>%
  tidy() %>%
  filter(lambda == model$lambda.1se)

Intercept <- coefs %>%
  filter(term == "(Intercept)") %>%
  pull(estimate)
```

Tidying our model

```
classifications <- tidy_books %>%
  inner_join(coefs, by = c("word" = "term")) %>%
  group_by(document) %>%
  summarize(score = sum(estimate)) %>%
  mutate(probability = plogis(Intercept + score))

classifications

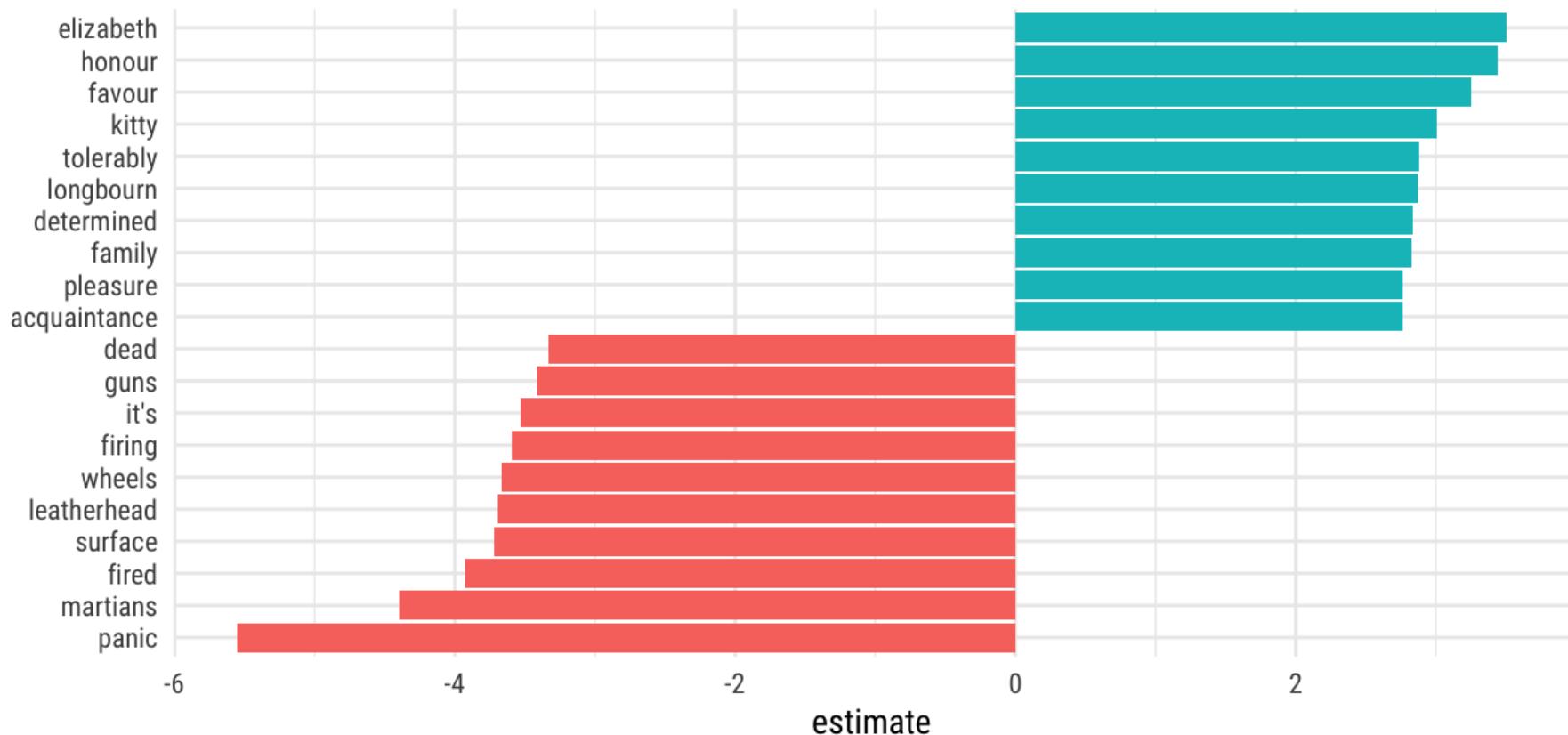
## # A tibble: 16,001 x 3
##   document  score probability
##       <int>   <dbl>      <dbl>
## 1         1 -2.34      0.110
## 2         3  0.205     0.611
## 3         6   1.85     0.890
## 4         7  -1.02     0.315
## 5         8  -1.25     0.268
## 6         9  -0.526    0.430
## 7        13  -0.238    0.502
## 8        15  -5.47    0.00533
## 9        19   0.373    0.650
## 10       21  -2.34      0.110
## # ... with 15,991 more rows
```

Understanding our model

```
coefs %>%
  group_by(estimate > 0) %>%
  top_n(10, abs(estimate)) %>%
  ungroup %>%
  ggplot(aes(fct_reorder(term, estimate),
             estimate,
             fill = estimate > 0)) +
  geom_col(show.legend = FALSE) +
  coord_flip()
```

Coefficients that increase/decrease probability

A document mentioning Martians is unlikely to be written by Jane Austen



ROC

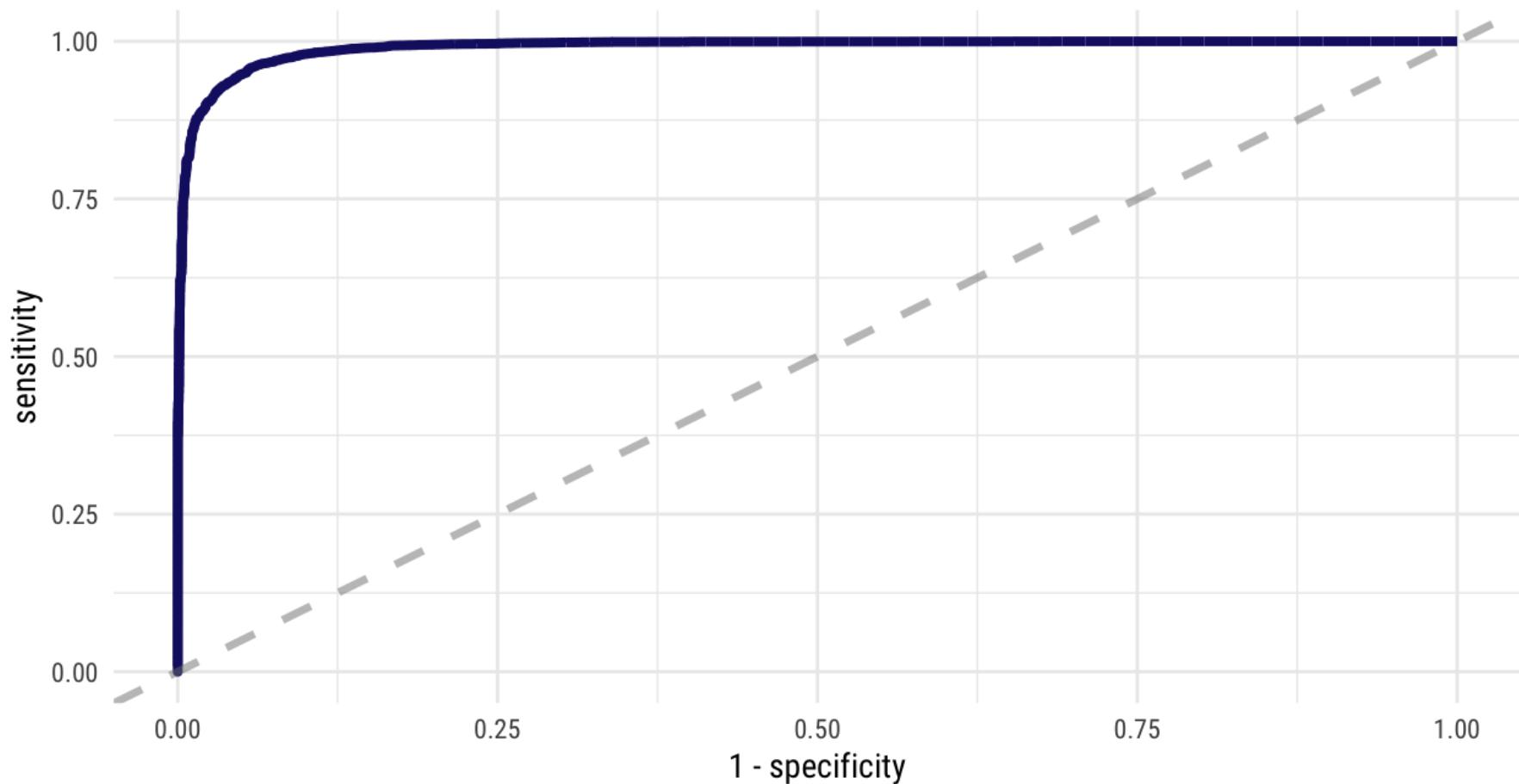
```
library(yardstick)

comment_classes <- classifications %>%
  left_join(books %>%
    select(title, document), by = "document") %>%
  mutate(title = as.factor(title))
```

ROC

```
comment_classes %>%
  roc_curve(title, probability) %>%
  ggplot(aes(x = 1 - specificity, y = sensitivity)) +
  geom_line(
    color = "midnightblue",
    size = 1.5
  ) +
  geom_abline(
    lty = 2, alpha = 0.5,
    color = "gray50",
    size = 1.2
  )
```

ROC curve for text classification



AUC for model

```
comment_classes %>%  
  roc_auc(title, probability)  
  
## # A tibble: 1 x 3  
##   .metric .estimator .estimate  
##   <chr>   <chr>       <dbl>  
## 1 roc_auc binary     0.990
```

Confusion matrix

```
comment_classes %>%
  mutate(
    prediction = case_when(
      probability > 0.5 ~ "Pride and Prejudice",
      TRUE ~ "The War of the Worlds"
    ),
    prediction = as.factor(prediction)
  ) %>%
  conf_mat(title, prediction)
```

```
##                                     Truth
## Prediction                           Pride and Prejudice  The war of the Worlds
##   Pride and Prejudice                10351                  484
##   The War of the Worlds                 264                  4902
```

Misclassifications

Let's talk about misclassifications. Which documents here were incorrectly predicted to be written by Jane Austen?

```
comment_classes %>%
  filter(
    probability > .8,
    title == "The War of the Worlds"
  ) %>%
  sample_n(10) %>%
  inner_join(books %>%
    select(document, text)) %>%
  select(probability, text)

## # A tibble: 10 x 2
##   probability text
##       <dbl> <chr>
## 1      0.858 ladies there being by no means the least active.
## 2      0.851 is wrong as well as any, but not what is possible to tortur...
## 3      0.972 She put her hand to her throat--swayed. I made a step forw...
## 4      0.962 "\"Be a man!\" said I. \"You are scared out of your wits! ...
## 5      0.832 "\"Take this!\" said the slender lady, and she gave my brot...
## 6      0.827 decorum were necessarily different from ours; and not only ...
## 7      0.906 "\"Half a mile, you say?\" said he."
## 8      0.910 breed. I tell you, I'm grim set on living. And if I'm not...
## 9      0.854 would be advisable to kill him, lest his actions attracted ...
## 10     0.919 winter. Its air is much more attenuated than ours, its oce...
```

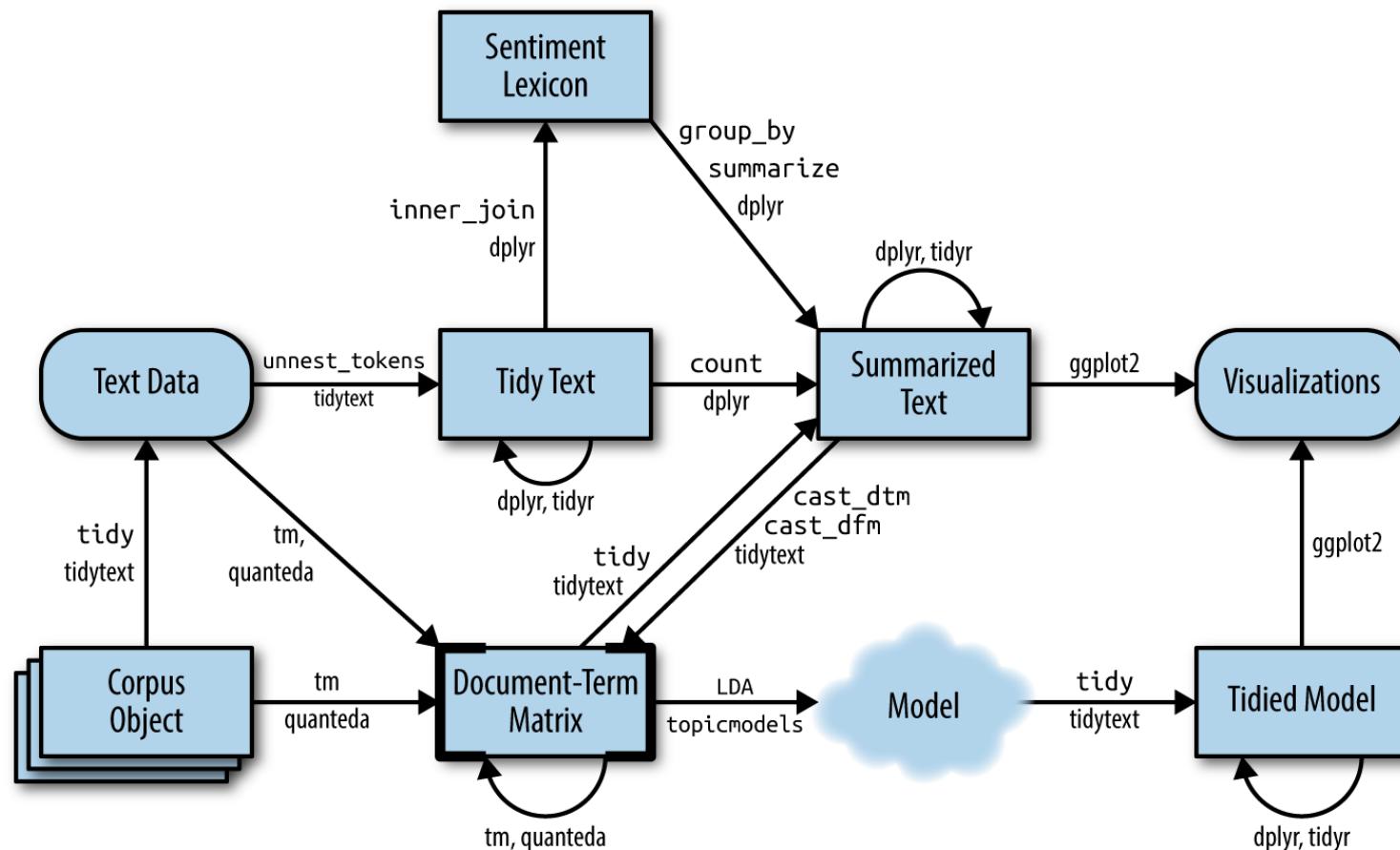
Misclassifications

Let's talk about misclassifications. Which documents here were incorrectly predicted to *not* be written by Jane Austen?

```
comment_classes %>%
  filter(
    probability < .3,
    title == "Pride and Prejudice"
  ) %>%
  sample_n(10) %>%
  inner_join(books %>%
    select(document, text)) %>%
  select(probability, text)

## # A tibble: 10 x 2
##   probability text
##   <dbl> <chr>
## 1 0.188 is so violent, that it would be the death of half the good ...
## 2 0.208 blush. He absolutely started, and for a moment seemed immov...
## 3 0.269 of contradictions and varieties, sighed at the perverseness...
## 4 0.220 suddenly came forward from the road, which led behind it to...
## 5 0.199 "\"A little sea-bathing would set me up forever.\""
## 6 0.286 it had just transpired that he had left gaming debts behind...
## 7 0.266 of the gates into the ground.
## 8 0.120 the happiest of men.
## 9 0.218 They travelled as expeditiously as possible, and, sleeping ...
## 10 0.279 the improvements it was receiving, he was happily employed ...
```

Workflow for text mining/modeling



Go explore real-world text!





Thanks!

-  tidytextmining.com
-  [@juliasilge](https://twitter.com/juliasilge)
-  [@juliasilge](https://juliasilge.com)
-  juliasilge.com
-  [@dataandme](https://twitter.com/@dataandme)
-  [@batpigandme](https://batpigandme.com)
-  maraaverick.rbind.io

Slides created with **remark.js** and the R package **xaringan**