



APSTA-GE 2352

Statistical Computing: Lecture 10

Klint Kanopka

New York University

NYUGreyArtGallery

SEYER CENTER

WASHINGTON St



Table of Contents

1. Statistical Computing - Week 10
 1. Table of Contents
 2. Announcements
2. A New Problem
3. Matrix Factorization
4. More Similarity Metrics
5. The Quantification of Text
6. Classical Text Analysis
7. Natural Language Processing with Deep Learning
 1. Word vectors
 2. Transformer Models
8. Language Models
9. Latent Semantic Analysis
10. Wrap Up

Check-In

- Pollev.com/klintkanopka

Announcements

- PS5 due in a week
- Upcoming PRIISM seminar talk with my friend Charlie Rahal from Oxford
 - Associate Professor in Data Science and Informatics @ University of Oxford
 - Wednesday, Nov 19 @11a
 - Talk is about algorithmic fairness and the potential consequences of using `set.seed()` incorrectly in your research
 - Check Slack for more info!
- Tomorrow the SCSS PhD students are giving talks
 - 10a in Kimball 301W
 - Should be super interesting, may give you ideas or places you can get involved in research!

A New Problem

Problem

- Today we want to try to understand why people hate movies
- To do this, we'll have a collection of the text of 12,500 **negative** reviews from IMDB
- Our goal is to try to understand what people hate about movies

Matrix Factorization

Matrix Factorization

- Less actual linear algebra than when we talked about PCA!
- Let's first think about factoring: Where have you used it in the past?
- Say we have a $m \times n$ real matrix, \mathbf{M}
- We can factor it into: $\mathbf{M} = \mathbf{U}\Sigma\mathbf{V}^\top$
 - \mathbf{U} is a $m \times m$ real orthogonal matrix and the columns are called the *left singular vectors*
 - \mathbf{V} is a $n \times n$ real orthogonal matrix and the columns are called the *right singular vectors*
 - Σ is a $m \times n$ positive diagonal rectangular matrix (aka: $\sigma_i = \Sigma_{ii}$ and $\Sigma_{i \neq j} = 0$), and the diagonal entries (called *singular values*) are arranged in decreasing order
 - We can write the whole original matrix as $\mathbf{M} = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^\top$
- This is the Singular Value Decomposition (SVD). Sometimes SVD refers to the *compact SVD*, where Σ is an $r \times r$ matrix that only contains the r non-zero singular values, and \mathbf{U} is $m \times r$ and \mathbf{V} is $n \times r$
 - The compact SVD just throws away a bunch of zeros

Singular Value Decomposition

- Does the SVD remind you of anything?
- As it turns out, the SVD of a data matrix is how PCA is usually computed
 - This is because of a relationship between the SVD and the eigenvalues of $\mathbf{M}^\top \mathbf{M}$
- Note that the SVD recovers our original matrix *exactly*! It is not an approximation!
- But, the SVD requires a lot of computational time for large matrices and sometimes we only want the first few (aka most important) singular vectors
- We introduce now the *truncated SVD*
 - The truncated SVD is the same setup as the compact SVD, with one modification:
 - We only keep the left and right singular vectors associated with the first r singular values!
 - This means the truncated SVD **is** an approximation if $r < \text{rank}(\mathbf{M})$

Truncated SVD

- The truncated SVD is *much* faster to compute than the full SVD
- Done in a method similar to power iteration
- The truncated SVD is written as: $\tilde{\mathbf{M}} = \mathbf{U}\Sigma\mathbf{V}^\top$
- The truncated SVD minimizes the *Frobenius Norm*: $||\mathbf{M} - \tilde{\mathbf{M}}||_F$
- Guess what the Frobenius Norm is

$$||A||_F = \sqrt{\sum_i^m \sum_j^n a_{ij}^2}$$

- This means that the result of the truncated SVD, $\tilde{\mathbf{M}}$, is the best rank r approximation of the original matrix, \mathbf{M}

More Similarity Metrics

- How do you know that two things are similar?
- In some ways, "similarity" is the inverse of "distance"
- You can use distance formulas (Euclidean, Manhattan)
- Minkowsky Distance generalizes the two:

$$D(A, B) = \left(\sum_{i=1}^n |a_i - b_i|^p \right)^{\frac{1}{p}}$$

More Similarity Metrics

- For vectors, we frequently use *cosine similarity*
 - Measures the angle between two vectors
 - Cares about the *direction* they point in, but not the *length*

$$S_c(A, B) = \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\| \|\mathbf{B}\|} = \frac{\sum_{i=1}^n a_i b_i}{\sqrt{\sum_{i=1}^n a_i^2} \cdot \sqrt{\sum_{i=1}^n b_i^2}}$$

Similarity between other things

- For sets, we often use Jaccard Similarity
 - The number of items that are common to both sets divided by the number of items across both sets
 - Very common with binary (or dichotomous) data!

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|} = \frac{|A \cap B|}{|A| + |B| - |A \cap B|}$$

Similarity between other strings

- For strings of text, there are a ton of options:
 - Hamming Distance counts the number of positions where the two strings are different
 - `hamburger` and `burgers` have a Hamming distance of 8
 - `here` and `there` have a Hamming distance of 5
 - Levenshtein distance measures the minimum number of single-character edits required to change one into the other
 - Counts insertions, deletions, and substitutions
 - `hamburger` and `burgers` have a Levenshtein distance of 4
 - `here` and `there` have a Levenshtein distance of 1

The Quantification of Text

Quick text vocabulary

- A *token* is an individual unit of text (usually a word)
- A *document* is a collection of tokens
- A *corpus* is a collection of documents
- A *vocabulary* is the collection of tokens represented in a corpus

Quantification of text

- All of the statistics and models we use to describe things act on *numbers*
- How might we turn text into numbers?
- The options boil down to basically two things:
 1. Count words (classical NLP)
 2. Embed words in a vector space (deep NLP)
- From here, we aggregate up to the document level in some way
- The most common classical way to represent a corpus for analysis is the *Document Term Matrix* (DTM)
 - One row per document in a corpus
 - One column per word in the vocab
 - Word counts in the individual cells
 - Treats documents as a *bag of words*

Preprocessing

- Raw text is kind of a mess!
- At the most basic level, we want to strip punctuation and convert words to lowercase
 - What other problems can arise from just counting words?
- **Stop words** are super common words
 - e.g., a, an, the
 - Typically non-informative about the actual content of text, so we remove them
- Sometimes the "same" word can appear in a different form
 - Singular and plural words (dog and dogs)
 - Adjectives and adverbs (quick, quicker, quickly)
 - **Stemming** is the process of stripping tokens down to stems so that these differences are counted together
- More important in classical methods
- Not needed for deep methods
- Often you'll have domain-specific words you also want to drop

Classical Text Analysis

Classical text analysis

- Any thoughts on what we could do with a DTM to learn about the collection of documents?
- We can look at the most common tokens in the corpus
- We can look at the most common tokens in a document, or *term frequency*

$$\text{tf}(t, d) = \frac{f_{t,d}}{\sum_{t' \in d} f_{t',d}}$$

- We can look at how many documents a word appears in, or *inverse document frequency*

$$\text{idf}(t, D) = \ln \frac{N}{d \in D : t \in D}$$

- A common way to weight is tf-idf:

$$\text{tf-idf}(t, d) = \text{tf}(t, d) \cdot \text{idf}(t, D)$$

Dictionary methods

- Dictionaries are collections of words that have been coded to have some particular meaning
 - Conceptually: specific words are associated with ideas, themes, or emotions
 - If you use more of those words, you're expressing more of that thing
- The most common collection of dictionaries is Linguistic Inquiry and Word Count (LIWC)
 - Has dictionaries for affect, emotion, social behavior, references to health and culture, and a ton of other stuff
 - Also counts mechanical properties of text construction
 - These are a *great* place to start with your own data!
 - LIWC costs \$129.95 for an academic three year license, but free alternatives exist
- You can also develop your own dictionaries for specific purposes
 - This can be really time consuming to do well
 - Take a look at *The Development and Psychometric Properties of LIWC-22* technical report for more information on how to do this

Topic Models

- Core idea: Topic models are basically principal component analysis performed on the DTM
- A singular value decomposition (SVD) on the DTM is called Latent Semantic Analysis (LSA)
 - This is functionally the same as doing a PCA on the DTM
 - Interpretation and use is identical - each "topic" is a dimension and words are more or less associated with each
- A better topic model is the Latent Dirichlet Allocation (LDA)
 - A Bayesian model (if you don't know what this means it doesn't super matter)
 - Structured as such:
 - A *document* is made up of a mixture of *topics*
 - A *topic* is made up of a mixture of *words*
 - Here, *mixture* means probability distribution

Topic Models

- Dealing with topic models is conceptually similar to dealing with PCA
- The interpretation you did in lab yesterday is the move!
- There are lots of types of topic models
 - Correlated topic models
 - Structural topic models
 - Allow for topic mixtures to vary as a function of other covariates
 - Use the excellent `stm` package for these

Natural Language Processing with Deep Learning

Word vectors

- The more modern approach is to not preprocess text, but instead *embed* it in a high dimensional vector space
- This means instead of word counts, documents are a sequence of high dimensional (typically $N \approx 300$) vectors
- Word2Vec was the classic implementation
 - "You shall know a word by the company it keeps"
 - Words that co-occur are pushed closer together in the vector space
 - Words that don't co-occur are spread farther apart
 - Typically trained on an enormous corpus (Google News, Google Books, etc)
- Word vectors capture latent features of language usage
 - See Garg, et al. (2018) for the most famous word embedding paper
- Word vectors can also be aggregated in different ways to get up to document-level embeddings

Transformer Models

- Transformers are a specific deep learning architecture developed by Google scientists
 - Vaswani, A., et al. (2017). Attention is all you need. *Advances in neural information processing systems*, 30.
- Bidirectional Encoder Representations from Transformers (BERT) provide contextual embeddings that allow words to have different representations based on context
- Seq2seq models map sequences of embeddings to other sequences of embeddings and drove most machine translation
- Large Language Models (LLMs) like ChatGPT are wildly over-parameterized transformer models that work on vector embeddings

Language Models

What is a language model?

- In general: a model that describes how text is *generated*
- Models should be *generative*, meaning that they can produce new text
- Models should also be *evaluative*, meaning they can quantify the degree to which something looks like "real" text
- This is, primarily (but not always), done by next-token-prediction
- Tons of ways to do this some work better than others
 - Transformer based models generally work best
 - Recurrent neural network models (RNNs) were the old state of the art
 - Before that, models were purely statistical and based on word and/or character frequency
- Part 3 of ps5 builds a character based language model—let's talk about what that is and how it works

A character-based language model

- **Core Idea:** Some character combinations are common (i.e., `th` , `on` , `er`) and some are not (i.e., `qz` , `xx` , `uu`), so we have some idea of what English words *should* look like
- If we train a model from some corpus by counting the instances of two letter pairs, we can compare two pieces new text and see which set of character combinations looks more likely to exist given our trained model

Building the language model

- Beginning from our training text:

1. Preprocess the text down to just individual letter information
2. Count pairs of characters: `teeth` gets a count for `te` , `ee` , `et` , and `th`
3. Put all of these counts into a table where you can retrieve them easily
4. Turn all of these counts into log probabilities

Using the language model

1. Take a string of text and separate it into character pairs:
 - onions becomes on , ni , io , on , ns
2. Look up each character pair in your log probability table
3. Add the log probabilities of each character pair together to get the log probability of the entire string
4. Real words should contain observed pairs of characters, and therefore be considered more likely (higher probability)

Latent Semantic Analysis

- How could we try to find themes in text using computational methods?
- **Setup:** I have a bunch of negative movie reviews, and I want to extract similar themes to understand why people hate movies
- First we need a plan to quantify the text
- Then what can we do to the text?
- **Latent Semantic Analysis** is the SVD of the document-term matrix
- Think of this as doing "PCA on word counts"

Exploring our text

```
1  imdb_neg_raw <- read_csv('data/neg-imdb.csv') |>
2    select(review = FileName, text = Content)
3
4  imdb_neg_raw$text[2]
```

An example review:

```
1  Well...tremors I, the original started off in 1990 and i found the movie quite
2  enjoyable to watch. however, they proceeded to make tremors II and III. Trust
3  me, those movies started going downhill right after they finished the first one,
4  i mean, ass blasters??? Now, only God himself is capable of answering the
5  question \"why in Gods name would they create another one of these dumpster
6  dives of a movie?\" Tremors IV cannot be considered a bad movie, in fact it
7  cannot be even considered an epitome of a bad movie, for it lives up to more
8  than that. As i attempted to sit through it, i noticed that my eyes started to
9  bleed, and i hoped profusely that the little girl from the ring would crawl
10 through the TV and kill me. did they really think that dressing the people who
11 had stared in the other movies up as though they we're from the wild west would
12 make the movie (with the exact same occurrences) any better? honestly, i would
13 never suggest buying this movie, i mean, there are cheaper ways to find things
14 that burn well.
```

Transforming the data into word counts

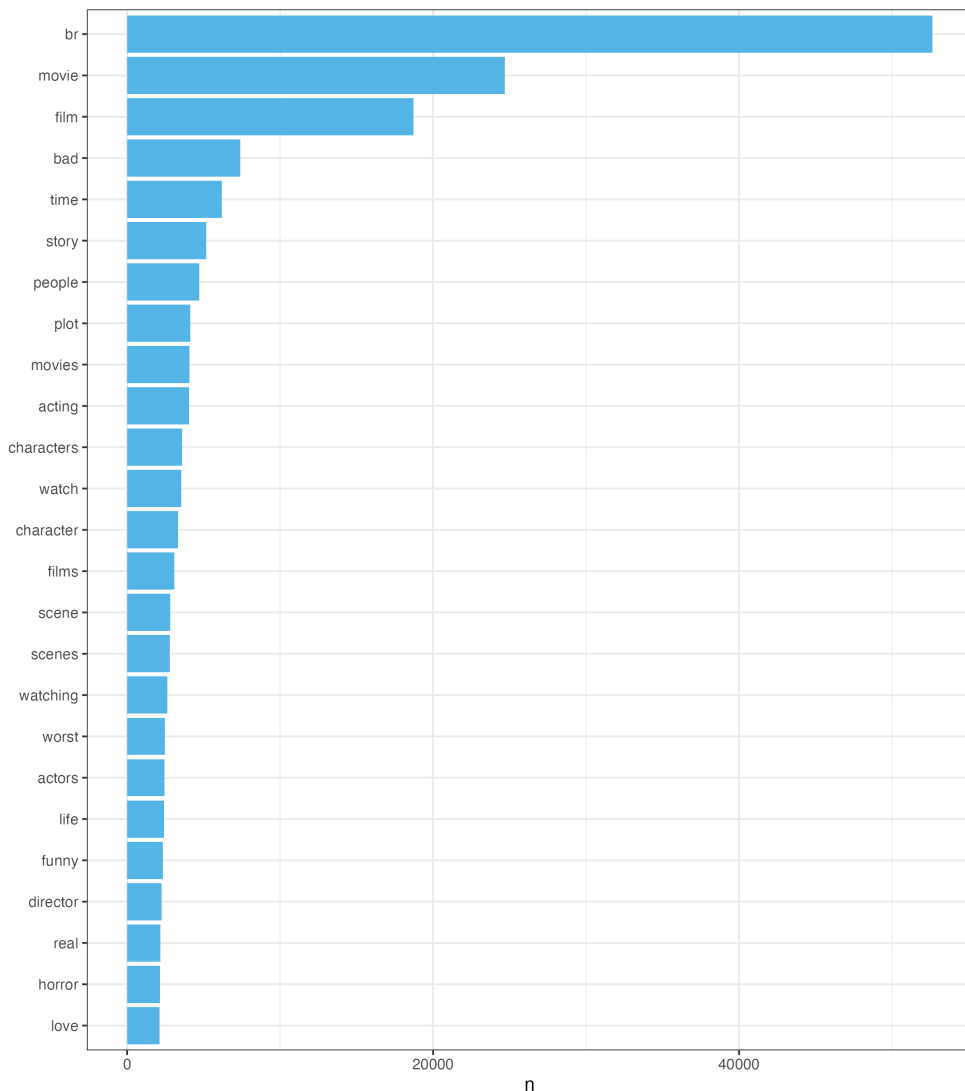
```
1 imdb_neg <- imdb_neg_raw |>
2   unnest_tokens(word, text)
3
4 imdb_neg
5
6 #   review      word
7 #   <chr>      <chr>
8 # 1 1821_4.txt working
9 # 2 1821_4.txt with
10 # 3 1821_4.txt one
11 # 4 1821_4.txt of
12 # 5 1821_4.txt the
13
14 count(imdb_neg, word, sort = TRUE)
15
16 #   word      n
17 #   <chr>  <int>
18 # 1 the    163174
19 # 2 a      79225
20 # 3 and    74347
21 # 4 of     68999
22 # 5 to     68969
```

Removing stop words

```
1 data(stop_words)
2 stop_words
3
4 imdb_neg <- imdb_neg |>
5   anti_join(stop_words)
6
7 imdb_neg
8 count(imdb_neg, word, sort = TRUE)
9
10 #      word      n
11 #   <chr>  <int>
12 # 1 br      52636
13 # 2 movie   24698
14 # 3 film    18717
15 # 4 bad      7389
16 # 5 time     6200
```

Visualizing most common words

```
1  imdb_neg |>
2    count(word, sort = TRUE) |>
3    head(25) |>
4    ggplot(aes(x = n, y = reorder(word, n))) +
5    geom_col(fill = okabeito_colors(2)) +
6    labs(y = NULL) +
7    theme_bw()
```

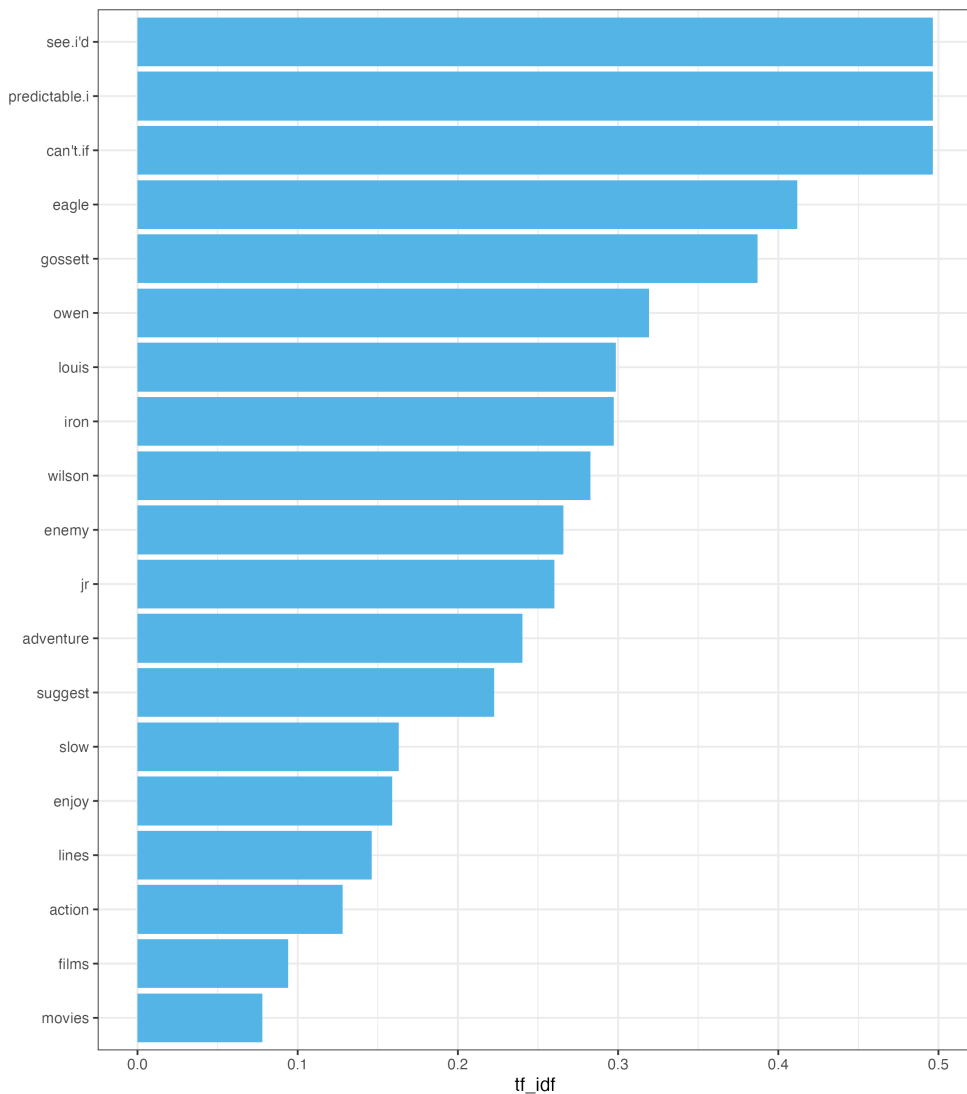


Count by reviews and drop uninformative words

```
1 imdb_counts <- imdb_neg |>
2   group_by(review, word) |>
3   summarize(n = n(), .groups = 'drop')
4
5 imdb_counts <- imdb_counts |>
6   filter(!word %in% c('br', 'movie', 'film'))
7
8 imdb_counts
9
10 #   review  word      n
11 #   <chr>  <chr>    <int>
12 # 1 0_3.txt absurd     2
13 # 2 0_3.txt audience   1
14 # 3 0_3.txt briefly    1
15 # 4 0_3.txt chantings  1
16 # 5 0_3.txt cinematography 1
```

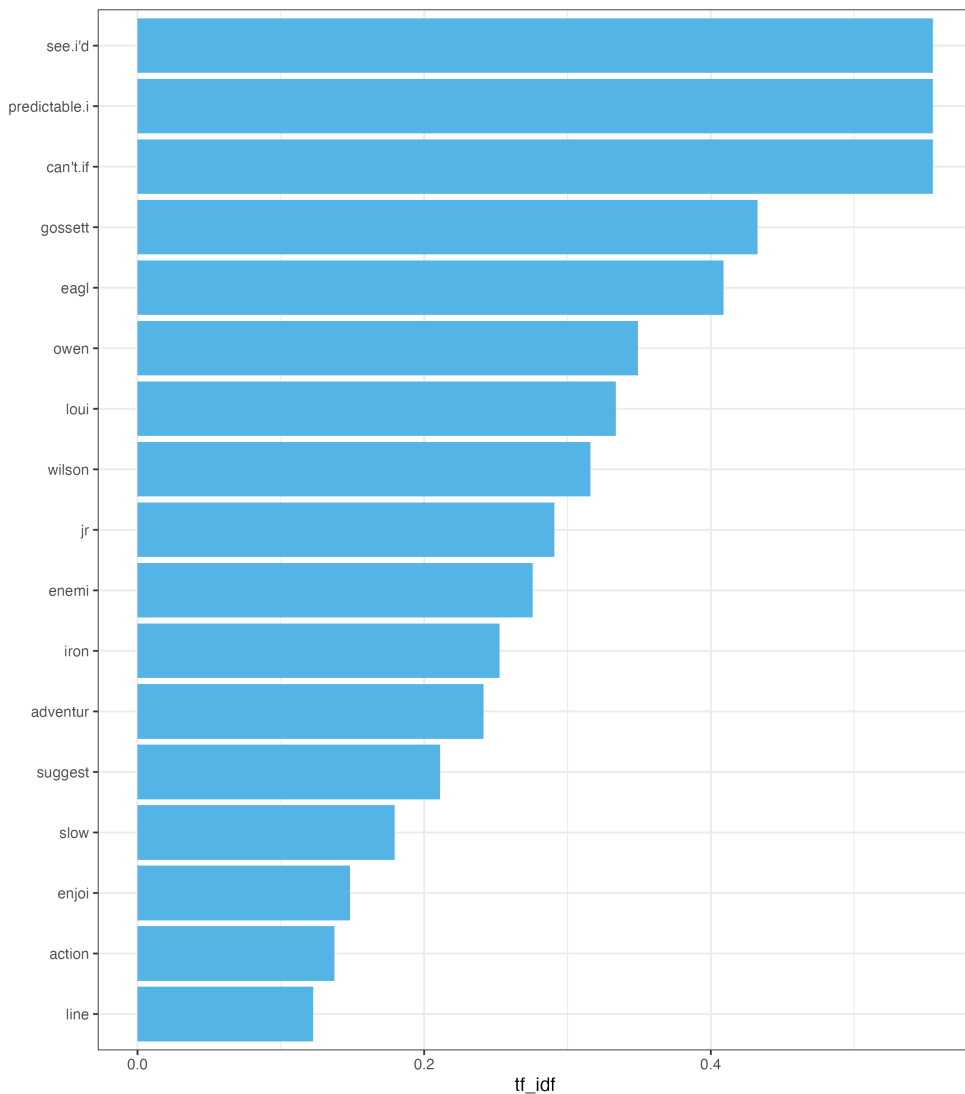
tf-idf weighting

```
1  tf_idf <- bind_tf_idf(imdb_counts,  
2                        word,  
3                        review,  
4                        n)  
5  
6  target <- sample(unique(imdb_neg$review), 1)  
7  
8  tf_idf |>  
9    filter(review == target) |>  
10   arrange(tf_idf) |>  
11   head(25) |>  
12   ggplot(aes(x = tf_idf,  
13             y = reorder(word, tf_idf))) +  
14   geom_col(fill = okabeito_colors(2)) +  
15   labs(y = NULL) +  
16   theme_bw()
```



Stemming

```
1  imdb_neg <- imdb_neg |>
2    mutate(stem = wordStem(word))
3
4  imdb_counts <- imdb_neg |>
5    group_by(review, stem) |>
6    summarize(n = n(), .groups = 'drop') |>
7    filter(!stem %in% c('br', 'movi', 'film'))
8
9  tf_idf <- bind_tf_idf(imdb_counts,
10                        stem,
11                        review,
12                        n)
13
14  tf_idf |>
15    filter(review == target) |>
16    arrange(tf_idf) |>
17    head(25) |>
18    ggplot(aes(x = tf_idf,
19              y = reorder(stem, tf_idf))) +
20    geom_col(fill = okabeito_colors(2)) +
21    labs(y = NULL) +
22    theme_bw()
```

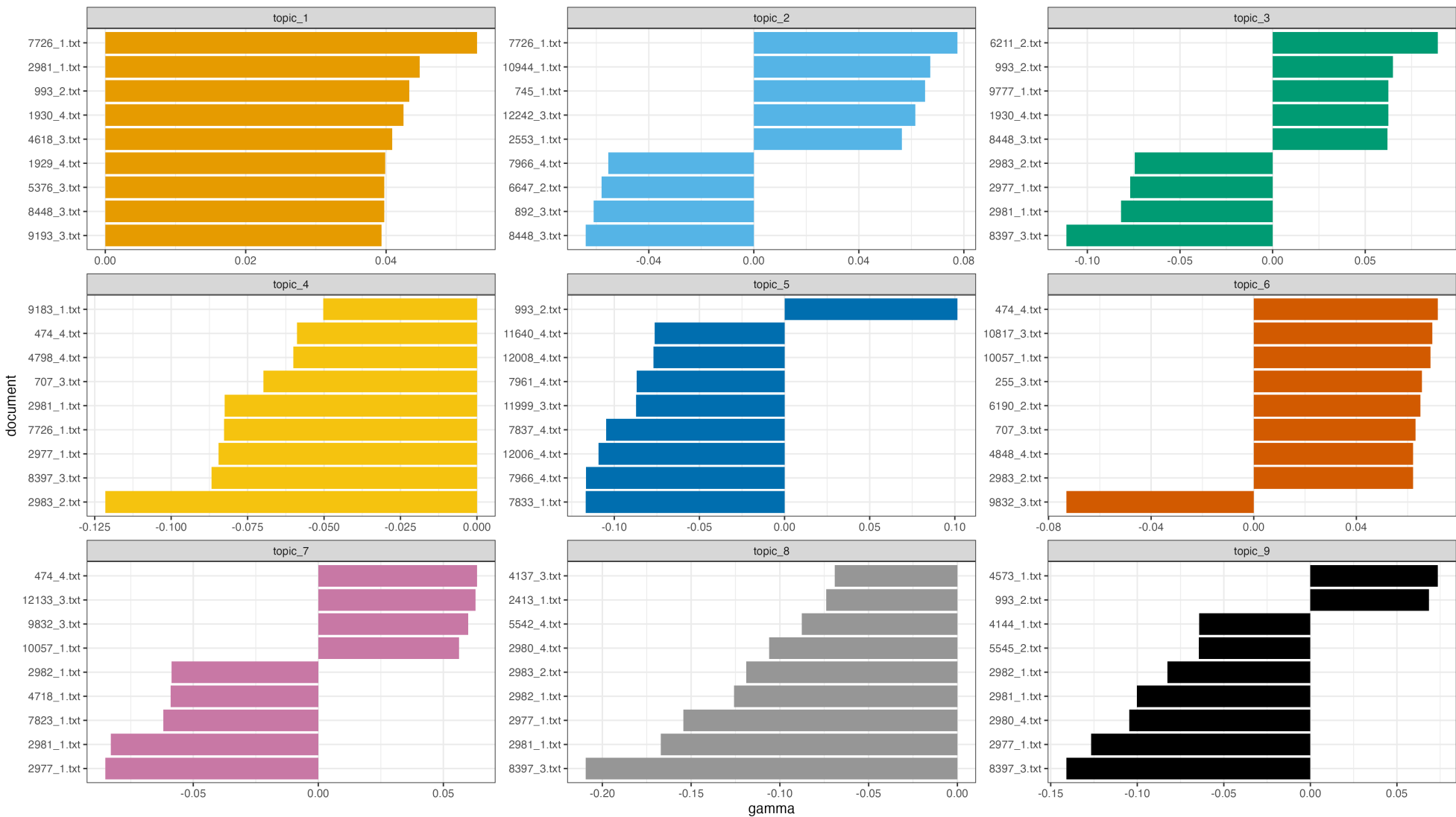


Construct DTM, conduct LSA, describe topics

```
1 dtm <- cast_dtm(imdb_counts, review, stem, n)
2
3 n_topics <- 9
4 imdb_lsa <- svds(as.matrix(dtm), n_topics)
5
6 imdb_topics <- as.data.frame(imdb_lsa$v)
7 colnames(imdb_topics) <- paste0('topic_', 1:n_topics)
8 imdb_topics$term <- colnames(dtm)
9
10 imdb_top_terms <- imdb_topics |>
11   pivot_longer(-term, names_to = 'topic', values_to = 'beta') |>
12   group_by(topic) |>
13   slice_max(abs(beta), n = 15) |>
14   ungroup() |>
15   arrange(topic, -beta)
16
17 imdb_top_terms |>
18   mutate(term = reorder_within(term, beta, topic)) |>
19   ggplot(aes(x = beta, y = term, fill = as.character(topic))) +
20   geom_col(show.legend = FALSE) +
21   facet_wrap(. ~ topic, scales = 'free') +
22   scale_fill_okabeito() +
23   scale_y_reordered() +
24   theme_bw()
```


Most prominent reviews by topic

```
1  imdb_reviews <- as.data.frame(imdb_lsa$u)
2  colnames(imdb_reviews) <- paste0('topic_', 1:n_topics)
3  imdb_reviews$document <- rownames(dtm)
4
5  imdb_top_reviews <- imdb_reviews |>
6    pivot_longer(-document, names_to = 'topic', values_to = 'gamma') |>
7    group_by(topic) |>
8    slice_max(abs(gamma), n = 9) |>
9    ungroup() |>
10   arrange(topic, -gamma)
11
12  imdb_top_reviews |>
13    mutate(document = reorder_within(document, gamma, topic)) |>
14    ggplot(aes(x = gamma, y = document, fill = as.character(topic))) +
15    geom_col(show.legend = FALSE) +
16    facet_wrap(. ~ topic, scales = 'free') +
17    scale_fill_okabeito() +
18    scale_y_reordered() +
19    theme_bw()
```



Use cosine similarity between right singular vectors to find similar reviews

```
1  U <- imdb_lsa$u
2  rownames(U) <- rownames(dtm)
3
4  norm <- sqrt(matrix(rowSums(U * U), byrow = F, nrow = nrow(U), ncol = nrow(U)))
5  cos_sim <- U %*% t(U) / (norm * t(norm))
6
7  cos_sim_df <- as.data.frame(cos_sim) |>
8    rownames_to_column('review_1') |>
9    pivot_longer(-review_1, names_to = 'review_2', values_to = 'cos_sim') |>
10   filter(review_1 != review_2) |>
11   filter(cos_sim < 0.98) |>
12   arrange(-cos_sim)
13
14 head(cos_sim_df, 5)
15
16 #   review_1   review_2   cos_sim
17 #   <chr>      <chr>      <dbl>
18 # 1 10223_4.txt 11884_4.txt 0.980
19 # 2 11884_4.txt 10223_4.txt 0.980
20 # 3 2297_4.txt  8846_2.txt  0.980
21 # 4 8846_2.txt  2297_4.txt  0.980
22 # 5 2505_2.txt  8018_1.txt  0.980
```

Let's look at an example:

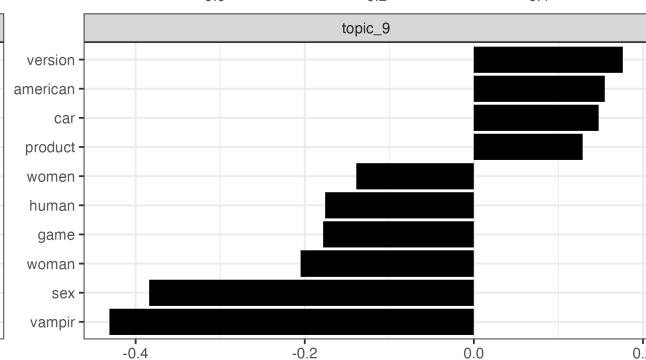
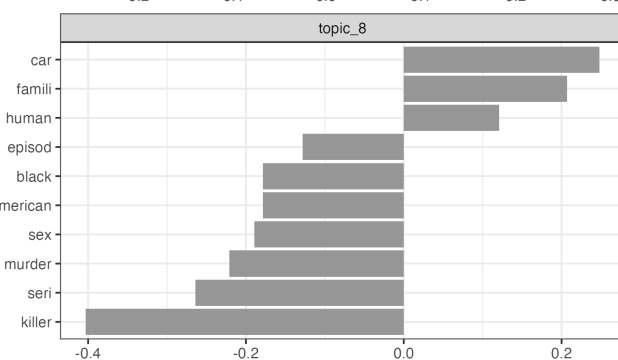
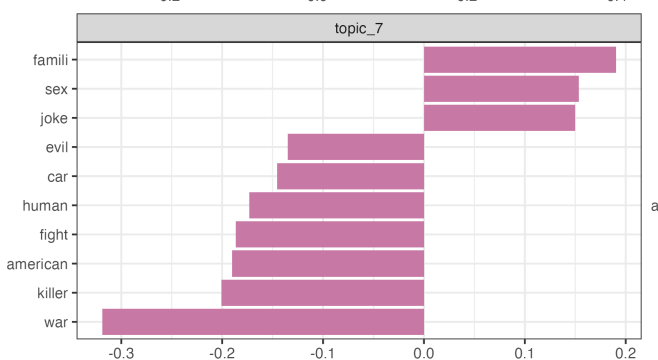
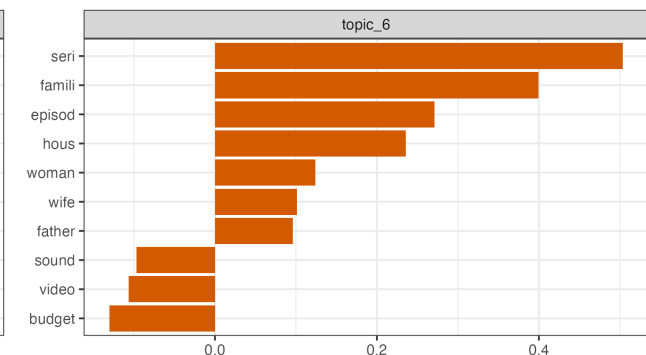
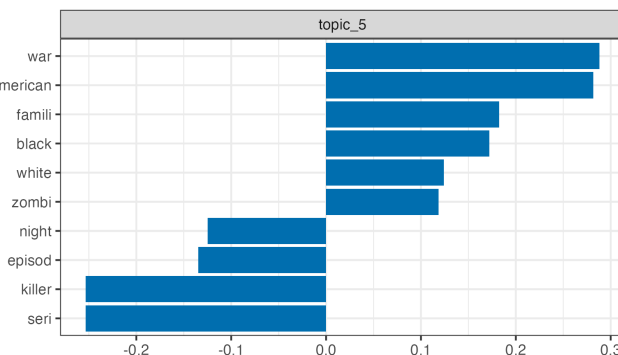
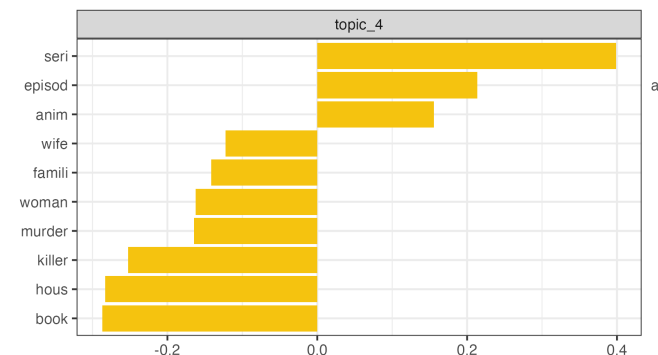
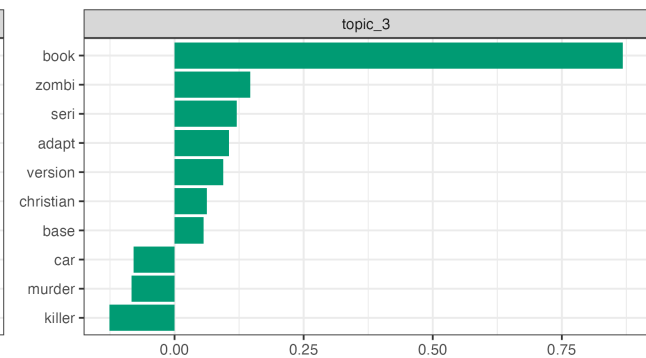
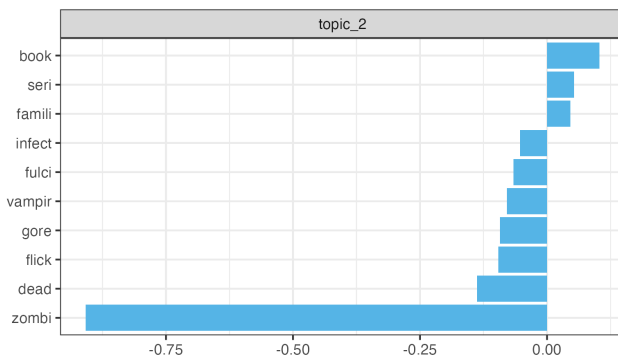
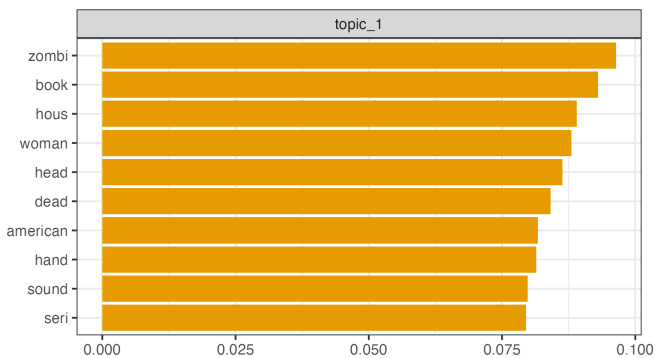
```
1 imdb_neg_raw |>
2   filter(review %in% c('10223_4.txt', '11884_4.txt')) |>
3   pull(text)
```

```
1 [1] "The worlds largest inside joke. The world's largest, most exclusive inside joke.<br /><br />Emulating the
2 brash and 'everyman' humor of office space, this film drives the appeal of this film into the ground by making
3 the humor such that it would only be properly appreciated by legal secretaries writing books. The audience is
4 asked to assume the unfamiliar role of a legal secretary, and then empathize with the excruciatingly dumb..."
5
6 [2] "I remembered the title so well. To me, it was a Flora Robson movie with Olivier and Vivien Leigh in
7 supporting roles. And it had Vincent Massey's voice from behind whiskers. Well Flora Robson was great. Her next
8 signature, for me, would be "\"55 Days at Peking\"". The same role but with different sumptuous gowns. And the
9 same voice. As for the Armada, it was a subtext. I like black-and-white films. Was everything done in..."
```

Let's filter out the common words:

```
1 filter_words <- tf_idf |>
2   select(stem, idf) |>
3   distinct() |>
4   filter(idf < 2.5)
5 filtered_imdb_counts <- anti_join(imdb_counts, filter_words, by = 'stem')
6 filtered_dtm <- cast_dtm(filtered_imdb_counts, review, stem, n)
7 f_imdb_lsa <- svds(as.matrix(filtered_dtm), n_topics)
8 f_imdb_topics <- as.data.frame(f_imdb_lsa$v)
9 colnames(f_imdb_topics) <- paste0('topic_', 1:n_topics)
10 f_imdb_topics$term <- colnames(filtered_dtm)
11 f_imdb_top_terms <- f_imdb_topics |>
12   pivot_longer(-term, names_to = 'topic', values_to = 'beta') |>
13   group_by(topic) |>
14   slice_max(abs(beta), n = 10) |>
15   ungroup() |>
16   arrange(topic, -beta)
17 f_imdb_top_terms |>
18   mutate(term = reorder_within(term, beta, topic)) |>
19   ggplot(aes(x = beta, y = term, fill = as.character(topic))) +
20   geom_col(show.legend = FALSE) +
21   facet_wrap(. ~ topic, scales = 'free') +
22   scale_fill_okabeito() +
23   scale_y_reordered() +
24   theme_bw()
```

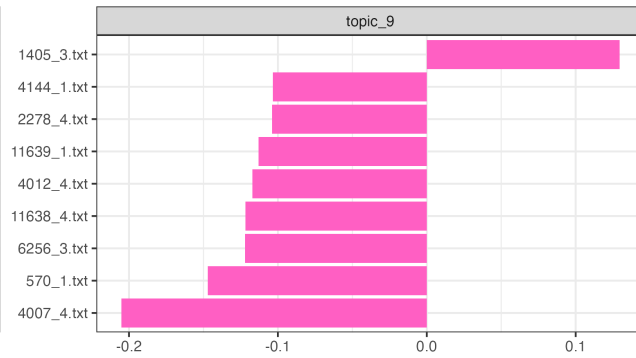
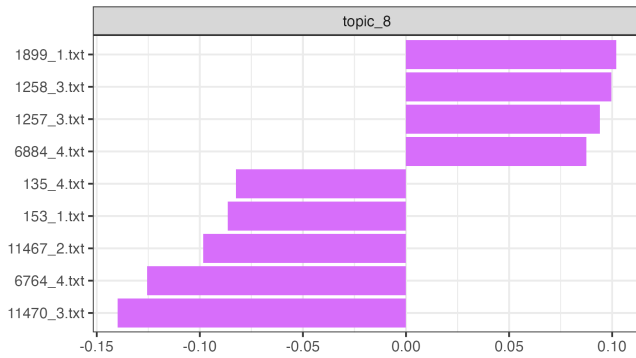
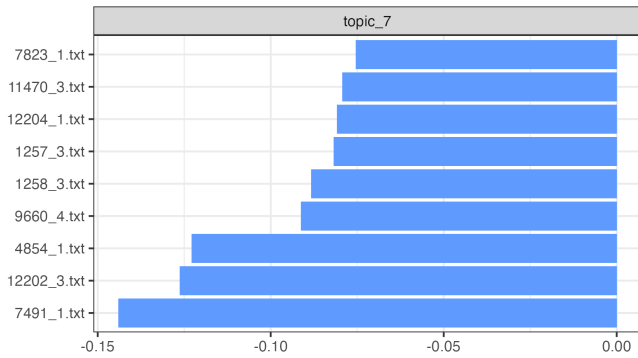
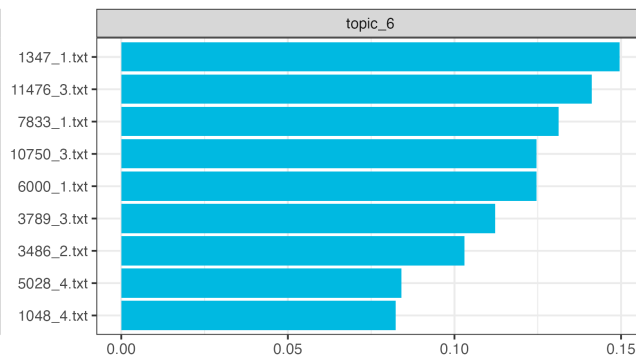
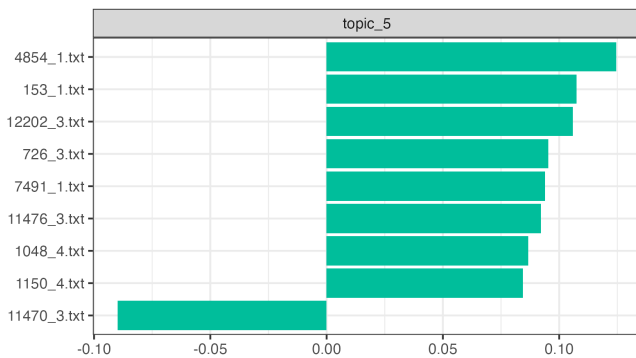
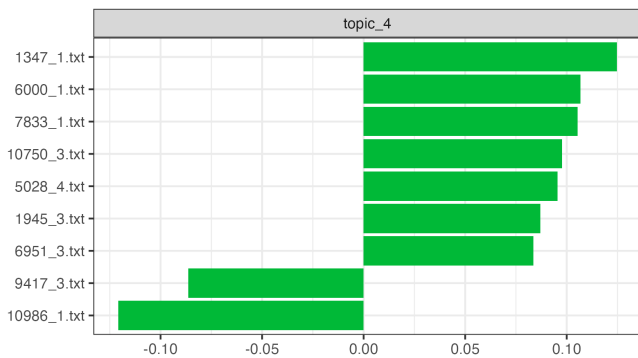
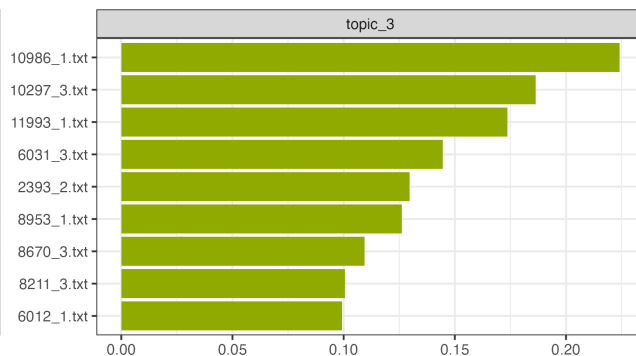
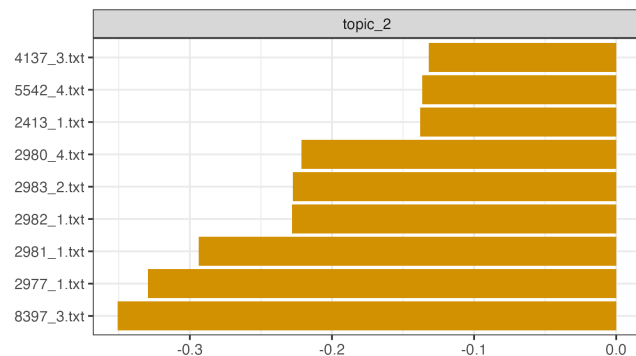
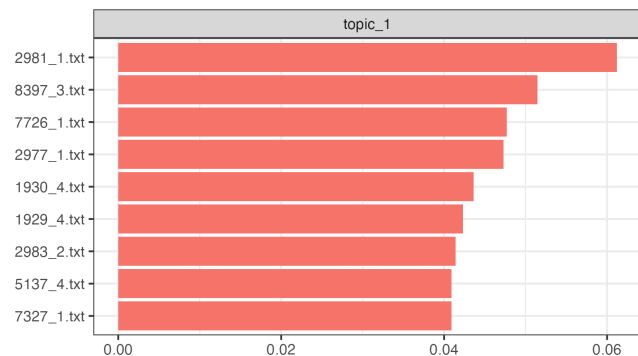
term



beta

Top reviews by topic

```
1 f_imdb_reviews <- as.data.frame(f_imdb_lsa$u)
2 colnames(f_imdb_reviews) <- paste0('topic_', 1:n_topics)
3 f_imdb_reviews$document <- rownames(filtered_dtm)
4
5 f_imdb_top_reviews <- f_imdb_reviews |>
6   pivot_longer(-document, names_to = 'topic', values_to = 'gamma') |>
7   group_by(topic) |>
8   slice_max(abs(gamma), n = 9) |>
9   ungroup() |>
10  arrange(topic, -gamma)
11
12 f_imdb_top_reviews |>
13   mutate(document = reorder_within(document, gamma, topic)) |>
14   ggplot(aes(x = gamma, y = document, fill = as.character(topic))) +
15   geom_col(show.legend = FALSE) +
16   facet_wrap(. ~ topic, scales = 'free') +
17   scale_y_reordered() +
18   theme_bw()
```



Example texts

■ Zombie topic:

```
1  imdb_neg_raw |>
2    filter(review %in% c('8397_3.txt')) |>
3    pull(text)
```

```
1  [1] "After reading the previous comments, I'm just glad that I wasn't the only person left confused, especially
2  by the last 20 minutes. John Carradine is shown twice walking down into a grave and pulling the lid shut after
3  him. I anxiously awaited some kind of explanation for this odd behavior...naturally I assumed he had something
4  to do with the evil goings-on at the house, but since he got killed off by the first rising corpse (hereafter
5  referred to as Zombie #1)..."
```

■ Book topic:

```
1  imdb_neg_raw |>
2    filter(review %in% c('10986_1.txt')) |>
3    pull(text)
```

```
1  [1] "When I saw that Mary Louise Parker was associated with this epic novel turned film, I was intrigued. Being
2  a fan of the book, I assumed she'd be playing Tony, Roz, or Charis, but more so, I was intrigued to see how they
3  would turn this very head-y, almost psychological (but not psychological thriller) novel in to a movie that
4  would be accessible to those who hadn't read the novel, and that would be at least mildly satisfying for those
5  who had. The book is a complex reflection of society, women, and modern life..."
```

Wrap Up

Recap

- Matrix factorization, via the SVD, shows up everywhere!
- It's a great way to break up more complicated data into easier to understand parts
- The ways in which we quantify text limit (or augment) how we can learn from it
- Most classical methods are *bag of words* models that just count instances of tokens
- Topic models function like PCA on the document term matrix
- Modern deep learning methods use vector embeddings of words to carry more information in each token

Final Thoughts

- [PollEv.com/klintkanopka](https://pollev.com/klintkanopka)