

Homework 3

This homework must be turned in on NYU Classes by **April 27, 2020, at 11pm**. Late work may be turned in up to 2 days late and will incur penalties of the equivalent of one third of a letter grade per day late.

It must be your own work, and your own work only—you must not copy anyone’s work, or allow anyone to copy yours. This extends to writing code. You may consult with others, but when you write up, you must do so alone.

Your homework submission must be a PDF or HTML report, containing all written answers and code, generated from RMarkdown. **Raw .R or .Rmd files will not be accepted.**

Please remember the following:

- Each question part should be clearly labeled in your submission.
- Do not include written answers as code comments. We will not grade code comments.
- The code used to obtain the answer for each question part should accompany the written answer.
- **Your code must be included in full, such that your understanding of the problems can be assessed.**

Please consult “RMarkdown Basics” on the course GitHub for help with RMarkdown. You can also use the “Sample RMarkdown HW” template on the course GitHub to get started. Using this template is not required.

You should continue to use quanteda version 1.5.2 for this assignment.

For the following exercises, it is recommended that you use the following packages: `topicmodels`, `lda`, and `stm`.

1. Applying `topicmodels` to the news corpus:

Note: Because topic models take a long time to run, you may save your fitted model objects or RStudio workspace to a .RData file after running the appropriate code to fit the models. When doing your final write-up, you may comment out the code that trains the model(s) and load the saved models from file. Please make sure the code to fit the models is clearly marked.

- (a) To decrease the time it takes to fit a topic model, we will limit our analysis to a subset of the immigration corpus. Create a subset of `data_corpus_immigrationnews` that only contains articles from the following news sources: `telegraph`, `guardian`, `ft`, `independent` and `express`. Create a table that shows how many documents are associated with each newspaper.
- (b) Create a document term matrix with your new immigration corpus in which punctuation and numbers are removed and words are stemmed and set to lower case. Also, remove a custom set of stopwords `custom_stopwords` that is relevant to this particular data set. Finally, use `quanteda`'s "`dfm_trim`" to remove words that occur fewer than 30 times or in fewer than 20 documents. Report the remaining number of features and the total number of documents in the DFM.
- (c) Preprocessing decisions can have substantive impacts on the topics created by topic model algorithms. Make a brief (1 paragraph) argument for or against removing rare terms from a dfm on which you plan to fit a topic model.
- (d) Fit a topic model with 25 topics using `LDA()`, with `method = "Gibbs"`. Increase the number of iterations to 3000 to ensure that the model describes the underlying data well and set the seed to 1234 so that you can replicate your results. Report the `@loglikelihood` of your topic model object.
- (e) Examine the top 10 words that contribute the most to each topic using `get_terms()`. You should save the top 10 words over all 25 topics, for later use. Next, find the most likely topic for each document using `topics()`. Create a table with one column of topic numbers one column containing the number of documents for which that topic is the most likely topic. Sort the table in decreasing order by the second column. Show this table in your answer. Give substantive labels to the top 5 topics (i.e. by looking at the most likely words within each of these topics). Briefly explain your choice of labels.
- (f) Examine the topics that contribute the most to each document, using the code from recitation to visualize the top two topics per document for the Independent and the Financial Times (`ft`) with separate graphs for each newspaper. Make sure that the documents are sorted by day of publication (the "`day`" variable in the `data_corpus_immigrationnews` corpus). Discuss your findings.
- (g) Finally, we can find the average contribution of a topic to an article from a particular newspaper, and compare newspapers on particular topics. For each of the 5 topics you've named, see how their prevalence varies among the different newspapers. To do so, estimate the mean contribution of each topic over each newspaper. Report the contribution of each of the top 5 topics to each of the 5 newspapers. Discuss your findings.

2. **Topic stability:** We want to see how stable these topics are, under two different topic parameter values.

- (a) Re-run the model from question 1 with a different seed. Report the @loglikelihood of your topic model object.
- (b) For each topic in the new model, find the topic that is the closest match in the original run in terms of cosine similarity of the topic distribution over words. Your answer should be a table.
- (c) Calculate the number of words in the top 10 words shared by each matched topic pair. Your answer should be a table.
- (d) Now run two more models, but this time, use only 5 topics. Again, find the average number of words in the top ten shared by each matched topic pair. How stable are the models with 5 topics compared to the models with 25 topics?

3. Topic Models with covariates:

The Structural Topic Model (STM) is designed to incorporate document-level variables into a standard topic model. Since, we have information about both the newspaper and the date of the articles, we can use an STM (from the `stm` package) to model the effects of these covariates directly.

- (a) Using only articles from the Guardian and Financial Times (ft), construct a numeric date variable from the “day” variable in the immigration news corpus. Use what preprocessing you believe to be appropriate for this problem. Discuss your preprocessing choice.
- (b) Fit an STM model where the topic content varies according to this binary variable, and where the prevalence varies according to both this binary variable and the spline of the date variable you’ve created. Be sure to use the spectral initialization and set `k=0`, which will allow the STM function to automatically select a number of topics using the *spectral learning* method. Keep in mind that this function is computationally demanding, so start with the minimum threshold document frequency threshold set to 10; if your computer takes an unreasonably long time to fit the STM model with this threshold, you can raise it to as high as 30. Report the number of topics selected in the fitted model. Also report the number of iterations completed before the model converged.
- (c) Identify and name each of the 5 topics that occur in the highest proportion of documents

using the following code:¹

```
plot(fit.stm, type = "summary")
```

- (d) Using the visualization commands in the `stm` package, discuss one of these top 5 topics. How does the content vary with the paper discussing that topic? How does the prevalence change over time?
4. **Non-Parametric Scaling - Wordfish:** Recall that the Wordfish algorithm allows us to scale political texts by a latent dimension. We will apply this function to analyze the UK manifestos.
- (a) First, create a corpus that is the subset of the `data_corpus_ukmanifestos` that contains only speeches by the Conservative ('Con') and Labor ('Lab') parties.
- (b) Quanteda's implementation of Wordfish, `textmodel.wordfish`, requires that we provide the indices for a pair of documents to globally identify θ , the latent dimension of interest (e.g. lower values of θ = more Liberal, higher values of θ = more Conservative). In this case, we are looking to estimate the latent left-right ideological dimension. Use the indices of the 1979 Labor and Conservative manifestos to do so. That is, set `dir = c(index of 1979 Labor manifesto, index of 1979 Conservative manifesto)`.
- (c) Which of the documents is the most left wing? Which is the most right-wing? Are these results surprising? Why or why not?
- (d) Re-create the "guitar plot" from recitation. Describe the parameters estimated by Wordfish that lie on the axes of the plot.
- (e) **Optional:** Estimate a linear regression with the Wordfish score as the dependent variable and binary variable indicating whether or not a Manifesto was from the Labor party. Include a binary control variable for each manifesto. If we use being Labor as a proxy for 'liberal' ideology (in the American sense of the word), how well did our Wordfish model do at capturing latent ideology?²
5. **Burstiness:** Here we evaluate the burstiness of several words using the `news_data` corpus of news headlines. To evaluate burstiness we will use the `bursts` package and the user-written function `bursty` from recitation that visualizes the results. You have been provided the `news_data` corpus.
- (a) Create a corpus. For each of the words "obama", "korea", and "afghanistan" use the `bursty` function to visualize the burst period(s) and levels. Also, for each of the plots

¹`fit.stm` Represents the output of the STM model you fit in the preceding question.

²If it did well, then our proxy variable for ideology should be significant at at least a 5% level.

include a brief interpretation about what the timing and level of the burst may indicate.

6. Dimension Reduction and Semantics: For this question use `news_data`. To reduce computation time, use the first 1000 headlines from the “WORLD NEWS” category.

- (a) Obtain the document feature matrix (DFM) of the corpus, removing stopwords, punctuation and lower-casing. Perform a principal components analysis on the resulting DFM and rank the words on the first principal component according to their loadings. Report the top 5 with the most positive loadings and the top 5 with the most negative loadings. Is the first principal component interpretable? If so, what would be your interpretation of what it is capturing?
- (b) Using the `lsa` package, estimate a latent semantic analysis model. According to the resulting term vector matrix, report the 5 nearest tokens to `america` and `corruption`. Did the model do a good job of capturing ‘meaning’? In other words, do the nearest neighbors for these words make sense?
- (c) Load the pretrained GloVe embeddings provided (`glove.rds`). Using these embeddings, find the nearest neighbors to `america` and `corruption`. Do these make sense? How do they compare to the nearest neighbors using LSA (keep in mind, these embeddings were estimated on a much larger corpus comprising Wikipedia and Google News articles)?