Unsupervised Text Analysis

Gov 51: Section 9

Sima Biondi Spring 2025

Overview

1 Housekeeping

- 2 Unsupervised learning overview
- 3 Implementation 1: polisci publishing Results
- 4 Implementation 2: text analysis for description Results

5 Summary

6 Office hours

• Reminder of deadlines

- Reminder of deadlines
 - $\cdot~$ April 11th \rightarrow first draft of poster

- Reminder of deadlines
 - $\cdot~$ April 11th \rightarrow first draft of poster
 - \cdot April 24th \rightarrow final draft of poster

- Reminder of deadlines
 - $\cdot~$ April 11th \rightarrow first draft of poster
 - $\cdot~$ April 24th \rightarrow final draft of poster
 - April 25: Pset 3 due

- Reminder of deadlines
 - $\cdot~$ April 11th \rightarrow first draft of poster
 - \cdot April 24th \rightarrow final draft of poster
 - April 25: Pset 3 due
 - \cdot April 29rd \rightarrow poster session

- Corpus a collection of texts that we want to analyze
 - Complete works of Charlotte Bronte, newspaper articles from the AP
- Document a unit within the corpus
 - Jane Eyre; an article from the AP

Unsupervised learning example 1: peer review bias

- Lots of criticism of the peer review model send your article in to two double blind reviewers
 - Despite double blind, questions about anonymity people post their papers online now
 - Reviewers can be quite biased as well!
- Two questions arise:
 - 1. Who is publishing in top political science journals? Are their ascriptive disparities?
 - 2. What is being published? Are some topics avoided?

- \cdot One advantage of unsupervised learning is pattern finding
 - Data could be high-dimensional, messy, huge observations (e.g. Jane Eyre) ...
- Patterns within the data are incredibly useful!
 - $\cdot\,$ For example, the topics found in a flagship political science journal
- One application is **topic modeling**

Latent Dirichlet Allocation (LDA)

- When we topic model, we need to make some assumptions about how text is generated
 - Recall basic assumptions about topics \rightarrow probability distribution of words
- Latent Dirichlet Allocation is just one model we can use to topic model
 - LDA relies on a similar assumption of the data generating process of text
- Steps:
 - 1. For each topic, draw a topic-word distribution \rightarrow how likely is a topic given a word?
 - 2. For each document, draw a document-topic distribution \rightarrow how likely is a document about certain topics?

"I had donuts this morning. I don't have diabetes yet."

LDA steps

- 1. For each topic, draw a topic-word distribution.
 - food: donuts (0.7), have (0.3), diabetes (0)
 - health: diabetes (0.7), have (0.3), donuts (0)
- 2. For each document in the corpus, draw a document-topic distribution.
 - Sentence 1: food (0.999), health (0.001)
 - Sentence 2: food (0.001), health (0.999)

- 'topicmodels' package is useful for implementing an LDA model
- In our example here, we will be taking a look at a sample of articles in the AP in 1992

library(topicmodels)
data("AssociatedPress")

Preview

Data is at the document-word level - for each document, each unique word is counted

head(tidy(AssociatedPress))

tibble: 6×3			
document <int></int>	term <chr></chr>	count <dbl></dbl>	
1	adding	1	
1	adult	2	
1	ago	1	
1	alcohol	1	
1	allegedly	1	
1	allen	1	

Function

- 'k': number of topics we want to specify
- 'control': setting a seed because probability distributions entail randomness

set a seed so that the output of the model is
predictable
ap_lda <- LDA(AssociatedPress, k = 2, control =
list(seed = 02138))
ap_lda

A LDA_VEM topic model with 2 topics.

Results

- 'beta' refers to the probability that a given word is related to a topic
- · Let's see which topic "harvard" is associated with

```
ap_topics <- tidy(ap_lda, matrix = "beta")
ap_topics |>
filter(term == "harvard")
```

topic <int></int>	term <chr></chr>	beta <dbl></dbl>
1	harvard	4.018027e-05
2	harvard	1.202947e-04

Instead of looking for specific words, let's visualize the most likely terms per topic

```
top_terms <- ap_topics |>
  group_by(topic) |>
  slice_max(beta, n = 10) |>
  ungroup() |>
  arrange(topic, -beta)

top_terms <- top_terms |>
  mutate(term = reorder_within(term, beta, topic))
```

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

Visualization



Document level visualization

- How about topic likelihoods at the document level?
- 'gamma' gives us the likelihood of a topic given the words in a document

```
ap_documents <- tidy(ap_lda, matrix = "gamma") |>
    arrange(document, topic)
head(ap_documents)
```

1 1 0.99932253
1 2 0.00067747
2 1 0.51357042
2 2 0.48642958
3 1 0.96200507
3 2 0.03799493

- 1. Who is publishing in top political science journals? Are their ascriptive disparities?
- 2. What is being published? Are some topics avoided?

- 1. Who is publishing in top political science journals? Are their ascriptive disparities?
- 2. What is being published? Are some topics avoided?

Saraceno (2020) did an analysis of publications in The Journal of Politics

Back to Saraceno (2020)

Interest Groups Female (Solo & Teams) Mixed-Gender Teams Male (Solo & Teams) Public Administration Partisanship Voting Behavior Political Economy Race & Ethnicity Approval Ratings Media & Politics Gender & Politics Electoral Politics Political Theory Law & Courts-Democratic Theory Legislative Studies Experimental Formal Theory Survey Methods International Affairs Comparative Studies Statistical 50 Ó 100 150 Number of Articles

Saraceno (2022)

Table 2. Topic Model Output

Topic Label

Interest Groups Political Theory Gender and Politics Political Economy Race and Ethnicity Democratic Theory The South Electoral Politics Survey Methods Postwar Politics International Affairs **Judicial** Politics Formal Theory Public Administration Statistical Legislative Politics **Comparative** Politics Partisanship Media and Politics Voting Behavior Approval Ratings Experimental Methods High-Frequency Words

Groups, interest, members, organizations, activity, influence Man, life, human, nature, thought, Hobbes, philosophy Women, social, participation, gender, female, men, politics Economic, income, fiscal, growth, welfare, market, inequality Black, white, racial, ethnic, school, minority, representation Rights, democratic, moral, public, liberal, justice, freedom States, southern, county, Negro, governor, Virginia, Carolina Elections, candidates, district, incumbent, office, primary Respondents, information, survey, attitudes, treatment Soviet, war, world, national, social, united, communist Conflict, foreign, military, international, war, aid, civil Court, supreme, cases, judicial, law, courts, justice, legal Model, decision, equilibrium, preferences, choice, probability Local, government, service, agencies, city, administer, board Models, data, variables, effect, results, analysis, significant Congress, House, legislative, committee, Senate, majority, bill Countries, international, institution, regime, corrupt, Latin Issues, ideological, opinion, liberal, conservative, polarization Media, coverage, exposure, negative, television, advertising Voter, campaign, candidate, turnout, electorate, ballot President, economic, support, approval, performance Experimental, subjects, control, condition, assigned, effect

- Reiterating caveats at the end of lecture
- Topic modeling is **not** a panacea!
 - Similar to other methods, it relies on assumptions, particularly about the DGP of text
 - Uncertainty is also a part of predictions similar in respect to regression predictions which have their own standard errors

library(tm) library(SnowballC)

- In what ways can we categorize and divide the Harvard Government faculty?
- Let's say we have a **corpus** with three variables in '.csv' form
 - 1. 'prof' name of faculty member
 - 2. 'phd' year of phd attainment
 - 3. 'bio' biography on website



Pre-processing

10

14

```
# make everything lowercase
corpus <- tm_map(corpus, content_transformer(tolower))</pre>
# remove white space (e.g. spaces)
corpus <- tm map(corpus, stripWhitespace)</pre>
# remove numbers
corpus <- tm_map(corpus, removeNumbers)</pre>
# remove stopwords
corpus <- tm_map(corpus, removeWords,</pre>
    stopwords("english"))
# stem words (e.g. remove "ing")
corpus <- tm map(corpus, stemDocument)</pre>
```

```
# Turning into a document term matrix
dtm <- DocumentTermMatrix(corpus)
dtm.mat <- as.matrix(dtm)
# Adding labels to each document
rownames(dtm.mat) <- df$prof</pre>
```

```
# Normalize by document size
tfidf <- weightTfIdf(dtm, normalize = TRUE)
tfidf.mat <- as.matrix(tfidf, normalize = TRUE)
# Adding labels to each document
rownames(tfidf.mat) <- df$prof</pre>
```

```
par(cex = 1.25)
library("wordcloud")
liu <- dtm.mat["naijia liu",]
liu.tfidf <- tfidf.mat["naijia liu",]</pre>
```





```
sort(liu, decreasing = TRUE)[1:5]
sort(liu.tfidf, decreasing = T)[1:3]
```

polit research current includ professor 2 2 1 1 1 demographi facebook imputation, 0.1745301 0.1745301

K-means algorithm

- K-Means clustering is simply an exercise in partitioning *n* observations into *k* clusters
- In a text context, this means comparing the similarity of words between documents
- Here, we are looking to cluster professors based on similar word usage in their bios!
- This is an iterative process the algorithm does initial groupings, then sees whether it can minimize error by another permutation

Descriptive stats

table(kconfour.out\$cluster)

1 2 3 4 5 11 14 5 7 11

x

x stephen ansolabehere nara dillon grzegorz ekiert peter a. hall jennifer hochschild alastair iain johnston taeku lee elizabeth j. perry michael rosen james m. snyder, jr richard tuck

danielle allen eric beerbohm daniel carpenter timothy colton katrina forrester claudine gay harvey c. mansfield eric nelson paul peterson stephen peter rosen michael sandel theda skocpol latanya sweeney daniel ziblatt

knitr::kable(df\$prof[kconfour.out\$cluster == 1])

X

stephen ansolabehere nara dillon grzegorz ekiert peter a. hall jennifer hochschild alastair iain johnston taeku lee elizabeth j. perry michael rosen james m. snyder, jr richard tuck

K-means group 2

knitr::kable(df\$prof[kconfour.out\$cluster == 2])

X

danielle allen eric beerbohm daniel carpenter timothy colton katrina forrester claudine gay harvey c. mansfield eric nelson paul peterson stephen peter rosen michael sandel theda skocpol latanya sweeney daniel ziblatt

knitr::kable(df\$prof[kconfour.out\$cluster == 3])

x jeffry frieden frances hagopian alisha c. holland torben iversen steven levitsky

knitr::kable(df\$prof[kconfour.out\$cluster == 4])

Х

melani cammett stephen chaudoin christina davis joshua d. kertzer christoph mikulaschek pia raffler yuhua wang

knitr::kable(df\$prof[kconfour.out\$cluster == 5])

X

matthew blackwell peter buisseret ryan enos chase h. harrison michael j. hiscox kosuke imai gary king naijia liu mashail malik stephanie ternullo dustin tingley

Overtime comparisons

10

• How similar are each bio to the professor with the earliest PhD, Harvey Mansfield?

```
# Isolate Harvey
harvey <- as.data.frame(tfidf.mat["harvey c.</pre>
    mansfield",])
# Isolate non-Harveys
nonhm.tfidf <-</pre>
    as.data.frame(tfidf.mat[rownames(tfidf.mat)
                           != "harvey c. mansfield", ])
# Sort everything chronologically by PhD attainment
nonhm.tfidf$year.index <-</pre>
    df$phd[df$prof!= "harvey c. mansfield"]
chron.tfidf <-
    nonhm.tfidf[order(nonhm.tfidf$year.index),]
years <- sort(unique(df$phd))</pre>
vears <- vears[-1]</pre>
```

For-loop

8

```
cosine <- function(a, b) {
  ## t() transposes a matrix ensuring that vector `a'
    is multiplied ## by each row of matrix `b'
  numer <- apply(a * t(b), 2, sum)
  denom <- sqrt(sum(a^2)) * sqrt(apply(b^2, 1, sum))
  return(numer / denom)
  }</pre>
```

```
avg.cosim <- rep(NA, length(years))
for (i in 1:length(years)) {
    decade <- subset(chron.tfidf,
                                (year.index == years[i]))
    decade <- decade[, names(decade) != "year.index"]
    similarity <- cosine(harvey, decade)
    avg.cosim[i] <- mean(similarity)
}</pre>
```

Plotting similarity to Harvey across time



- Pre-processing text in an easy-to-implement way
- Also, pre-pre-processing when our data isn't already a document term matrix
- Learned one way to group text based on similarity (e.g. k-means algorithm)
- Using for-loops and our own cosine similarity function, we can plot similarity over time

- Bring all your questions!
- Happy to help on code, identification, etc!