

# **Homework 8: Topic Modeling using Latent Dirichlet Allocation**

**Student:** Patrick Walsh  
**Professor:** Dr. Ami Gates  
**School:** Syracuse University  
**Date:** 3/19/2024

## INTRODUCTION

In today's data-driven world, extracting meaningful insights from vast amounts of text data is crucial for informed decision-making. The "110" dataset, encompassing the floor debate transcripts of the 110th Congress (House only), presents a valuable opportunity to leverage advanced analytics techniques like topic modeling.

Topic modeling offers a powerful way to uncover hidden patterns and themes within large text datasets. For businesses, this translates into actionable intelligence that can drive various strategic initiatives such as:

- **Policy Analysis:** By identifying prevalent topics in congressional debates, policymakers can gain a comprehensive understanding of legislative priorities, enabling them to make informed decisions on policy formulation and amendment.
- **Public Opinion Monitoring:** Analyzing debates can provide insights into public sentiment and concerns, aiding businesses in understanding societal trends and aligning their strategies accordingly.
- **Competitive Intelligence:** Topic modeling can uncover key issues and discussions within specific political or demographic groups, offering businesses valuable insights into their competitors' priorities and strategies.
- **Content Curation:** Understanding the dominant topics in congressional debates helps in curating relevant content for stakeholders, customers, and the public, enhancing engagement and brand perception.
- **Risk Management:** Detecting emerging topics or controversial discussions allows businesses to proactively address potential risks and challenges, mitigating negative impacts on their operations.

## ANALYSIS

### Dataset

The dataset used in this analysis, referred to as "110," comprises text documents containing the floor debate of the 110th Congress, focusing on the House of Representatives. The dataset is organized into four subfolders, where "m" denotes male

speakers, "f" denotes female speakers, "d" denotes Democrats, and "r" denotes Republicans. The objective of this analysis is to gain insights into the topics discussed during the floor debates using natural language processing techniques.

```
Female Democrats:
['110_baldwin_x_wi.txt', '110_bean_x_il.txt', '110_berkley_x_nv.txt', '110_boyda_x_ks.txt', '110_brown-waite_x_fl.txt', '110_brown_x_fl.txt', '110_capps_x_ca.txt', '110_castor_x_fl.txt', '110_clarke_x_ny.txt', '110_davis_x_ca.txt']
Female Republicans:
['110_bachmann_x_mn.txt', '110_biggett_x_il.txt', '110_blackburn_x_tn.txt', '110_bono_x_ca.txt', '110_capito_x_wv.txt', '110_cubin_x_wy.txt', '110_davis_jo-ann_va.txt', '110_drake_x_va.txt', '110_emerson_x_mo.txt', '110_fallin_x_ok.txt']
Male Democrats:
['110_abercrombie_x_hi.txt', '110_ackerman_x_ny.txt', '110_allen_x_me.txt', '110_altmire_x_pa.txt', '110_andrews_x_nj.txt', '110_arcuri_x_ny.txt', '110_baca_x_ca.txt', '110_baird_x_wa.txt', '110_barrow_x_ga.txt', '110_becerra_x_ca.txt']
Male Republicans:
['110_aderholt_x_al.txt', '110_akin_x_mo.txt', '110_alexander_x_la.txt', '110_bachus_x_al.txt', '110_baker_x_la.txt', '110_barratt_x_sc.txt', '110_bartlett_x_md.txt', '110_barton_x_tx.txt', '110_bilbray_x_ca.txt', '110_bishop_x_ut.txt']
```

The dataset can be found here:

<https://drive.google.com/drive/folders/1hf3ALdntNIm6r-te2z8yAM26e33F0WFt>

## Data Preparation

The data preparation phase involved combining and cleaning the text data to create a document-term matrix (DTM) suitable for topic modeling. This process included removing irrelevant columns, special characters, and numbers, ensuring that the data is ready for analysis. The text data was vectorized using CountVectorizer() and converted to a dataframe.

## Tokenize & Vectorize

```
: 1 from sklearn.feature_extraction.text import CountVectorizer
2
3 # create CV objects
4 CV = CountVectorizer(input='filename', stop_words='english', encoding='latin-1')
5
6 # create document term matrix
7 dtm = CV.fit_transform(master_list)
8
9 print(dtm[0])
```

```
(0, 18407)    94
(0, 18419)    94
(0, 37175)    95
(0, 7443)     94
(0, 20182)     4
(0, 18111)     4
(0, 22696)    10
(0, 22013)     5
(0, 18000)     0
```

```
1 import re
2
3 def clean_cols(df):
4     # drop columns with numbers
5     columns_to_drop = [col for col in df.columns if any(char.isdigit() for char in col)]
6     df = df.drop(columns = columns_to_drop)
7
8     # drop columns with special characters
9     pattern = '^([a-zA-Z])+$' # regex pattern to match only letters from a-z or A-Z
10    columns_to_drop = [col for col in df.columns if not re.match(pattern, col)]
11    df = df.drop(columns = columns_to_drop)
12
13    # drop columns that have a length of less than 3 or more than 15
14    columns_to_drop = [col for col in df.columns if len(col) < 3 or len(col) > 15]
15    df = df.drop(columns = columns_to_drop)
16
17    # drop columns that have three or more consecutive occurrences of the same letter
18    pattern_consecutive = r'([a-zA-Z])\1\1'
19    columns_to_drop = [col for col in df.columns if re.search(pattern_consecutive, col)]
20    df = df.drop(columns=columns_to_drop)
21
22    return df
23
24 cleaned_df = clean_cols(df)
25 display(cleaned_df)
```

	aachen	aap	aapg	aapi	aapis	aaron	aaronsen	aaronson	aarp	aarti	...	zurich	zury	zvi	zweig	zwilling	zwyer	zwyers	zydeco	zygote	zyuganov
0	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0
1	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0
2	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0
3	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0
4	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
424	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0
425	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0
426	0	0	0	0	0	0	0	0	0	2	...	0	0	0	0	0	0	0	0	0	0
427	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0
428	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0

429 rows × 57065 columns

## Models

The analysis utilized Latent Dirichlet Allocation (LDA) as the main modeling technique. LDA is a probabilistic topic modeling approach that identifies underlying topics within a collection of documents.

The LDA model was set to 100 max iterations and a learning method of ‘online’ rather than ‘batch’ so that words and documents were processed one at a time rather than in batches.

```
def run_LDA(df, num_topics, colnames):
    lda_model = LatentDirichletAllocation(n_components=num_topics, max_iter=100, learning_method='online')

    # Start the timer
    start_time = time.time()

    # Fit the model with progress bar
    with tqdm(total=len(df)) as pbar:
        lda_fitted = lda_model.fit_transform(df)
        pbar.update(len(df))
```

## RESULTS

The LDA model was trained with varying numbers of topics (5, 10, 20, and 30) to explore different levels of granularity in topic identification. The top words representing each topic were extracted, providing insights into the main themes discussed during the floor debates.

### 5 Topics

Topic #0	Topic #1	Topic #2	Topic #3	Topic #4
cornwells simpler coroners friedberg rightsize sedan plentiful 1318 firewalls naacp beverage handbasket smutty odious sludge	simpler cornwells friedberg coroners naacp europeans rightsize highland defines sludge size plentiful 3014 testicular shiver	simpler cornwells friedberg coroners friedberg 1318 rightsize plentiful sludge sedan befalls rumsfeld su lovelace incinerator breastfeeding	truncheons hemant ix madigan ah cut lifting bleck centcom hydrographic cottoned beachside baffled chains badminton	friedberg simpler cornwells sedan rightsize coroners plentiful handbasket naacp 1318 smutty bunk 3014 bossed incinerator

## 10 Topics

Topic #0	Topic #1	Topic #2	Topic #3	Topic #4
simpler cornwells coroners friedberg 1318 sedan plentiful rightsize sludge befalls incinerator smutty lovelace madigan breastfeeding	simpler friedberg coroners cornwells rightsize naacp 1318 plentiful bunk 3014 odious rumsfeld sludge smutty sedan	friedberg simpler coroners cornwells europeans naacp rightsize highland defines size sludge plentiful shiver 3014 testicular	simpler friedberg coroners cornwells sludge plentiful 1318 3014 naacp bunk lifting befalls firewalls sedan rumsfeld	cornwells coroners simpler rightsize friedberg naacp plentiful odious firewalls sludge bagpipes handbasket 3014 beverage carrhart
Topic #5	Topic #6	Topic #7	Topic #8	Topic #9
simonson 1318 rightsize simpler coroners cornwells incinerator friedberg rumsfeld sedan su sludge 3014 naacp bunk	simpler cornwells coroners friedberg naacp rightsize 1318 sludge plentiful bunk smutty 3014 sedan odious bossed	friedberg cornwells simpler coroners rightsize sludge naacp 1318 incinerator plentiful sedan 3014 trophy smutty testicular	handbasket rightsize simpler cornwells coroners friedberg naacp testosterone su sedan rumsfeld madigan sludge odious transparency	naacp rumsfeld su antiaircraft carrhart 3014 highland transparency rightsize simpler cornwells coroners 3011 thornhill imperfectability

## 20 Topics

Topic #0	Topic #1	Topic #2	Topic #3	Topic #4
coroners simpler cornwells friedberg rightsize 3014 sludge bunk naacp 1318 plentiful trona smutty highland sedan	decamped feisal erosions arnolds cynthiana david ash presse heins cumbersomeness privacy adolph technet nice jihadists	lowenberg spews godless memb anticipatory spermatozoa shaun blink berrigan pretensions kasparian airs fates mementos bowls	tricare klobuchar sponsoring kama chaff emporium ilston medium revel paiute gain abridging 1916 glastonbury simpler	simpler coroners naacp cornwells rightsize 1318 friedberg sedan plentiful incinerator firewalls smutty beverage 3014 sludge
Topic #5	Topic #6	Topic #7	Topic #8	Topic #9
simpler cornwells coroners friedberg rightsize naacp europeans sludge plentiful defines highland size 3014 bunk testicular	coroners cornwells simpler naacp friedberg rightsize handbasket plentiful sedan sludge 3014 smutty bunk su 1318	coroners simpler cornwells friedberg 1318 rightsize sedan plentiful naacp sludge smutty bunk 3014 trona incinerator	naacp su rumsfeld rightsize transparency 3014 carrhart 3011 simpler testosterone cornwells coroners highland handbasket exaggerate	coroners simpler cornwells friedberg rightsize 1318 sedan smutty plentiful lovejoy naacp firewalls 3014 bunk sludge
Topic #10	Topic #11	Topic #12	Topic #13	Topic #14
coroners rightsize cornwells simpler friedberg naacp 3014 plentiful smutty bunk rumsfeld sedan sludge trophy sin	hammock madeline harerre appellant agronomic fertile antimatter tias hammondsport batteries connectiveness socked asu tsos resent	friedberg coroners simpler rightsize 1318 cornwells rumsfeld sedan 3014 sludge plentiful naacp smutty bunk trona	simpler coroners cornwells rightsize friedberg naacp madigan 3014 sedan plentiful sludge 1318 rumsfeld bunk smutty	cornwells simpler coroners friedberg plentiful rightsize 3014 1318 naacp smutty sludge carrhart bunk sedan lovelace
Topic #15	Topic #16	Topic #17	Topic #18	Topic #19
coroners friedberg simpler cornwells sedan naacp plentiful smutty 1318 sludge 3014 rightsize incinerator madigan sin	simpler cornwells coroners friedberg 1318 rightsize sedan plentiful sludge incinerator smutty befalls lovelace madigan handbasket	coroners rightsize cornwells simpler friedberg naacp sedan 1318 plentiful smutty bunk handbasket ghim sludge rumsfeld	coroners rightsize cornwells friedberg rightsize naacp 1318 sedan sludge bunk plentiful madigan incinerator humane firewalls	cornwells coroners simpler rightsize friedberg 1318 naacp bunk plentiful smutty befalls madigan testicular sedan sludge

## 30 Topics

Topic #0	Topic #1	Topic #2	Topic #3	Topic #4
coroners friedberg simpler cornwells rightsize naacp plentiful 1318 smutty incinerator madigan sludge sedan bunk rumsfeld	simpler cornwells coroners rightsize friedberg sludge incinerator 1318 plentiful madigan odious trophy handbasket highland naacp	soothes gear frieze odious 1937 bourgeois transitioning hellenic 3011 astronautical letpp nate doster brewers saroyan	su rumsfeld naacp rightsize transparency handbasket preeclampsia fulfills exaggerate testosterone antiaircraft simpler coroners madigan cornwells	cornwells coroners simpler rightsize friedberg 1318 bunk sedan 3014 sludge naacp befalls madigan humane ghim
Topic #5	Topic #6	Topic #7	Topic #8	Topic #9
friedberg simpler coroners cornwells 3014 rightsize sludge naacp plentiful sedan handbasket 1318 carrhart incinerator odious	balkanizing masquerading athleticism hendry raceptives biorefining jsf carriers transformative allstate evacuating gunpowder dave pffs abhors	simpler cornwells friedberg coroners sludge plentiful 3014 sedan rightsize naacp 1318 incinerator smutty highland handbasket	simpler cornwells friedberg coroners sedan plentiful naacp 1318 rightsize incinerator madigan sludge befalls firewalls beverage	friedberg simpler cornwells coroners rightsize naacp 1318 sludge trophy plentiful europeans madigan sedan 3014 smutty
Topic #10	Topic #11	Topic #12	Topic #13	Topic #14
lifting melding hendry trainmen godlessness firewalls backhanded banc bunkhouses beverage heterosexual bars 1842 allstate blockade	coroners friedberg simpler cornwells naacp rightsize plentiful 1318 incinerator sludge sedan handbasket smutty bunk madigan	cornwells coroners friedberg simpler rightsize 3014 sludge naacp sedan smutty plentiful 1318 bunk rumsfeld incinerator	coroners cornwells simpler friedberg rightsize plentiful naacp 1318 sludge madigan trona smutty 3014 firewalls rumsfeld	coroners simpler cornwells friedberg sedan rightsize sludge plentiful 1318 rumsfeld smutty incinerator beverage bunk su
Topic #15	Topic #16	Topic #17	Topic #18	Topic #19
cornwells coroners simpler friedberg rightsize plentiful sedan 1318 naacp sludge incinerator smutty 3014 madigan su	cornwells coroners friedberg 1318 rightsize simpler sludge naacp madigan bunk sedan smutty incinerator plentiful firewalls	simpler cornwells coroners friedberg rightsize naacp bunk 3014 handbasket 1318 plentiful smutty sludge madigan sedan	simpler cornwells coroners rightsize friedberg incinerator 1318 sludge plentiful bunk lovelace naacp beverage carrhart sedan	coroners simpler 1318 friedberg cornwells rightsize naacp plentiful sedan bunk incinerator 3014 beverage sludge befalls
Topic #20	Topic #21	Topic #22	Topic #23	Topic #24
simpler cornwells coroners friedberg 1318 rightsize plentiful sedan sludge incinerator smutty befalls lovelace bunk breastfeeding	coroners cornwells friedberg simpler naacp 1318 3014 incinerator sedan bunk rightsize plentiful rumsfeld simonson odious	cornwells friedberg coroners simpler rightsize 1318 handbasket 3014 naacp plentiful rumsfeld madigan sedan incinerator firewalls	cornwells simpler coroners rightsize 1318 plentiful friedberg naacp incinerator bunk sludge sedan trophy madigan handbasket	simpler friedberg cornwells coroners 1318 plentiful incinerator rightsize naacp sludge beverage carrhart sedan smutty bunk
Topic #25	Topic #26	Topic #27	Topic #28	Topic #29
simpler cornwells coroners friedberg naacp europeans rightsize highland 3014 size sludge defines plentiful testicular bunk	juveniles thinnest kahaukua kyl incredibly 3494 interceptors berrigan herodotus 2671 rewind hirohito elan enabling apportion	authenticity carriers punching interceded metropolitan thermodynamics purportedly extrapolations taped herodotus incredibly majumdar sorenson 3229 324	cornwells friedberg coroners simpler rightsize naacp 1318 sludge trophy sedan plentiful su carrhart ghim handbasket	eyeshade ploesti dressing simpler friedberg cornwells coroners 1318 rightsize plentiful sedan naacp smutty handbasket bunk

## CONCLUSIONS



The analysis of the 110th Congress floor debate dataset using topic modeling techniques has yielded significant insights:

- **Identified Topics:** The analysis successfully identified and categorized topics discussed during the congressional debates, providing a structured view of the key themes and issues addressed.
- **Trends and Priorities:** The identified topics shed light on the legislative priorities and societal trends during the 110th Congress, offering valuable historical context for policymakers and researchers.
- **Public Perception:** By understanding the prevalent topics and discussions, businesses and policymakers can gauge public sentiment and concerns, aiding in decision-making and policy formulation.
- **Strategic Implications:** The insights gained from topic modeling have strategic implications for various sectors, including policy development, public relations, risk management, and content strategy.

In conclusion, leveraging advanced analytics techniques like topic modeling on textual data sets such as the "110" dataset provides a robust framework for understanding complex narratives, uncovering hidden insights, and driving informed decision-making across diverse domains.