Text as data

Section 8

Sima Biondi Spring 2025

Gov 51: Data Analysis and Politics

1 Housekeeping

2 Missing data implementation



• Remember to sign-up for office hours with me or Pranav as a group to discuss projects

- Remember to sign-up for office hours with me or Pranav as a group to discuss projects
- Upcoming deadlines:

- Remember to sign-up for office hours with me or Pranav as a group to discuss projects
- Upcoming deadlines:
 - April 3 (today): Problem Set II due

- Remember to sign-up for office hours with me or Pranav as a group to discuss projects
- Upcoming deadlines:
 - April 3 (today): Problem Set II due
 - April 4 (tomorrow): 1-pager due, see https://naijialiu.github.io/Gov_51/final.html for more information

- Remember to sign-up for office hours with me or Pranav as a group to discuss projects
- Upcoming deadlines:
 - April 3 (today): Problem Set II due
 - April 4 (tomorrow): 1-pager due, see https://naijialiu.github.io/Gov_51/final.html for more information
 - April 11: Initial result due, within one page write up

• 'mice' package - "Multiple Imputation by Chained Equations"

```
library(mice) # recall install.packages("mice") for
    first use
library(NHANES)
data(NHANES)
nhanes <- NHANES[c("Age", "SmokeNow", "TotChol")]</pre>
```

md.pattern(nhanes)





Evaluating imputation



Regressions with imputed data

	(1)
(Intercept)	4.157***
	(0.028)
Age	0.018***
	(0.000)
SmokeNowYes	-0.009
	(0.022)
Num.Obs.	10000
R2	0.141
R2 Adj.	0.141

+ p < 0.1, * p < 0.05, ** p < 0.01, *** p < 0.001

Comparing with listwise deletion

model2 <- lm(TotChol ~ Age + SmokeNow, data = nhanes)</pre>

	(1)
(Intercept)	4.775***
	(0.072)
Age	0.006***
	(0.001)
SmokeNowYes	-0.022
	(0.042)
Num.Obs.	3056
R2	0.010
R2 Adj.	0.009

+ p < 0.1, * p < 0.05, ** p < 0.01, *** p < 0.001

Multiple imputation

• Questions?

• A large goal of this course is to give you a framework to understand different methodologies and the challenges they tackle.

- A large goal of this course is to give you a framework to understand different methodologies and the challenges they tackle.
- Missing data

- A large goal of this course is to give you a framework to understand different methodologies and the challenges they tackle.
- Missing data

- A large goal of this course is to give you a framework to understand different methodologies and the challenges they tackle.
- $\cdot\,$ Missing data \rightarrow types of missingness drives our solutions.
- Text as data

- A large goal of this course is to give you a framework to understand different methodologies and the challenges they tackle.
- $\cdot\,$ Missing data \rightarrow types of missingness drives our solutions.
- Text as data

- A large goal of this course is to give you a framework to understand different methodologies and the challenges they tackle.
- $\cdot\,$ Missing data \rightarrow types of missingness drives our solutions.
- $\cdot \,\, \text{Text as data} \to ?$

• Ban, Grimmer, Kaslovsky, and West (2022) - committee hearings and the effect of women in Congress

- Ban, Grimmer, Kaslovsky, and West (2022) committee hearings and the effect of women in Congress
 - Finds that less interruptions with more women, more substantive conversations.

- Ban, Grimmer, Kaslovsky, and West (2022) committee hearings and the effect of women in Congress
 - Finds that less interruptions with more women, more substantive conversations.
- King, Pan, and Roberts (2013) social media in censored contexts

- Ban, Grimmer, Kaslovsky, and West (2022) committee hearings and the effect of women in Congress
 - Finds that less interruptions with more women, more substantive conversations.
- King, Pan, and Roberts (2013) social media in censored contexts
 - Censorship in China does not target individual criticisms, but instead attempts at collective action.

- Ban, Grimmer, Kaslovsky, and West (2022) committee hearings and the effect of women in Congress
 - Finds that less interruptions with more women, more substantive conversations.
- King, Pan, and Roberts (2013) social media in censored contexts
 - Censorship in China does not target individual criticisms, but instead attempts at collective action.
- Communication relies on or can be represented by text

- Ban, Grimmer, Kaslovsky, and West (2022) committee hearings and the effect of women in Congress
 - Finds that less interruptions with more women, more substantive conversations.
- King, Pan, and Roberts (2013) social media in censored contexts
 - Censorship in China does not target individual criticisms, but instead attempts at collective action.
- Communication relies on or can be represented by text

- Ban, Grimmer, Kaslovsky, and West (2022) committee hearings and the effect of women in Congress
 - Finds that less interruptions with more women, more substantive conversations.
- King, Pan, and Roberts (2013) social media in censored contexts
 - Censorship in China does not target individual criticisms, but instead attempts at collective action.
- Communication relies on or can be represented by text \rightarrow many political applications!

Key question: how do humans come up with sentences?

• One theory is that given a topic (or multiple!), there is a certain probability that words appear

- One theory is that given a topic (or multiple!), there is a certain probability that words appear
- In a given *document*, there could be a number of topics

- One theory is that given a topic (or multiple!), there is a certain probability that words appear
- In a given *document*, there could be a number of topics
- Those topics then dictate the likelihood of which words appear

- One theory is that given a topic (or multiple!), there is a certain probability that words appear
- In a given *document*, there could be a number of topics
- Those topics then dictate the likelihood of which words appear

- One theory is that given a topic (or multiple!), there is a certain probability that words appear
- In a given *document*, there could be a number of topics
- Those topics then dictate the likelihood of which words appear



Figure 1: Opposite of this!

Takeaways of the DGP for text data:

• Pretty rigid framework, but some ground truth

Takeaways of the DGP for text data:

- Pretty rigid framework, but some ground truth
- Also undergirds the logic under Chat GPT and other large language models (LLMs)

Takeaways of the DGP for text data:

- Pretty rigid framework, but some ground truth
- Also undergirds the logic under Chat GPT and other large language models (LLMs)
- Why GPTZero and other programs to detect AI usage are somewhat easily able to detect because text generated using this framework is rigid!

Where do we begin analysis on text data?

• Text is just a collection of words - the order and structure do not matter particularly when we care about topical relevancy

- Text is just a collection of words the order and structure do not matter particularly when we care about topical relevancy
- Method **assumes** that the frequency of words can provide us information about the context in the text

- Text is just a collection of words the order and structure do not matter particularly when we care about topical relevancy
- Method **assumes** that the frequency of words can provide us information about the context in the text

- Text is just a collection of words the order and structure do not matter particularly when we care about topical relevancy
- Method **assumes** that the frequency of words can provide us information about the context in the text

Process:

1. Tokenization - dividing text into individual words

- Text is just a collection of words the order and structure do not matter particularly when we care about topical relevancy
- Method **assumes** that the frequency of words can provide us information about the context in the text

Process:

- 1. Tokenization dividing text into individual words
- 2. Counting counting the frequency they show up

- Text is just a collection of words the order and structure do not matter particularly when we care about topical relevancy
- Method **assumes** that the frequency of words can provide us information about the context in the text

Process:

- 1. Tokenization dividing text into individual words
- 2. Counting counting the frequency they show up
- 3. Vectorization representing the text as a vector of word frequencies

Example: Congressional hearings about TikTok



"Tiktok is fun" - Sima Biondi (Not in Congress-CA)

"Tiktok is fun" - Sima Biondi (Not in Congress-CA)

Process:

1. Tokenization: {TikTok, network, fun, etc.}

"Tiktok is fun" - Sima Biondi (Not in Congress-CA)

Process:

- 1. Tokenization: {TikTok, network, fun, etc.}
- 2. Counting frequency of text appearances: {2, 1, 1, etc.}

"Tiktok is fun" - Sima Biondi (Not in Congress-CA)

Process:

- 1. Tokenization: {TikTok, network, fun, etc.}
- 2. Counting frequency of text appearances: {2, 1, 1, etc.}
- 3. Vectorization \rightarrow

i	TikTok	network	fun
Richard Hudson	1	1	0
Sima Biondi	0	0	1

i	Does	TikTok	access	the	home	Wi-Fi	network	is	fun	China	
1	1	1	1	1	1	1	1	0	0	0	
2	0	1	0	0	0	0	0	1	1	0	

• Topic modeling

i	Does	TikTok	access	the	home	Wi-Fi	network	is	fun	China	
1	1	1	1	1	1	1	1	0	0	0	
2	0	1	0	0	0	0	0	1	1	0	

- Topic modeling
- Sentiment analysis

i	Does	TikTok	access	the	home	Wi-Fi	network	is	fun	China	
1	1	1	1	1	1	1	1	0	0	0	
2	0	1	0	0	0	0	0	1	1	0	

- Topic modeling
- Sentiment analysis
- Text classification

i	Does	TikTok	access	the	home	Wi-Fi	network	is	fun	China	
1	1	1	1	1	1	1	1	0	0	0	
2	0	1	0	0	0	0	0	1	1	0	

- Topic modeling
- Sentiment analysis
- Text classification
- Also a great application of Lasso, since data can get VERY large.

 \cdot Can be supervised or unsupervised

- Can be supervised or unsupervised
 - Supervised: use labeled data to guide model in identifying topics

- Can be supervised or unsupervised
 - Supervised: use labeled data to guide model in identifying topics
 - Unsupervised: no labeled data, based on distribution of words in document

- Can be supervised or unsupervised
 - Supervised: use labeled data to guide model in identifying topics
 - Unsupervised: no labeled data, based on distribution of words in document
- Oftentimes, will require some human interpretation of the generated categories

- Can be supervised or unsupervised
 - Supervised: use labeled data to guide model in identifying topics
 - Unsupervised: no labeled data, based on distribution of words in document
- Oftentimes, will require some human interpretation of the generated categories
- Example: Terman (2017)

- Can be supervised or unsupervised
 - Supervised: use labeled data to guide model in identifying topics
 - Unsupervised: no labeled data, based on distribution of words in document
- Oftentimes, will require some human interpretation of the generated categories
- Example: Terman (2017)
 - Examines portrayals of Muslim woman in American media

- Can be supervised or unsupervised
 - Supervised: use labeled data to guide model in identifying topics
 - Unsupervised: no labeled data, based on distribution of words in document
- Oftentimes, will require some human interpretation of the generated categories
- Example: Terman (2017)
 - Examines portrayals of Muslim woman in American media
 - Results suggest that US news media propagate the perception that Muslims are distinctly sexist

• Multiple imputation is easily implemented through the 'mice' package.

- Multiple imputation is easily implemented through the 'mice' package.
- Introduces a number of easy to use functions that help with descriptive statistics and modeling.

- Multiple imputation is easily implemented through the 'mice' package.
- Introduces a number of easy to use functions that help with descriptive statistics and modeling.
- Text as data is incredibly popular and powerful tools in social science.

- Multiple imputation is easily implemented through the 'mice' package.
- Introduces a number of easy to use functions that help with descriptive statistics and modeling.
- Text as data is incredibly popular and powerful tools in social science.
- Data generating process of text is an important foundation for understanding how to tackle text as data.

- Multiple imputation is easily implemented through the 'mice' package.
- Introduces a number of easy to use functions that help with descriptive statistics and modeling.
- Text as data is incredibly popular and powerful tools in social science.
- Data generating process of text is an important foundation for understanding how to tackle text as data.
- Bag-of-words is a simple, but powerful model to analyze text.