

# Data Expedition Proposal

*Becky Tang and Graham Tierney*

*8/5/2019*

1. Sponsoring faculty members: Fan Li (STA 440) and Maria Tackett (STA 199)
2. Title of dataset(s)
  - Democratic Debate Transcripts, June 2019
  - Democratic Candidate Tweets
  - State of the Union Transcripts, 1970-present
3. Description: The following datasets were chosen for their quality of being open to different kinds of analyses which explore real and relevant social science questions.
  - Democratic Debate Transcripts, June 2019
    - Each row in the dem\_debates.csv file represents a candidate’s continuous block of speech before being interrupted by another candidate, with variables specifying the night, placement, speaker, and the text of each block. We also provide text files of the raw transcripts from both debates.
    - Source(s)
    - Why this dataset?
    - How dataset was put together
    - The dem\_debates.csv file has 1136 observations of 4 variables.
  - Democratic Candidate Tweets
    - The dem\_cand\_tweets.csv file contains tweets from four major Democratic presidential candidates: Joe Biden, Kamala Harris, Bernie Sanders, and Elizabeth Warren. Other variables include the timestamp of each tweet, if/which hashtag was used, and how many retweets and favorites a tweet garnered (at the time of scraping).
    - Similar to the Democratic Debate Transcripts data, the twitter data can be used to compare and contrast the four candidates based on the language of their tweets. A benefit to this dataset is that each tweet is a complete body of text, allowing for easier determination or classification of the sentiment in a tweet as compared to in a book or speech.
    - The data were scraped via the rtweet package in R.
    - The dem\_cand\_tweets.csv file has 4799 observations of 90 variables.
  - State of the Union Transcripts, 1970-present
    - The data associated with this project are transcripts from State of the Union addresses beginning with President Nixon’s 1970 address. Also included is an Excel file ‘stateofunion.xlsx’ specifying the year, president, and the president’s political party corresponding to each address.
    - A major theme presented in *The New Jim Crow* is that mass incarceration in America began when President Nixon coined the term “War on Drugs”. We can use the State of the Union addresses to investigate how different presidents have championed this issue.
    - The data consist of 50 transcripts saved as text files which vary in length/size.
    - The transcripts were obtained from the website The American Presidency Project. We simply copy-and-pasted each transcript into a .txt file, although in the future we could use a web scraper.
4. Potential classroom exercises: We can ask questions at the following levels of granularity
  - Document specific (ex. What are the major themes/topics presented in Document Z? What is the overall sentiment associated with Document Z?)
  - Author/speaker specific (ex. How does Person X’s language change over time? How does Person X feel about Issue Y?)
  - Comparisons across documents and authors (ex. How often does Person X speak about Issue Z in comparison to Person Y? How similar is the language used by different authors or in different documents? Given text from Document Z, can we determine who the author is?)

5. Techniques: All analysis will be performed in R, and the installation of several R packages will be necessary. Therefore, pre-loading packages with a virtual environment or RStudio in the cloud would be helpful.

- Sentiment Analysis
  - Sentiment analysis provides a simple way to reduce the dimensionality of text data. Certain words are mapped to a given (short) list of sentiments. Each document is then scored by the number or proportion of words of the given sentiment. Some sentiment dictionaries map words to a continuous positive or negative scale, and then a document's sentiment is the average or total of the score.
  - These scores can be easily visualized in histograms and density plots for different categories of documents or across different levels of some document feature (such as time of creation).
  - As EDA, we can apply PCA using the sentiment scores to observe any clustering of documents.
  - Logistic regression can be used to classify documents based on their sentiment scores, with k-fold cross-validation to assess accuracy. In the case of State of the Union addresses, students can try to predict the party of the author of a given paragraph via the sentiment information. They can also evaluate which Presidents are most distinctive, by comparing the accuracy rate of their logistic regression predictions for each President.
- Multiple testing for word proportions
  - A method to identify words that are used “significantly” differently across authors.
  - If documents in category A use word  $i$  at frequency  $p_{iA}$  and documents in category B use word  $i$  at frequency  $p_{iB}$ , a simple two-sample t-test can be used to test the null hypothesis that  $p_{iA} = p_{iB}$ .
  - However, since the number of unique words in the corpus is quite large, multiple testing corrections are necessary. The Bonferroni correction is useful if one is trying to assess whether there is any difference in word use between the categories. The Benjamini and Hochberg (1995) correction is more useful to identify the subset of words that differentiate the two categories. This method was used in Airolidi et al. (2006).
- Bag-of-words generative models
  - The general assumption of these approaches is that authors pick words according to some unknown set of frequencies. The likelihood of an observed word is thus proportional to the unknown frequency. The order of words is generally ignored in these models.
  - With a Dirichlet prior on the frequency vector, the posterior can be easily computed via conjugacy. Students with more statistics background can be assigned to compute these posteriors and their properties as homework problems. For instance, the probability that a text with unknown authorship comes from a certain author is proportional to the prior probability times the posterior-predictive likelihood for that author. This probability can be used to measure classification accuracy and quantification of uncertainty.
  - This method generally requires labeled data from two or more authors. Examples are found in Airolidi et al. (2006) and Gentzcow, Shapiro, and Taddy (2019).
- Topic modeling with LDA (advanced)
  - LDA is an unsupervised generative model that does not require labeled data. A “topic” is modeled by a frequency vector over all possible words. If a word  $W$  comes from topic  $t$ , then  $P(W = w_i) = \beta_{ti}$ .
  - For a fixed number of topics, a corpus of documents can be generated by the following process. Choose the length of the document (number of words) via  $N \sim \text{Pois}(\lambda)$ . Choose each document's topic frequency by  $\theta_d \sim \text{Dir}(\alpha)$ . For each word  $w_i$  in the document  $d$ , draw  $z_{di}$   $\text{Categorical}(\theta_d)$ , then draw word  $w_i \sim \text{Categorical}(\beta_{z_i})$ . Both full and collapsed Gibbs samplers can be programmed by students who have taken STA 360. All results follow from conjugacy.
  - After estimation, one can inspect the topic-word frequencies to see what concepts they relate to and the document-topic frequencies to see which document come mostly from each topic. Documents can be classified into the topic their words most frequently come from.

6. Source(s)

- The American Presidency Project (<https://www.presidency.ucsb.edu/>)