

Food for Thought: Analyzing Public Opinion on the Supplemental Nutrition Assistance Program

Miriam Chappelka 1, Jihwan Oh 2 Dorris Scott 3, and Mizzani Walker-Holmes 4

College of Computing, Georgia Institute of Technology

Introduction

The Supplemental Nutrition Assistance Program (SNAP), formerly known as food stamps, is a federal program that helps low income individuals purchase food. The Atlanta Community Food Bank (ACFB) aspires to eliminate hunger in its service area by 2025, and one strategy the food bank is using to achieve this goal is to raise awareness about the importance of SNAP. Their audience is stakeholders who contribute to the Atlanta Community Food Bank (who may be skeptical of the food bank's support of SNAP) and politicians (who can influence SNAP policy). We are assisting the food bank by analyzing public opinion of SNAP on social media and news outlets, as well as tracking Georgia politicians' voting records on issues relating to food insecurity. This project focuses on utilizing natural learning processing tools, sentiment analysis, machine learning, and text mining to capture public opinion on the Supplemental Nutrition Assistance Program on a national and state level. One objective of this project is to explore how discourse regarding SNAP varies geographically. While the ACFB has hypotheses based on their experiences, they do not have any quantitative measures to support their conjectures as of yet. After analyzing the sentiment of the data gathered from social media and news outlets, geospatial analytics will be used to identify geographic variation in SNAP sentiment. In addition to better understanding public opinion on SNAP, the ACFB is also interested in the voting records of Georgia politicians in Congress and in the Georgia General assembly. Having easy access to representatives' voting records on bills regarding food insecurity will help the food bank prepare for policy meetings with these politicians. Ultimately, this research will produce a tool that communicates dominant narratives and opinions about SNAP so that the ACFB Advocacy team can better communicate to stakeholders about SNAP. This research is being conducted in conjunction with the Atlanta Community Food Bank and the Data Science for Social Good program at the Georgia Institute of Technology.

Methods

Data Collection

A variety of data sources were collected for this study. News articles were collected using the webhose.io service. Webhose.io is a commercial website which allows for easy web scraping of articles through keyword searches. 2,239 news articles about SNAP, published between May 14, 2017 and July 13, 2017, were scraped using this service. The scraped articles were geo-coded based on the location of the news outlet, and classified based on the type of news outlet. Out of 2,239 articles scraped through Webhose.io, duplicates were removed which reduced the article number down to 1,081 articles. Tweets were collected for a one-month period using the streamR package in R, which accesses the Twitter Streaming API. The Streaming API allows access to around one percent of tweets that are being tweeted in real time. The collection of the tweets was based on search terms related to SNAP: "SNAP," "food stamp," "food stamps," and "EBT." The tweets were selected if they had any meaningful content regarding SNAP and were further sorted based on if they were geo-tagged. There were approximately 700 tweets about food stamps that were used for this analysis. Finally, the voting records of Georgia state representatives were collected through Open States, a site that collects data on state representatives. Bills were selected if they contained the phrases "food stamps", "SNAP", "food bank", "food desert," "hunger," "food insecurity," or "georgia peach card". Bills with no votes were removed, and votes by representatives' no longer in office were removed.

Text Mining and Sentiment Analysis

Sentiment analysis was used to assess the discourse regarding SNAP. Sentiment analysis is a form of text analysis that determines the subjectivity, polarity (positive or negative) and polarity strength (weakly positive, mildly positive, strongly positive, etc.) of a text. In other words, sentiment analysis tries to gauge the tone of the writer. There are two main approaches in classifying the sentiment of a given text: supervised classification and unsupervised classification. Supervised classification requires

2018 ASEE Southeastern Section Conference

labeled data and its features must be extracted from the data. Examples of features are part of speech tags, most frequent words, reading level, and name entity tags. Labels are nominal data. With these features and labeled data, any type of supervised learning approach can be used. It creates a model that is suitable for the data set with the label, so that it can predict with a new dataset without the label. This model is totally dependent on the dataset and its characteristics. When the characteristics in the dataset are similar, supervised learning classification tends to perform well. This applies for the Twitter data set, where the length and diction of the tweets are similar to one another. For the Twitter data, the scikit-learn package from Python was used to perform supervised classification. Unsupervised classification was performed on the news articles. Unsupervised classification is different from supervised learning where the model is independent from the data, but it follows specific rules that it has in place. In this case, it uses a pre-existing lexicon, a dictionary that contains more information than just its meaning, and syntactic data, set of rules regarding the syntax of the sentence structure, to determine its sentiment. This method creates a numerical value or a probability of the sentiment rather than a nominal classification. This form of classification was used to analyze the news articles because the text has varying length, style, dictions, and form depending on the writer, which requires a bigger dataset to perform supervised classification. The Vader and AFINN packages in Python were used to conduct unsupervised sentiment analysis. Vader is short for Valence Aware Dictionary Sentiment Reasoner, and is a lexicon and rule-based sentiment analysis tool. AFINN is a dictionary of words that rates connotation severity from -5 to 5. The actual sentiment score was given as the sum of the word score within a sentence. The Vader tool gauges the overall syntactical sentiment more so than the word usage. Conversely, AFINN gauges the type of words that are being used and their intensity. Additionally, sentences with key words (words relating to SNAP) were given a higher weight so that sentiment towards this issue would be amplified. Each article was tokenized to the sentence level, and each sentence was given a sentiment score according to the two sentiment analysis tools {NLTK}. Then, the scores were aggregated for each article with the weight that was assigned to each sentence. This aggregated score represents the sentiment of the article. To take into account of impact of the article, each article was then aggregated in regards to the traffic level of the website and the reading level of the article.

Sentiment Analysis Methods

Additionally, information on the arguments and topics in these articles would be very useful to the ACFB. To do this, preliminary topic modeling (Latent Dirichlet Allocation) has been performed to extract the topical words from the set of text. It returns a set of words with probabilistic weight on each of the word to indicate its importance. Bigram collocation has been used to detect sets of two words that are most frequent and meaningful. Term frequency inverse document frequency (TFIDF) was used to detect important words across all the documents. Name Entity Recognition (NER) from the Stanford Natural Language Processing Group and gensim were used to detected key people or locations mentioned in the articles. After generating all the statistics, each word within TFIDF, bigram collocation and NER was multiplied with the weight that was computed with each of the documents. Then, all the words were aggregated into a list. Using this list, a word cloud can be generated to visualize meaningful words. Word clouds are especially of interest to our partners at the food bank. Along with the word cloud, its aggregation by each date will help the viewer understand the subject of the sentiment to better decipher the public opinion about SNAP.

Spatial Analysis

Based on the hot spot analysis that was conducted on the AFINN sentiment scores of 1,250 of the 2,239 news outlets, the news outlets with negative AFINN scores were more concentrated compared to the news outlets that had positive AFINN scores. Many of the news outlets that have a negative sentiment on SNAP were in the Midwest, especially in Indiana, Michigan, and Illinois. On the other hand, news outlets with positive AFINN scores were more dispersed, with a concentration of positive AFINN scores in the South and Southeast. This could be due to the high enrollment of individuals on SNAP such as in the District of Columbia, Mississippi, and Tennessee which have the highest number of individuals on SNAP in the nation.

Deliverables

The results of the sentiment analysis, text mining, and aggregation of voting records will be contained in an on-line application which was created using the Shiny web framework in R. This application will

2018 ASEE Southeastern Section Conference

allow the ACFB to better understand reporting and public opinion on SNAP through interactive visualizations such as word clouds, maps, charts, and graphs. The “Background” section of the application gives an overview of the SNAP program and the importance of the program in various contexts. “The Word on SNAP” section will provide visualizations of how SNAP is discussed in social media and media outlets, such as the interactive word cloud that is displayed in Figure 4. This section also includes an interactive map called “SNAP InfoMap” in which users can see the location and types of news outlets reporting on SNAP and the affiliated sentiment score attached to each outlet. Users are also able to explore how the location of the news outlets correlates to the socioeconomic characteristics that are related to the program such as the percentage of households that are on SNAP (see Figure 6). In addition to the word cloud and the interactive map, a sentiment analysis tool was created to show the average AFINN and Vader scores for the news outlets and tweets and how the sentiment on SNAP changes through a specified time period. For example, when President Trump announced a budget cut on SNAP, most of the articles for higher trafficked websites had a negative sentiment score. The “Politician Tracking on SNAP” section will allow one to look up the voting record of Georgia legislators on bills related to SNAP on the state level, as shown in figure 5. The word cloud uses TFIDF in order to show which words are prominent in a set of articles which is related to the size of the word in the visualization.