

Identifying Political Bias in the Media

Introduction

The media is a service industry and has a responsibility to be factual and without bias. Every news-writer and news agency is bound by a certain set of ethics and unwritten laws that are based on this responsibility.

The first and foremost is the responsibility of accuracy in reporting. The facts have to be precise, specific, and clear. Reporting is about stating facts without the injection of opinions.

Although these rules and regulations are well understood, the deviation from them is sometimes very easily evident, i.e. inclination towards and away from some group can be noticed. This is particularly true when it comes to politics.

This is where we (the team) would like to assess how the major news providers are performing, i.e. we would like to see if the major news providers of the country are biased towards or away from certain political figures.

Data Selection

News Providers

The selection of news providers were based on their popularity and the number of unique visitors the news provider receives on their website. The ranking was based on each news providers Alexa Rank. The list was altered as drawing data from some of the news providers was infeasible. The final list is as follows.

1. New York Times
2. Washington Post
3. Fox News
4. USA Today and
5. New Yorker

Political Figures

The political figures picked were based on the popularity of individual figures based off of polls conducted by <https://www.crowdpac.com/>. Two from each of the leading political parties were picked and for an added perspective, the two political parties were included as well.

The final list of the political figures are as follows.

1. Hillary Clinton
2. Bernie Sanders
3. Democratic Party
4. Donald Trump
5. Ben Carson and
6. Republican Party

Data Collection

Data collection involved scraping data manually from individual pages of news articles as a single repository of the data was not available. The cost of simply purchasing the data was not feasible. Scraping a web page is a method by which single elements of data are extracted from the webpage. The single elements are identified by their html tags.

A list of all URLs to be scraped was created using a search engine optimization tool known as SEOquake. The tool interacted with the Google search engine to return relevant results, i.e. results having a certain key word from specific websites within a given time period.

The job of extracting the specific text body from all the URLs was done using the tool Kimono Labs. Once specific web elements were selected, Kimono Labs was capable of creating a template and applying it across multiple webpages to return a flat file containing the extracted data for a chosen political figure from a specific news provider. 30 such flat files were generated; each containing a certain combination of the political figure and news agency.

Data Cleaning

The flat file containing the extracted data contained the following attributes.

1. Date of the article.
2. Text of the article
3. Index (Row identifier)
4. URL - URL of the article

Microsoft Excel

Each article had several rows of data (each paragraph in the article is assigned a different row). The different rows of data for a single article were concatenated into one single cell. This was done using the MoreFunc add-in in Microsoft Excel.

Another Microsoft Excel add-in "PowerQuery" was used to collate the 30 different files containing data from the different news providers and for different political figures.

R

R was used for finer cleansing of the data. Abbreviated words and irrelevant text (such as "Mr.", "Mrs.", Quotes etc.) were removed from all the article text using R.

Analyzing the data

Sentiment analysis was conducted on the text to identify bias. A positive or negative sentiment on a certain political figure does not indicate bias but a consistent sentiment can be used to identify bias. Another way to identify bias is to spot consistency in associating the political figure with a certain sentiment irrespective of the sentiment of the article itself.

Application of the NLP to build sentiment scores.

The following steps describes how the model is designed to function:

Step 1: Read data from the flat csv file that contains the news article, the date and URL.

Step 2: Break the article down to individual sentences and words.

Step 3: Use the AFINN lexicon to calculate the mean of all the sentiments of the words from the article. This is the text score of the article. (AFINN is a dictionary of approximately 2500 words that contain the strength of the positive and negative sentiment associated with words)

Step 4: Use the AFINN lexicon to calculate the mean of all the sentiments of the words from the sentences that mention the political figure. This is the sentence score of the article.

Step 5: Calculate the intensity of the conflict in the sentence score and the text score for every article. This is the net score. This measures the intensity of bias towards the mentioned political figure in the article.

```
10 nlzr <- function(file_name_vec,key_word){
11   fil <- read.csv(file_name_vec)
12   txt <- fil$Text
13   exwords <- c("Mr.", "Ms.", "Mrs.", "etc", "\n")
14   txt <- gsub(x=txt, pattern = paste(exwords, collapse = "|"), replacement = "")
15   txt <- tolower(txt)
16   sentence <- strsplit(txt, split=".", fixed=TRUE)
17   word <- strsplit(txt, split=" ", fixed=TRUE)
18
19   cname <- c("words", "value")
20   afinn <- read.delim2("AFINN-111.txt", header=FALSE, col.names = cname)
21
22   for(i in 1:length(txt))
23   {
24     word[[i]] <- gsub("[[:punct:]]", "", word[[i]])
25
26     sentence[[i]] <- gsub("[[:punct:]]", " ", sentence[[i]])
27
28     inds <- which(afinn$words %in% word[[i]])
29     ifelse(length(inds) == 0, txt_score <- 0, txt_score <- mean(afinn$value[inds]))
30     fil$Text_Score[i] <- txt_score
31
32     pos <- grep(key_word, sentence[[i]])
33     hold_pos <- sentence[[i]][pos]
34     sent_score <- c(0)
35     if(length(pos) != 0){
36       for (j in 1:length(hold_pos)){
37         word_from_sent <- strsplit(hold_pos[j], split=" ", fixed=TRUE)
38         inds <- which(afinn$word %in% word_from_sent[[length(word_from_sent)]])
39
40         if(length(inds)>0){
41           sent_score <- sent_score + mean(afinn$value[inds])
42         }
43       }
44     }
45
46     if(length(hold_pos) == 0){
47       fil$Sentence_Score[i] <- 0
48     } else{
49       sent_score <- sent_score/(length(hold_pos))
50     }
51     fil$Sentence_Score[i] <- sent_score
52
53     net_score <- c(0)
54     ifelse(sign(txt_score) == sign(sent_score),
55           net_score <- sign(sent_score)*abs(abs(txt_score)-abs(sent_score)),
56           net_score <- sign(sent_score)*abs(abs(txt_score)+abs(sent_score)))
57
58     fil$Net_Score[i] <- net_score
59   }
60   ret <- fil[,c(1,6,7,8,4,5)]
61   return(ret)
```

Conclusion

A stagnant or single sided text score indicates bias. The same applies to the sentence score. A very high net score (above a chosen threshold) indicates bias. This score needn't be consistent to be an accurate indicator of bias.

The complete R Code is attached separately.

Visualization of the result

The GoogleVis package in R is used to represent the relevant information. The package served as a convenient tool to visualize multi dimensional data. A set of sample visualization images are attached as below:



