

In light of the upcoming 2020 election and becoming more familiar with the ever-expanding field of Natural Language Processing, we wanted to take a deeper look into how our political leanings affect the types of words we use.

Emily Laber-Warren’s article *Unconscious Reactions Separate Liberals and Conservatives* suggests that conservatives are more anxious than liberals because they typically favour traditional ideals, stability, clear answers, and structure; conservatives tend to be less accepting of change. In *Fear and Anxiety Drive Conservatives’ Political Attitudes*, Bobby Azarian outlines how conservatives are more focused on negative situations, fear new experiences, and are more reactive to fear than liberals, a claim buttressed by medical data from University College London. Scientists found that self-identifying conservative students have larger amygdalas than self-identifying liberal students. The amygdala, which is involved in emotion processing, is especially reactive to fearful stimuli. An oversized amygdala could intensify one’s sensitivity to a potential threat. As Carrie Sheffield explains in *Are Political Conservatives Inherently Fearful and Angry?*, conservatives are most moved by jarring realities and more determined to act to prevent further harm.

DATASET

After reading these articles, we decided to focus on liberals’ and conservatives’ use of negative words. We found a political news dataset containing 87,157 English news articles, from which we randomly generated a sample of 5,000 articles. The dataset included the text and news source of each article. Our sample contained articles from 10 news sources, and we determined how conservative or liberal they were using a media bias chart.

```
Channel: "H",
"url": "http://www.cbs4.com/story/3024907/southern-company-supports-operation-migration-annual-journey-for-eighth-conservative-year",
"category": "H09",
"title": "Southern Company supports Operation Migration's annual journey for eighth consecutive year",
"performance_score": 0,
"site": "cbs4.com",
"participants_count": 1,
"title_full": "Southern Company supports Operation Migration's annual journey for eighth consecutive year - CBS News 4 - San Diego, CA News Station - CBS Channel 4",
"span_score": 0,
"site_type": "news",
"published": "2015-10-12T03:00:00-00:00+03:00",
"media_count": 0
```

FIGURE 1: SAMPLE FORMAT OF DATASET

HYPOTHESIS

Articles in our exploratory research brought forward the notion that conservatives tend to be more negative, hence we wanted to test:

★ H0: Conservative news sources have the same proportion of articles with negative words such as “fail”, “hate”, “critic”, and “scare” as liberal news sources.

★ HA: Conservative news sources have a greater proportion of articles with negative words such as “fail”, “hate”, “critic”, and “scare” than liberal news sources.

We chose the words “fail”, “hate”, “critic”, and “scare” based on selecting words with a negative sentiment from the 5,000 most frequently occurring words in our corpus. We had more words in consideration, but discarded words that appear very infrequently in our dataset after performing some preliminary exploration. Words that occur multiple times in an article will only be counted once since we are focusing on the proportion of articles rather than the total number of occurrences of each word, which would overrepresent lengthier articles.

METHODOLOGY

We first put all of the data into a Pandas dataframe so that it is easier to work with. We assigned each article a bias_score according to its news source. For example, a NYT article has a bias_scorer of -1, Fox News is 1, Huffington Post is -2, etc. according to a news-bias datasource with which we joined the article dataset. We then stemmed and tokenized each article’s text. We added a binary variable for each of the test words (eg. containsHATE).

NOTE: A higher bias score does not necessarily indicate more bias. A higher np.abs(bias_score) indicates more bias, and bias < 0 = liberal, bias > 0 = conservative.

```
bias_dict = { 'isNYT': -1,
              'isFOX': 1,
              'isCNN': -1,
              'isWSJ': 1,
              'isCBS': -0.5,
              'isNYMAG': -1,
              'isUSNEWS': -0.5,
              'ishuffpo': -2,
              'dailycaller': 2,
              'isCBN': 2,
              'isREUTERS': 0 }
```

FIGURE 2: BIAS SCORES OF LIBERAL AND CONSERVATIVE NEWS SOURCES

LIBERAL VS.

SENTIMENT HYPOTHESIS TESTING

Conservative

RESULTS

The OLS results indicate that the data do not support our alternative hypothesis. The negative value of the const term in Figure 6 indicates there are more liberal than conservative articles in the dataset overall. Each independent variable is a binary indicator for whether a given word appears in an article. The R-squared is 0.029, which is a bit low even for a social science experiment, where we know that 0.1 is a common benchmark. This is not too concerning on its own however, because R-squared is more a commentary on tightness than correlation. The F-stat is quite strong, so there does seem to be a nonzero correlation. For all individual variables, the t-stats are all very strong with the exception of ‘president’, our control variable, and ‘scare’ which is still significant to a ~92% confidence interval. What was surprising however, are the findings themselves.

Regarding the control, we hypothesized that ‘president’ is a word that will appear with very minimal bias in a political dataset, and that is indeed the case. The word “president” has a t-stat of -0.342, which implies no correlation, and the effect size is less than 50% the magnitude of any other variable that we regressed with. The addition of the word ‘president’ did not change any of the other results, so it is not multi-collinear or correlated with any of the actual test variables, and we think it is useful to include as an easily interpretable spot check that there is nothing severely wrong with the data. Regarding the effect magnitude and direction, the direction of the effects are exactly opposite to what our alternative hypotheses predicted. The negative coefficients for our test words mean that the presence of these words are correlated with liberal articles.

Given that the effect direction is the same for all of our variables, there is a temptation to ask if the data support an alternative alternative hypothesis. First of all, this is very dangerous as a practice, because fitting a hypothesis to results is not really an honest execution of the scientific method. That fact aside, there is the question of effect sizes. The coefficients for the variables hover right around -0.3. However, the domain of the dependent variable bias_score ranges from -2.0 to 2.0 inclusive. This means that an effect size of magnitude 0.3 is not a particularly compelling value. It is certainly far from no effect, but the implications of such a low magnitude score are not irreproachable.

The t-statistics show that the independent variables (with the exception of the control, and 1 other word depending on significance threshold) have a strong correlation with the dependent variable. The low R-squared value shows that the data are not very tight around this line. The direction of the effects are opposite to the direction predicted by the alternative hypothesis, but the effect sizes are not so great that they immediately suggest we adopt and test the alternative, alternative hypothesis. Some of the challenges we faced were selecting our initial words to make sure they were representative of words with a negative sentiment in that they occurred often. We also had to take the time to decide which statistical test would be “correct” science.

A note on interpretation: A coefficient for variable VAR of value VAL means that if an article contains VAR, the bias score is const + VAL. More specifically, the fact that ‘hate,’ ‘fail,’ ‘scare,’ and ‘critic’ all have negative coefficients means that if any of those words are in an article, it is predicted to have a negative bias score (indicating it is liberal).

TESTING

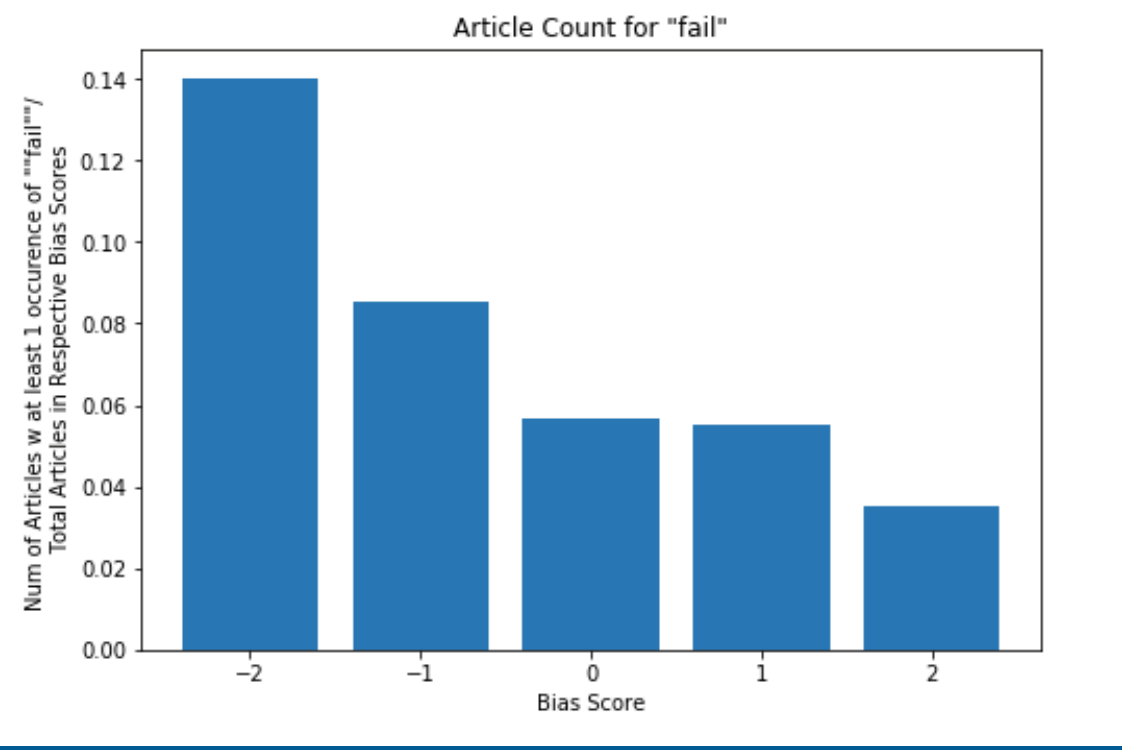
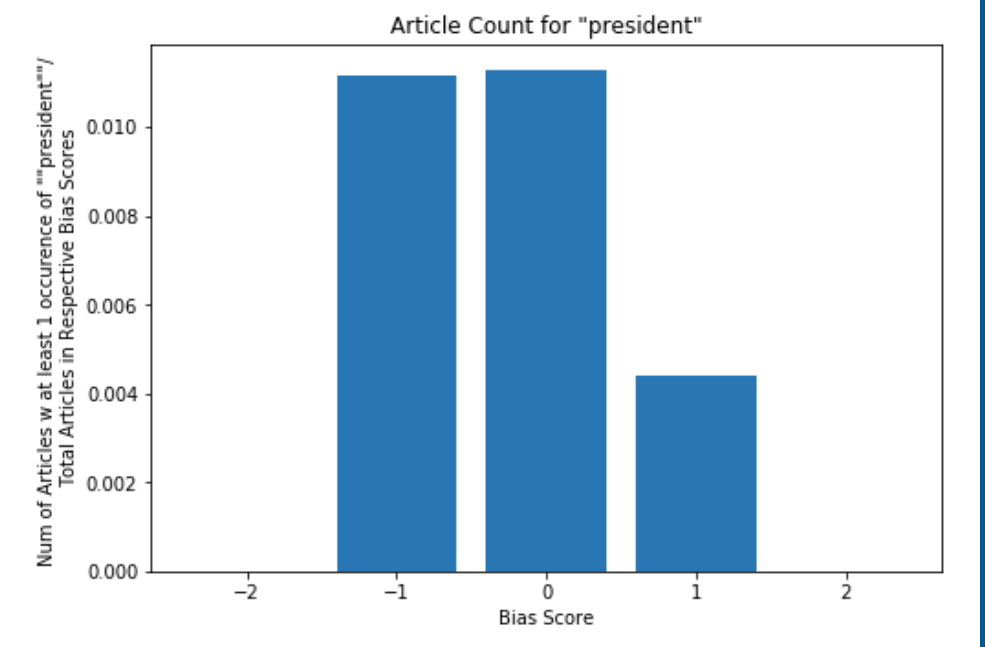
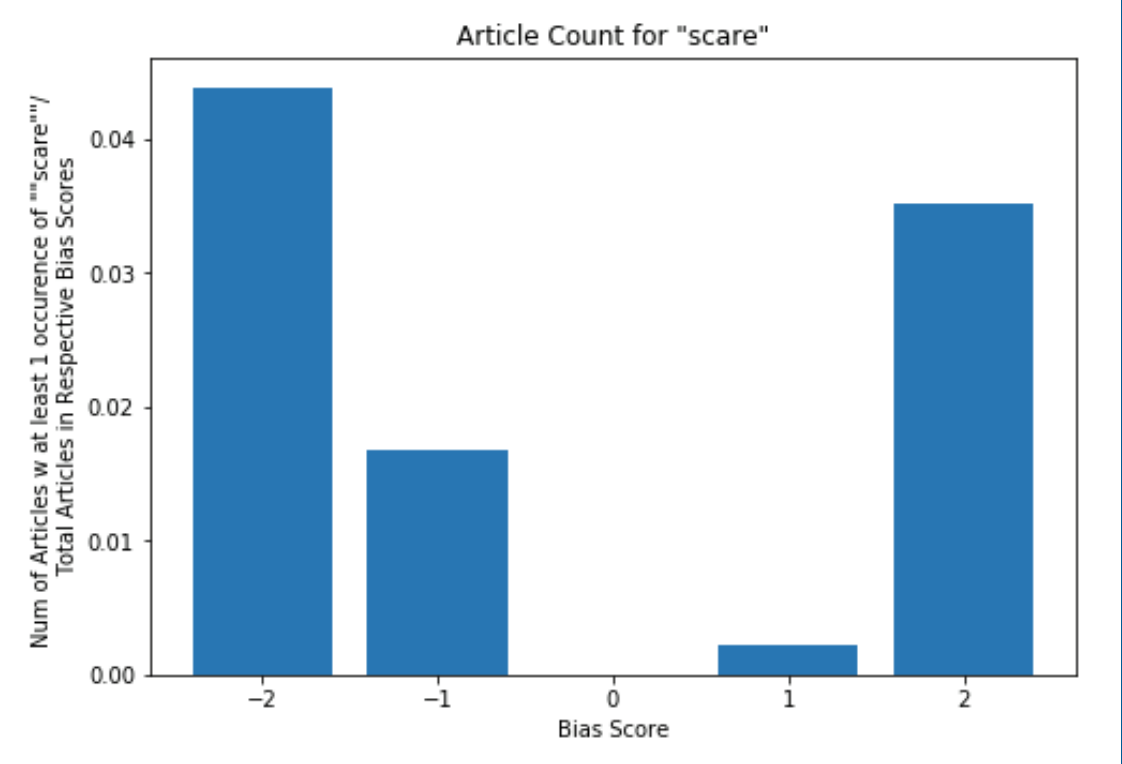
We performed a multivariable regression using the StatsModels OLS package. One of the many reasons that we chose this type of analysis is that there seemed to be more liberal articles in our dataset than conservative articles, and rather than “throw” them out, we can account for this bias in the data with the constant term in the OLS regression. Also, while we do not want to count multiple instances of the word “hate” in an article, for example, we do think that the existence of multiple target words suggests a stronger indication of article tone. Multiple regression is a method of testing multiple variables, as the alternative of performing multiple single regressions creates a whole slew of new problems such as omitted variable bias.

We also performed a sanity check test using the NLTK Vader SentimentIntensityAnalyzer. This involves identifying the net positivity (positive sentiment - negative sentiment) of each article, multiplying this term by the article’s bias score, and taking the mean of all liberal articles, all conservative articles, and all articles for this score.

NOTE:
1: This is a sanity check and is not meant to test the hypothesis directly.
2: Negativity is not a direct proxy for our exact hypothesis.

That being said, one might expect that net_positivity would be negatively correlated with anger and fear, and so if conservative articles (which we give a bias_score>0) are angrier, then the average net_positivity * bias_score across all articles would be negative. For example, if a NYT article has a net_positivity of 0.1, then net_positivity * bias_score = -0.1. If Fox has an article with a net_positivity of -0.1, then net_positivity * bias_score = -0.1.

If our hypothesis was supported and anger was correlated with negativity, then we expect a negative number - however, this was not the case. Across the dataset, the overall bias_score * net_positivity is positive, with both liberal and conservative article groups having positive means. Thus, in our data, liberal articles had a net negative tone, and conservative articles had a net positive tone.



FIGURES 3-5: BIAS OF “SCARE”, “FAIL”, AND CONTROL TERM “PRESIDENT”

Dep. Variable:	bias_score	R-squared:	0.029			
Model:	OLS	Adj. R-squared:	0.027			
Method:	Least Squares	F-statistic:	11.37			
Date:	Tue, 05 May 2020	Prob (F-statistic):	7.85e-11			
Time:	12:58:54	Log-Likelihood:	-2660.9			
No. Observations:	1878	AIC:	5334.			
Df Residuals:	1872	BIC:	5367.			
Df Model:	5					
Covariance Type:	nonrobust					
	coef	std err	t	P> t	[0.025	0.975]
const	-0.3103	0.025	-12.180	0.000	-0.360	-0.260
hate	-0.3378	0.098	-3.430	0.001	-0.531	-0.145
fail	-0.1578	0.061	-2.590	0.010	-0.277	-0.038
critic	-0.2040	0.040	-5.054	0.000	-0.283	-0.125
scare	-0.3177	0.179	-1.775	0.076	-0.669	0.033
president	-0.0738	0.216	-0.342	0.732	-0.497	0.349

FIGURE 6: OLS REGRESSION RESULTS