

1/1/2022

Hate Speech Data Analysis

A Brief Study of Hateful

Sentiments on Twitter During The

2016 & 2020 Presidential Elections

Using Machine Learning

Khoa Cao

I. Abstract

The widespread use of social media has made it easier than ever for individuals to express their thoughts and opinions online. Unfortunately, this also means that hateful sentiments, such as racism, sexism, and bigotry, can be spread quickly and easily. Due to the COVID-19 pandemic, the ensuing economic downturn, and the government's restrictions on outdoor activities, many people's economic status and, specifically, mental health have suffered (Cullen et al., 2020). Paired with social media use skyrocketing during that same time, sites such as Twitter and Facebook have been chosen by many as a place to share their thoughts and escape from their struggles (Cinelli et al., 2020). Many argue that due to the pandemic and a myriad of other reasons, hateful sentiments have risen between the 2016 and 2020 presidential elections. The purpose of this independent study project is to investigate the change in the prevalence of hateful sentiments in social media and to understand the motivations behind them.

To achieve its goals, the project will utilize a combination of qualitative and quantitative research methods. The study will first conduct a content analysis of social media posts to quantify the occurrence of hateful sentiments. Twitter will be the focus of this study since it is the main channel for political discussion, specifically during election time. It has been used frequently by most presidential candidates in the past two election cycles. This analysis will be supplemented by third-party psychological research, in order to gain a deeper understanding of people's motivations and attempt to explain the analysis.

This independent study project will be of interest to researchers and practitioners working in the fields of social media, internet studies, and media psychology. By shedding light on this issue, the project aims to inform future efforts to mitigate the spread of hateful sentiments on social media. Its results will have implications for policymakers, as well as for social media platforms and tech companies, as they seek to create a more inclusive and respectful online environment.

II. Data

a. Datasets

I attempted to collect my own data using the Twitter API to get tweets with the specific tags and timeframe that my analysis required. However, the tool required an application that

proved unsuccessful, and I was only allowed one dataset of a maximum of 100,000 data points. Therefore, most datasets in this analysis have been taken from independent collectors on Kaggle.com, a social platform for data science projects and research. Originally, these datasets were scraped from the Twitter archive using Twitter API – a Twitter-created tool for developers and researchers.

The datasets contain tweets and their associated data: text, creator, time, id, location, etc. The tweets collected for this study contain tags relevant to the 2016 and 2020 elections and mentions of the candidates themselves. Due to limitations on data collection, each dataset will vary in scope, time range, and specifications (tags and mentions). Each dataset's name is original to its creator, who will be cited.

The specification for each dataset used in this study is listed below. “.csv” denotes the file type: comma-separated values. The word clouds are generated by Python's WordCloud library showing the frequency of each word based on size.

1. Train.csv

This is the training dataset for the Machine Learning algorithm, consisting of two columns: “tweets” and “labels.” It has 31,962 tweets classified as either “neutral,” labeled “0,” or “hateful,” labeled “1.” Analytics Vidhya created it for a machine learning hackathon. Due to the ambiguity and subjectivity associated with classifying hate speech, this will accept their definitions as true and will compare its results to similar studies.

```
In [4]: df.info()
<class 'pandas.core.frame.DataFrame'>
Int64Index: 31962 entries, 0 to 31961
Data columns (total 2 columns):
#   Column  Non-Null Count  Dtype
---  -
0    tweet   31962 non-null   object
1   labels  31962 non-null   object
dtypes: object(2)
memory usage: 749.1+ KB
```

Figure 1: train.csv dataset info

Here are some examples of hateful and neutral tweets in the training set.

9	welcome ! i'm !	Neutral
10	â ¤ ireland consumer price index (mom) climbed previous .% .% may blog silver gold forex	Neutral
11	selfish. orlando standwithorlando pulseshooting orlandoshooting biggerproblems selfish heabreaking values love	Neutral
12	get see daddy today!! gettingfed	Neutral
15	ouch...junior angryð É junior yugyom omg	Neutral
16	thankful paner. thankful positive	Neutral
18	friday! ð smiles around via ig user: cookies make people	Neutral
19	know, essential oils made chemicals.	Neutral
20	people blaming ha conceded goal fat rooney gave away free kick knowing bale hit there.	Neutral
21	sad little dude.. badday coneofshame cats pissed funny laughs	Neutral
22	product day: happy man wine tool who's weekend? time open & drink up!	Neutral
24	tgif ff gamedev indie dev indiegame dev squad!	Neutral
25	beautiful sign vendor \$!! upsideofflorida shopolassyas love	Neutral
26	smiles media !!ð ð pressconference antalya turkey ! sunday throwback love! ð ð â¤¸î ¨	Neutral
27	great panel mediatization public service	Neutral
28	happy father's day ð ð ð ð	Neutral

Figure 2: *train.csv* Neutral Tweets

579	yes lets this,suppoing openly ,prowar anti islamic,homophobic,rapist,who advocates same,hypocrite	Hateful
609	porn vids www.smallgirlsex.com	Hateful
617	latest obsidian radio daily! thanks latesnews	Hateful
621	might libtard if... libtard sjw liberal politics	Hateful
622	overwhelming evidence company trump keeps echoes sentiments & ideals dumptrump	Hateful
649	"nigger?" lifelessons white kid grew 's. blogpost whitepeople respectâ	Hateful
692	carolyn cooper ugly, poor, ignorant black!	Hateful
733	no, definitely mexican. fakenewsale	Hateful
735	standing racism hate americad curse congress haters deplorable hateâ	Hateful
746	trump used hate, putin win white house trump presses attack khan family g.o.p. leaders	Hateful
749	everytime wear soccer shis joie fries says look mexican fuck ð	Hateful
763	danger white liberalism oveurning (applies liberal men feminism too)	Hateful

Figure 3: train.csv Hateful Tweets

2. 2016_US_election_tweets.csv

This dataset contains 100,000 tweets from 2016-08-30 till 2017-02-28 containing mentions of candidates: Hillary Clinton, Donald Trump, and Bernie Sanders. It stores 18 different data points for each tweet, including id, text, and time. This study is only interested in the posted time and content of each tweet.

```
In [7]: runcell(2, 'C:/School Work/Courses/22 - 23 School
Analysis/HateSpeechDetection.py')
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 100000 entries, 0 to 99999
Data columns (total 18 columns):
#   Column                Non-Null Count  Dtype
---  -
0   id                     100000 non-null   int64
1   candidate_id           100000 non-null   int64
2   tweet_id               100000 non-null   float64
3   polarity               100000 non-null   float64
4   subjectivity           100000 non-null   float64
5   retweet_count          100000 non-null   int64
6   favorite_count         100000 non-null   int64
7   device                 100000 non-null   int64
8   retweeted_status_id    44607 non-null    float64
9   lang                   91451 non-null    object
10  state                  3279 non-null     object
11  text                   55393 non-null    object
12  created_at             100000 non-null    object
13  inserted_at            100000 non-null    object
14  updated_at             100000 non-null    object
15  tw_user_id             11060 non-null    float64
16  latitude                0 non-null        float64
17  longitude              0 non-null        float64
dtypes: float64(7), int64(5), object(6)
memory usage: 13.7+ MB
None
```



```

In [9]: runcell(2, 'C:/School Work/Courses/22 - 23 School Year/Analysis/HateSpeechDetection.py')
<ipython-input-9-963c3f554fac>:1: DtypeWarning: Columns (1) have dtype object, using dtype object on import or set low_memory=False.
runcell(2, 'C:/School Work/Courses/22 - 23 School Year/Analysis/HateSpeechDetection.py')
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 971157 entries, 0 to 971156
Data columns (total 21 columns):
#   Column                Non-Null Count  Dtype
---  -
0   created_at            971088 non-null    object
1   tweet_id              971073 non-null    object
2   text                  971073 non-null    object
3   likes                 971045 non-null    object
4   retweet_count         970933 non-null    float64
5   source                970057 non-null    object
6   user_id               970929 non-null    object
7   user_name             970917 non-null    object
8   user_screen_name      970933 non-null    object
9   user_description      869663 non-null    object
10  user_join_date         970779 non-null    object
11  user_followers_count   970917 non-null    object
12  user_location          675839 non-null    object
13  lat                    445702 non-null    object
14  long                   445705 non-null    object
15  city                   227180 non-null    object
16  country                442732 non-null    object
17  continent              442749 non-null    object
18  state                  320614 non-null    object
19  state_code             300414 non-null    object
20  collected_at           970765 non-null    object
dtypes: float64(1), object(20)
memory usage: 155.6+ MB
None

```

Figure 8: 2020_hashtag_donald_trump.csv dataset info

b. Data cleaning

Before using our datasets, we must get rid of noise within the datasets. Tweets often contain stop words – insignificant symbols, words, or phrases that serve no purpose in our analysis and that crowd out critical phrases – such as punctuation, common phrases such as “the,” “a,” or “and,” URLs, retweet tags, or mentions of other users. Python’s Regular Expressions (Re) module provides regular expressions matching operations. The relevant operation to this study is the substitution method - `re.sub()` – which allows for the identification and removal of stop words and custom phrases without affecting the rest of the text.

```

def cleanTwts(text):
    text = str(text).lower()
    text = re.sub(r'@[A-Za-z0-9]+', '', text)
    text = re.sub(r'#', '', text)
    text = re.sub(r'RT[\s]+', '', text)
    text = re.sub(r'https?:\/\/\S+', '', text)
    text = re.sub('<.*?>+', '', text)
    text = re.sub('\n', '', text)
    text = re.sub('\w*\d\w*', '', text)
    text = [word for word in text.split(' ') if word not in stopwords]
    text = ' '.join(text)
    return text

#Clean tweets
df["tweet"] = df["tweet"].apply(cleanTwts)

```

Figure 9: Tweet-cleaning method

III. Analysis

a. Hate Speech Detection Algorithm

This study employs a Natural Language Processing (NLP) library, TextBlob, and a machine learning (ML) algorithm to evaluate the tweets' sentiment (positive, negative, or neutral) and its hatefulness (hateful or neutral). Classifying sentiments and hate speech can be extremely subjective, ambiguous, and contextually dependent (Luvell and Barnes, 2022), the two different approaches are meant as checks for one another.

TextBlob is a simplified text processing and analysis tool built for Python. It provides a simple API for diving into NLP tasks such as sentiment analysis, text classification, and even spelling correction. Even though TextBlob only has only an average 70% confidence score on its prediction, what makes it more powerful for small-scale projects such as this one is the fact that it is popular, so it is robust and widely documented, prebuilt, and well maintained. This means that it can be used by anyone and tailored to any project without having to train it, which, without a sufficiently sized training set, could result in low accuracy scores.

The second tool used by this study is an ML algorithm, namely the Decision Tree (DT) Classifier implemented Scikit Learn (sklearn) library for Python. Simplistically, a DT is a non-parametric supervised learning method used for classification and regression. It is a model that predicts the value of a target variable by learning simple decision rules inferred from the data

features. A tree can be seen as a piecewise constant approximation. DTs are one of the most powerful tools used by data scientists to categorize qualitative data, such as sentiments and hate speech. To name some of its myriad advantages, it requires little data preparation, has low memory and time costs, and performs well even if the true model from which the data were generated disagrees with its assumptions. The first and last advantages are critical to this study, as Twitter data is riddled with noise and generally poorly formatted. Typical data preparation steps, such as data normalization/standardization, missing value treatment, outlier capping, etc., are not required for the decision tree, making it a ‘go-to’ algorithm for data scientists. Furthermore, other algorithms, such as linear regression and the popular Naïve Bayes theorem, require various specific assumptions that need to be fulfilled in order to work properly. However, DTs are non-parametric; thus, we need not make any significant assumptions or consider data distribution.

b. Methodology

After cleaning and importing the datasets into Python as separate data frames, we extract only the data of interest – time and text. Next, we run the data through the sentiment analysis tool to categorize them into three groups: “Positive sentiment,” “Negative sentiment,” and “Neutral sentiment.” The DT will also make predictions on the data and classify it into two groups: “Hateful” and “Neutral” (again, based on each of the tools’ lexicons and training data).

IV. Results

a. Algorithm Predictions

As mentioned previously, the sentiment analysis and hate speech classification algorithms separate texts into distinct groups: positive, negative, or neutral for the former; hateful or neutral for the latter. We are particularly interested in the percentages of tweets in each of these groups. Since the number of tweets in each dataset varies greatly (so the number of tweets in each timeframe varies), percentages give a fairer and more accurate understanding of the change in sentiments over time. Since the sentiment analysis and DT classification are separate analyses, their figures add up to 200%; it is the figures for each analysis that add up to 100%.

For the 2016 dataset, TextBlob identified 10.74% of the tweets to be of “negative sentiment.” For the 2020hashtag_donaldtrump and 2020hashtag_joe Biden datasets, this figure increased, to 16.37% and 12.2%, respectively. The DT Classifier also showed increases in hateful sentiments between 2016 and 2020. In 2016, only 6.65% were “negative.” However, the 2020 figures were much higher than that of 2016: 29.62% of tweets relating to Donald Trump and 13.88% of tweets relating to Joe Biden were identified as “hateful.” It is notable that the percentages of tweets relating to Trump that are hateful is double that of the tweets relating to Biden.

The figures for each analysis are shown below.

- 2016_US_election_tweets
 - o TextBlob Sentiment Analysis
 - Positive tweets: 14943, 14.94% of total
 - Negative tweets: 10743, 10.74% of total
 - Neutral tweets: 74314, 74.31% of total
 - o Decision Tree Classifier Predictions
 - Hateful tweets: 6650, 6.65% of total
 - Neutral tweets: 93350, 93.35% of total
- 2020hashtag_donaldtrump
 - o Decision Tree Classifier Predictions
 - Hateful tweets: 287652, 29.62% of total
 - Neutral tweets: 683505, 70.38% of total
 - o TextBlob Sentiment Analysis
 - Positive tweets: 254625, 26.22% of total
 - Negative tweets: 158971, 16.37% of total
 - Neutral tweets: 557561, 57.41% of total
- 2020hashtag_joe Biden
 - o Decision Tree Classifier Predictions
 - Hateful tweets: 107847, 13.88% of total
 - Neutral tweets: 669231, 86.12% of total
 - o TextBlob Sentiment Analysis:

- Positive tweets: 222178, 28.59% of total
- Negative tweets: 94818, 12.2% of total
- Neutral tweets: 460082, 59.21% of total

Here are some examples of neutral and hateful tweets identified by the DT classifier. The preview only shows a portion of the text so it may not be representative of the whole text.

biden made i trumpisnotamerica i	Neutral
watching setting dvr. let's give bonus ratings!! joe Biden	Neutral
ensorship hunter Biden Biden Bidenemails Bidenemail corruption	Neutral
is wrong?!!!" cory booker's brilliant final questioning trump nominee amy coney barrett amyconeybarrett corybooker barrett booker trump kamalaharris joe Biden scotus supremecourtconfirmation	Neutral
ny post censorship censored twitter manipulate us election favor joe Biden trump.but ccp china porn twitter? that's always fine . sick?	Neutral
Biden	Neutral
proof bidens crooked. twitter suspend sharing bidencrookedbidenukraniacollusion democrats hunter Biden emails	Neutral
bi allegedly obtained hunter Biden computer, data Ukraine dealings, report claims joe Biden hunter Biden	Neutral
joe Biden point man	Neutral
omments this? "do democrats understand ruthless china is?" china hunter Biden joe Biden Biden Harris trump ealdonaldtrump wto coronavirus trade	Neutral
joe Biden admitting	Neutral
'm going share things like Biden more. too. Biden cares Biden chickent trump kamalaharris	Neutral
hunter Biden hunter Biden emails joe Biden joe Biden must step down	Neutral
effort find truth allegations allowing people share link article hunter Biden, popped ny post hunter Biden Biden	Neutral
joe Biden calls liar insults overweight. Biden shows low iq daily. maybe Joe hold town halls kindergarten. he'll amongst equals one say anything takes afternoon nap	Neutral
vote wisely.and wisely mean Joe.. elsewe see t.v. for...joe Biden	Neutral
Biden crime family joe Biden hunter Biden hunter Biden emails	Neutral
vote joe Biden	Neutral

Figure 10: Hateful Tweets Examples

biden made i trumpisnotamerica i	Neutral
watching setting dvr. let's give bonus ratings!! joe Biden	Neutral
ensorship hunter Biden Biden Bidenemails Bidenemail corruption	Neutral
is wrong?!!!" cory booker's brilliant final questioning trump nominee amy coney barrett amyconeybarrett corybooker barrett booker trump kamalaharris joe Biden scotus supremecourtconfirmation	Neutral
ny post censorship censored twitter manipulate us election favor joe Biden trump.but ccp china porn twitter? that's always fine . sick?	Neutral
Biden	Neutral
proof bidens crooked. twitter suspend sharing bidencrookedbidenukraniacollusion democrats hunter Biden emails	Neutral
bi allegedly obtained hunter Biden computer, data Ukraine dealings, report claims joe Biden hunter Biden	Neutral
joe Biden point man	Neutral
omments this? "do democrats understand ruthless china is?" china hunter Biden joe Biden Biden Harris trump ealdonaldtrump wto coronavirus trade	Neutral
joe Biden admitting	Neutral
'm going share things like Biden more. too. Biden cares Biden chickent trump kamalaharris	Neutral
hunter Biden hunter Biden emails joe Biden joe Biden must step down	Neutral
effort find truth allegations allowing people share link article hunter Biden, popped ny post hunter Biden Biden	Neutral
joe Biden calls liar insults overweight. Biden shows low iq daily. maybe Joe hold town halls kindergarten. he'll amongst equals one say anything takes afternoon nap	Neutral
vote wisely.and wisely mean Joe.. elsewe see t.v. for...joe Biden	Neutral
Biden crime family joe Biden hunter Biden hunter Biden emails	Neutral
vote joe Biden	Neutral

Figure 11: Neutral Tweets Examples

b. Plotting Hateful Sentiment vs Time

We can plot the DT's analysis over time using Matplotlib, a plotting library for Python. This allows us to track the change in hateful sentiment approaching the election dates.

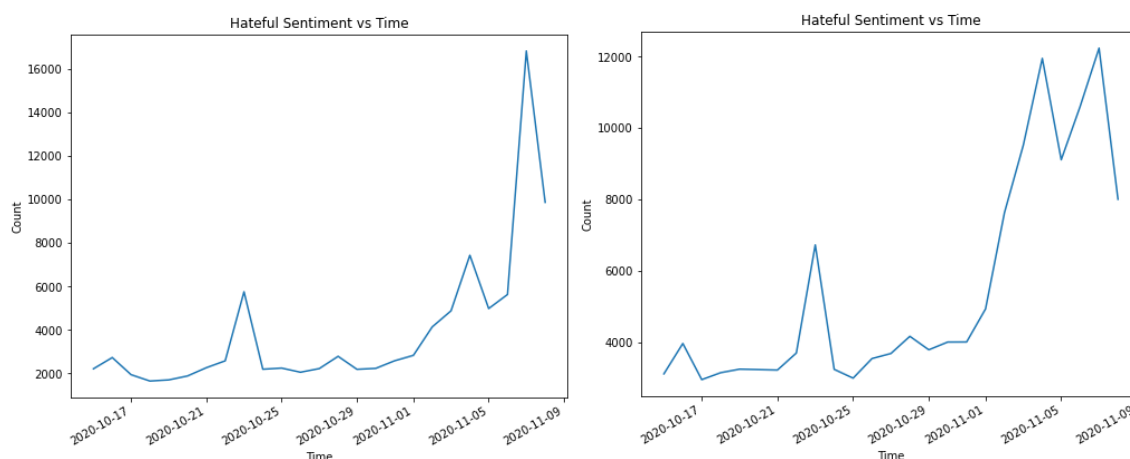


Figure 12: Tweets Relating to Joe Biden (left) and Donald Trump (right)

The graph illustrates a gradual rise in hateful sentiments as election day (November 3rd, 2020) approaches, with figures skyrocketing on election day and peaking a few days afterward. The three highest peaks of each graph correspond to the three critical events of the election: the third presidential debate (2020-10-22), Election Day (2020-11-03), and President Biden's first speech (2020-11-07). The highest volume of hateful tweets in both graphs is on the day of Biden's speech.

c. Accuracy

To test the accuracy of TextBlob sentiment analysis and DT classifier, we split the training sets of each tool into two equal datasets (around 15,000 data points each). Using the first half as the training set, we let the algorithms perform predictions on the other half. We then simply calculate the percentage of accurate predictions: TextBlob is 70% accurate, and the Sklearn DT is more than 98% accurate. It is crucial to note that these accuracy scores compare the algorithms' predictions with their training data, which is subject to their creators' interpretations and opinions.

V. Conclusion

Two separate analyses on more than 1.7 million tweets surrounding the 2016 and 2020 elections were conducted. These analyses split the tweets into distinct groups based on their sentiments and hatefulness. The change in the volume of tweets in each group is then plotted in a time-series graph. The data has led to the following observations.

As evident from the numerical data from the sentiment analysis and hate speech analysis, the percentages of negative and hateful sentiments increased between the 2016 U.S. presidential election to the 2020 election. Not only so, but these sentiments also seemed to increase as election day in November approached. These results prove that negative and hateful sentiments rise as the critical day nears and as the election becomes a hot topic. Our results agree with a 2023 study, *Offline events and Online hate*, on hate speech and its correlation to significant real-world events (Lupu et al., 2023).

This begs the question: why do these sentiments gain popularity around election time, and what factors affect their rise? The media and politicians' attention plays an important role. The volume of hateful sentiments likely depends on the media's coverage of the associated events. For example, the rise in hateful sentiments between the last two elections could be explained by the Covid-19 pandemic. As news outlets cover the pandemic and politicians use it as their main campaign focus, anti-Asian discriminatory sentiments rose (Chen and Alexander, 2020). Controversy accompanies real-world events in increasing hate speech. Recent religious, racial, and ethnic events relating to figures such as Vice President Kamala Harris, the Iranian General Qasem Soleimani, and George Floyd sparked sharp rises in gender, racial, and religious hate speech. Following George Floyd's altercation with the police, hateful speech rose by 250% (Faguy, 2023). This rise in interest and powerful social media recommendation algorithms led to a boom in hateful sentiments online (Laub, 2019). Since social media platforms frequently polish recommendation algorithms to tailor content to maximize engagement and ad revenues, they sometimes allow hate groups to reach their target audiences more quickly than ever. This could explain the rise in hateful sentiments between the 2016 and 2020 elections.

This study faces many potential flaws due to its limited scope. The first limitation is my inability to gather data independently from the Twitter API. Instead, I had to rely on pre-collected data from third-party sources, which did not allow me to tailor the data – the tags, timeframe, and amount - to my needs. For example, the 2016 dataset was minimal, as the data

was unpopular and data science was still in its infancy. Thus, many user accounts have since been deleted, leaving only 100,000 data points to work with, creating significant potential for errors. Furthermore, the dataset used to train the DT classifier only contained about 32,000 data points, and an independent party developed it, so it was not peer-reviewed. Thus, the DT's classifications might disagree with those of a psychologist. For example, the algorithm categorized a tweet about the rise in bigotry and racism as hateful since it contained the words "hate" and "die."

Before understanding the relationship between hateful sentiments and the U.S. elections, much work must be done. Such a question requires interdisciplinary collaboration and research. However, online moderators and social media companies must mitigate hate speech on their platforms to protect their users. This study has proven that Machine Learning is a powerful tool to help identify hate speech. However, the painstaking work of defining, interpreting, and classifying language must be left to humans.

VI. References

<https://www.kaggle.com/datasets/manchunhui/us-election-2020-tweets>

https://www.kaggle.com/datasets/dv1453/twitter-sentiment-analysis-analytics-vidya?select=train_E6oV3lV.csv

<https://data.world/alexfilatov/2016-usa-presidential-election-tweets>

W Cullen, G Gulati, B D Kelly, Mental health in the COVID-19 pandemic, *QJM: An International Journal of Medicine*, Volume 113, Issue 5, May 2020, Pages 311–312, <https://doi.org/10.1093/qjmed/hcaa110>

Anderson, Luvell and Michael Barnes, "Hate Speech", The Stanford Encyclopedia of Philosophy (Spring 2022 Edition), Edward N. Zalta (ed.), URL = [<https://plato.stanford.edu/archives/spr2022/entries/hate-speech/>](https://plato.stanford.edu/archives/spr2022/entries/hate-speech/).

Cinelli, M., Quattrocioni, W., Galeazzi, A., Valensise, C. M., Brugnoti, E., Schmidt, A. L., ... & Scala, A. (2020). The COVID-19 social media infodemic. *Scientific reports*, 10(1), 1-10.

Lupu Y, Sear R, Velásquez N, Leahy R, Restrepo NJ, et al. (2023) Offline events and online hate. *PLOS ONE* 18(1): e0278511. <https://doi.org/10.1371/journal.pone.0278511>

Chen, H Alexander et al. "Anti-Asian sentiment in the United States - COVID-19 and history." *American journal of surgery* vol. 220,3 (2020): 556-557.
doi:10.1016/j.amjsurg.2020.05.020

Faguy, A. (2023, January 25). *Real-World Events Drive Increases In Online Hate Speech, Study Finds*. Forbes; Forbes. <https://www.forbes.com/sites/anafaguy/2023/01/25/real-world-events-drive-increases-in-online-hate-speech-study-finds/?sh=24a911f13d6d>

Laub, Z. (2019, June 7). *Hate Speech on Social Media: Global Comparisons* | Council on Foreign Relations. Council on Foreign Relations. <https://www.cfr.org/backgrounder/hate-speech-social-media-global-comparisons>