

Data Intensive Computing (CSE 587) – Lab 2

PROBLEM STATEMENT:

Collect data about from at least three sources, one opinion-based social media in twitter, research data in New York Times, and the third is the common crawl data for the same topic or key phrase, and similar time periods. Process the three data sets collected individually using classical big data methods. Compare the outcomes using popular visualization methods.

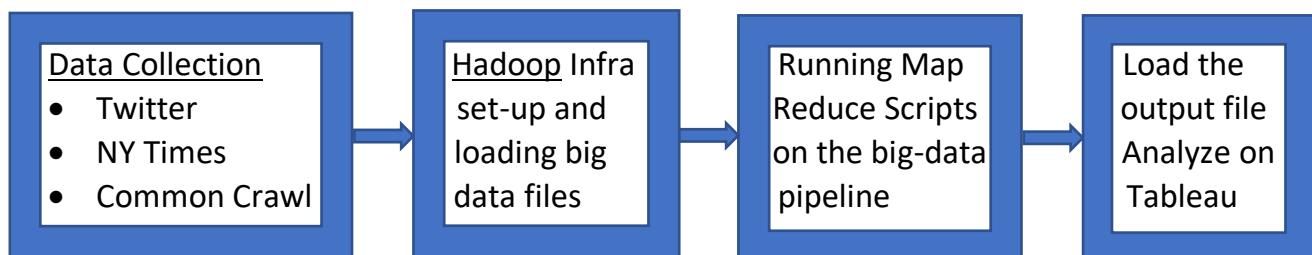
SOLUTION:

The topic of interest chosen for data collection and visualization is ***“IMMIGRATION”***.

Topic Related Key-words: US Immigration, U.S. Immigrant, US Civil wars, U.S. Citizenship and Immigration Services, USCIS, immigrant visa, Migrator, Migration, Refugees, Pilgrim, Green Card, Emigrate, Emigration, Illegal immigration, Legal immigration, border wall, etc.

The entire Data Flow is as below:

Data is collected and scrapped from various sources (Twitter, New York Times and Common Crawl). The data is then loaded onto the Hadoop infrastructure and Map-Reduce scripts are run on it. The word count and word co-occurrence files are then loaded into the Tableau for visualizing the most frequent words as *‘word-cloud’*.



→ Prototype data collection

Social Media – Tweets from Twitter

1. TwitterR library is loaded and respective keys are set-up and authorized.
2. Tweets are collected using “searchTwitter” command.
3. Tweets are collected (> **20,000 tweets**) for each day and stored separately and collectively.

Research Data - New York Times

1. articleAPI key is set-up using NYTimes developer’s account.

2. Articles are selected for various sub-topics using “api.search” command for the entire year 2019.
3. URLs are collected for each topic (Total URLs – 257) and the data is scrapped from each URL using “BeautifulSoup”. More than **500 articles** have been collected for NY Times.

Common Crawl

1. Code Reference: -
<https://www.bellingcat.com/resources/2015/08/13/using-python-to-mine-common-crawl/>
2. Libraries are loaded and an “**index-list**” from common crawl site is chosen. The domain passed on the index list is “**americanimmigrationcouncil.org**”.
3. All the URLs are collected (3671 hits with 1393 links discovered).
4. Data is then scrapped from each URL using “BeautifulSoup”. More than **500 articles** have been collected for Common Crawl.

→ Set up big data infrastructure

Hadoop Infrastructure is set-up for storing and analyzing the big data pipeline.

Install a Hadoop infrastructure in the one of the following ways:

- (i) cloudera docker image from the document provided, or
- (ii) virtual machine (VM) image for data storage in HDFS and Hadoop infrastructure, or
- (iii) amazon aws/ any other.

For **NY Times** data, the VM is installed for data storage for loading data and running map reduce files. README.txt file contains the instructions on setting the VM.

For **Twitter** and **Common Crawl**, Amazon AWS is utilized.

→ Data Pre and Post Processing

Pre-processing is done on each of the data files generated from the data collection. Using python, Mapper and Reducer scripts have been written to remove all the punctuations (?, ! etc.), unwanted symbols/characters (@, # etc.), **stop words**(said, would etc.). The output of mapper is then fed to the Reducer script where the data is grouped using key, value pair where keys are the words and the value is the count of each word. Reducer gives the collective **count of the words** as output. Using this file as input a python code is used to calculate the top 50 words and prepended it in the respective files.

In case of co-occurrence the top 50 words are read from the input file and each of these words are searched in the input file articles. These words and their **co-occurring words** are then treated as keys and fed to the reducer.

The output obtained for both “word-count” and “word-cooccurrence” is then converted into a csv file and is further used for analyzing and visualization.

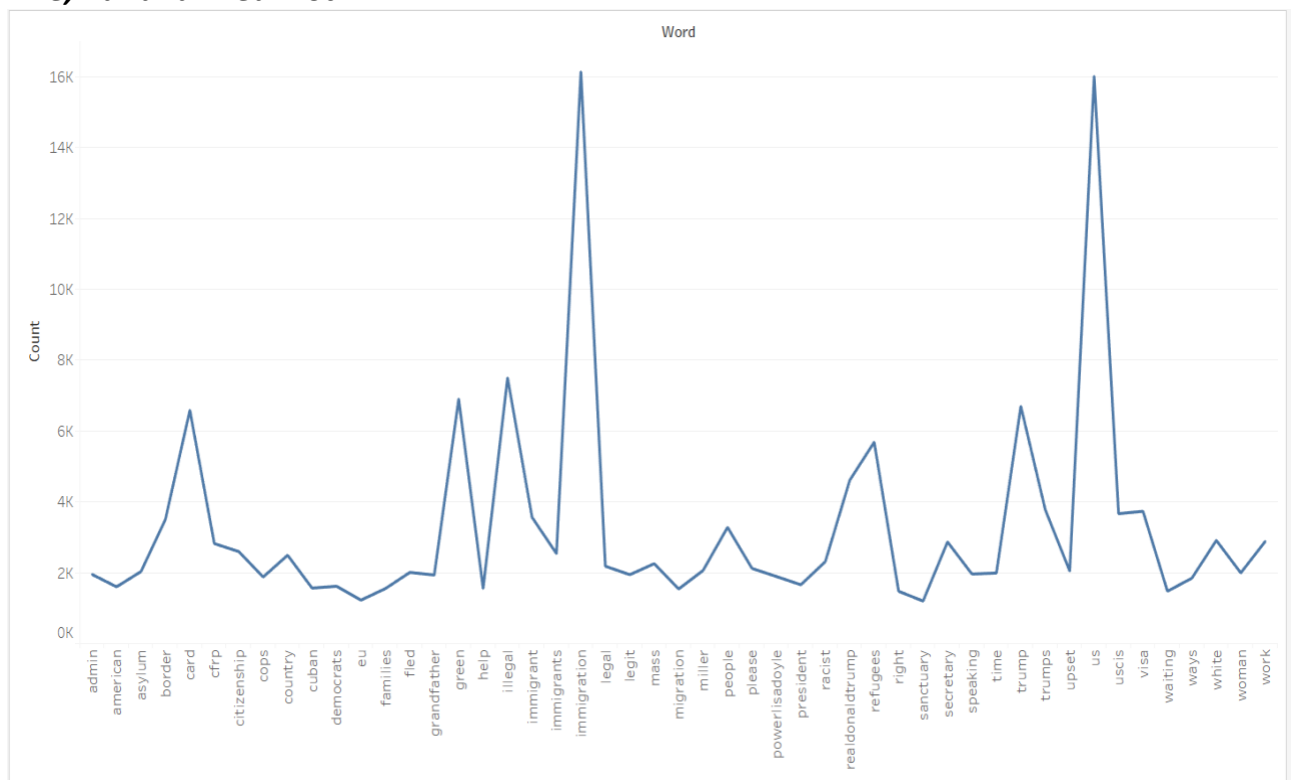
→ Analyze and visualize

Twitter Top 50 Word Count

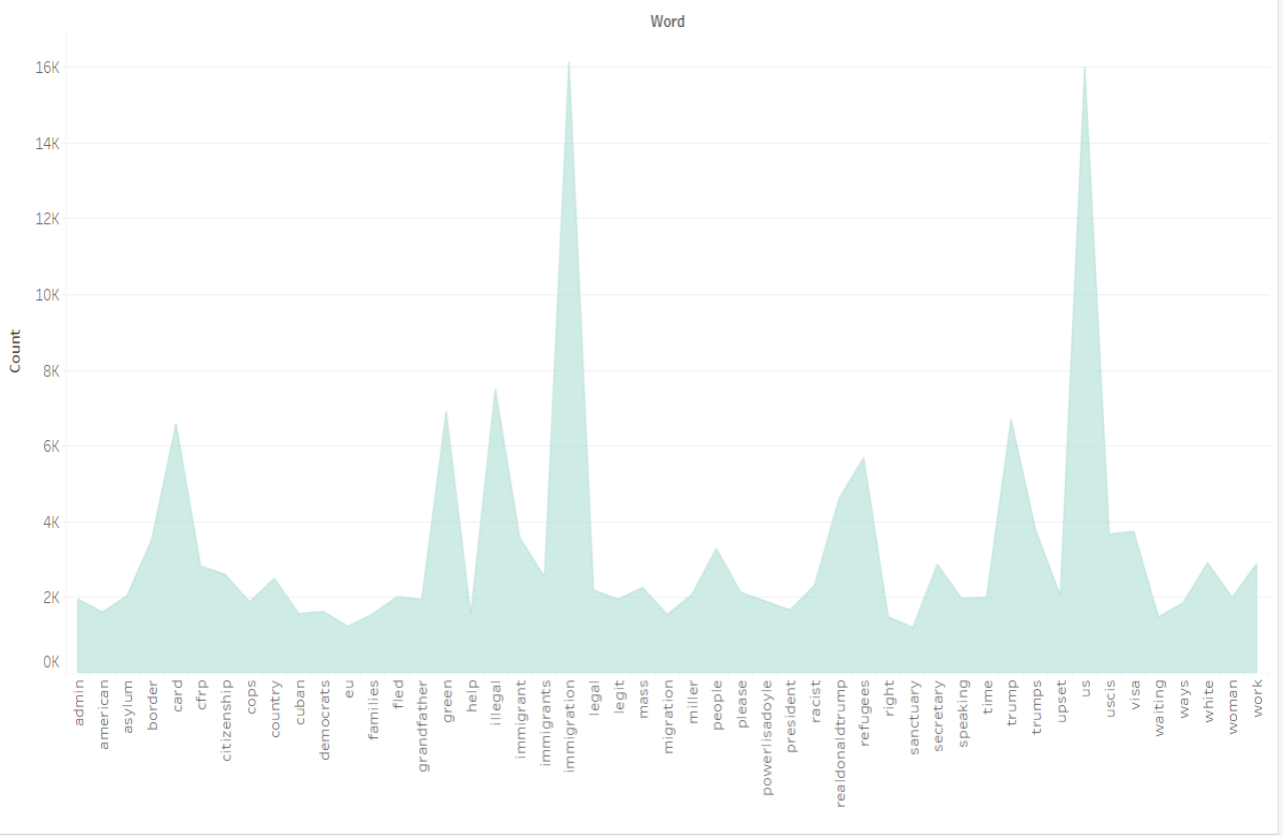
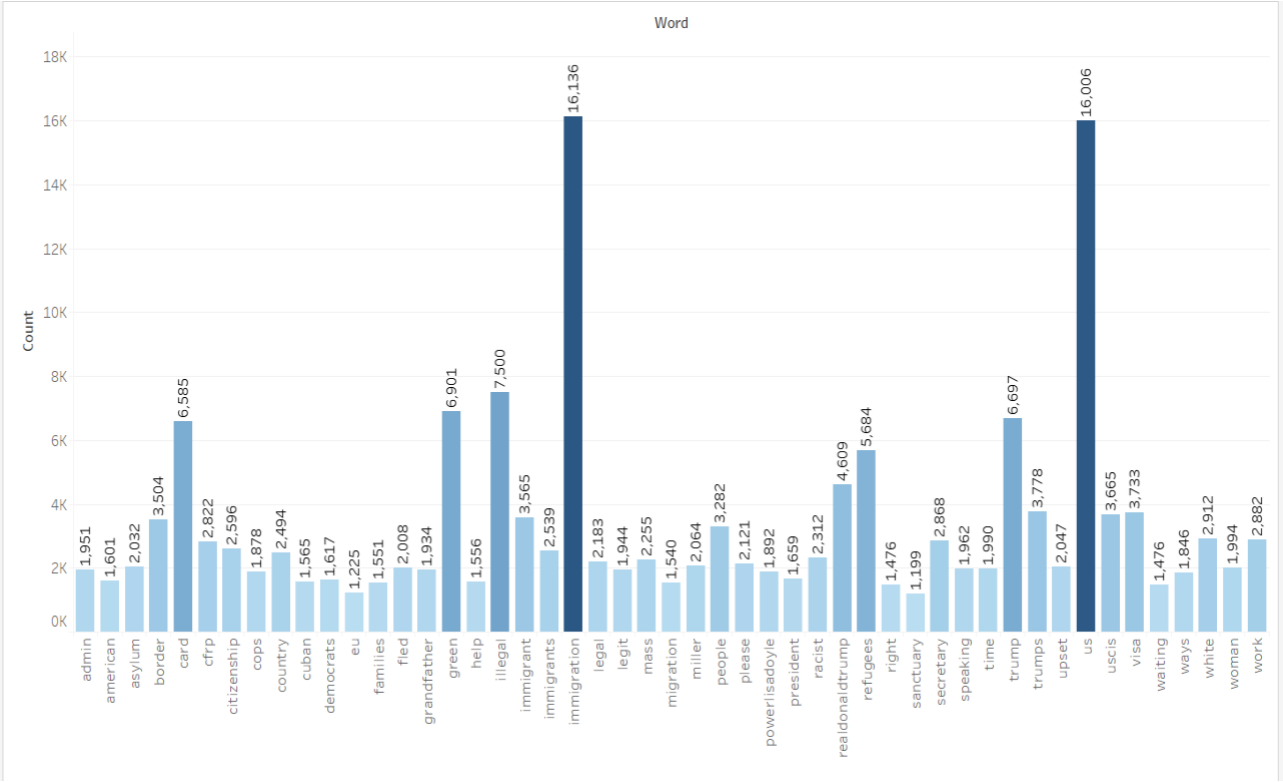
Word Cloud



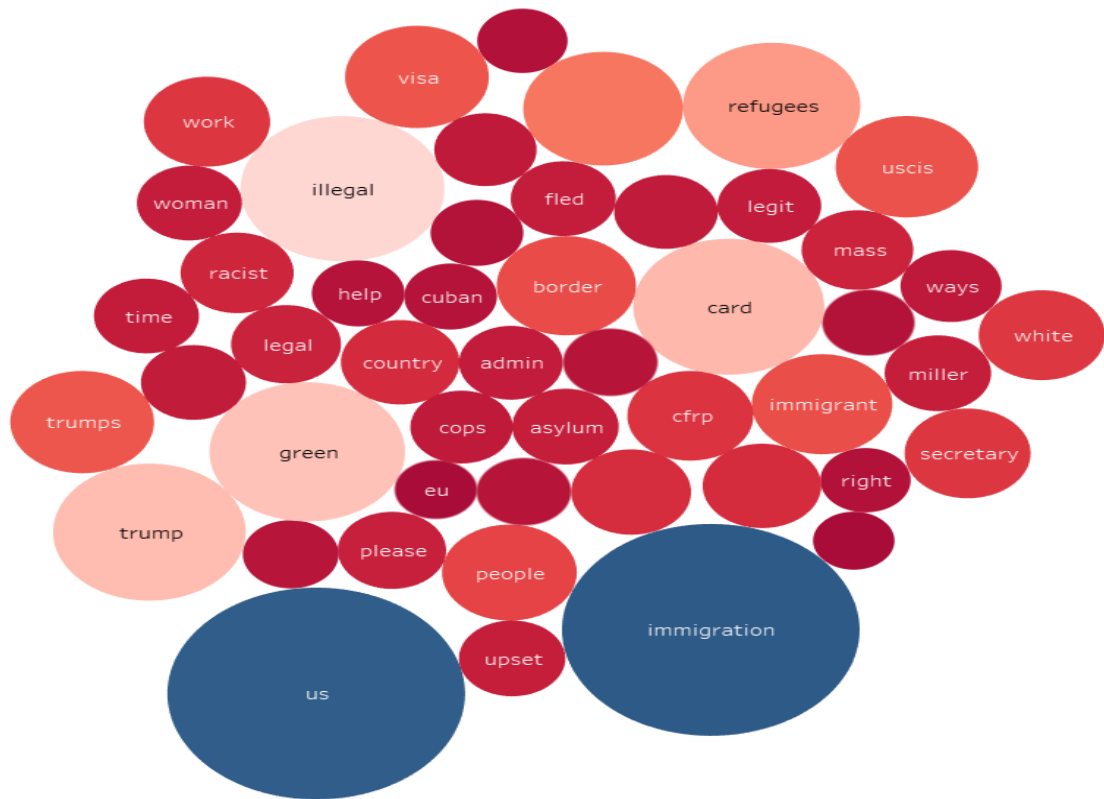
Line, Bar and Area Plot



TEAM MEMBERS: - PRADEEP KUMAR JOSHI, ABHISHEK GOSWAMI

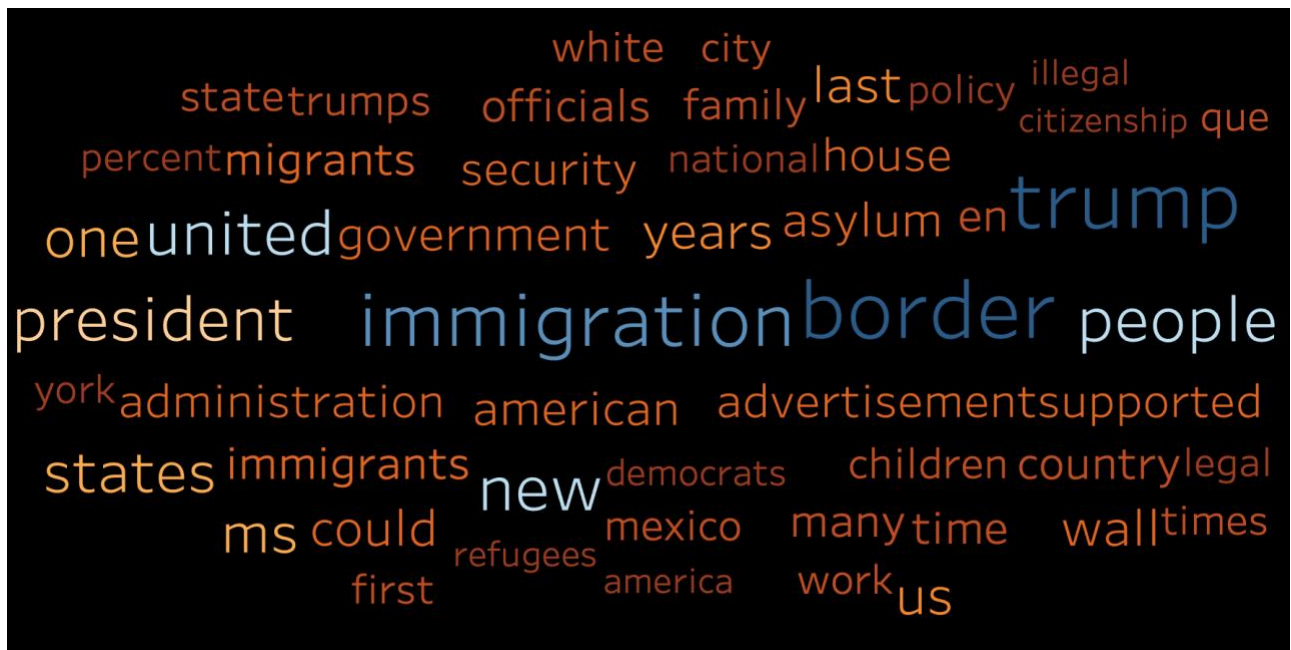


Round Plot

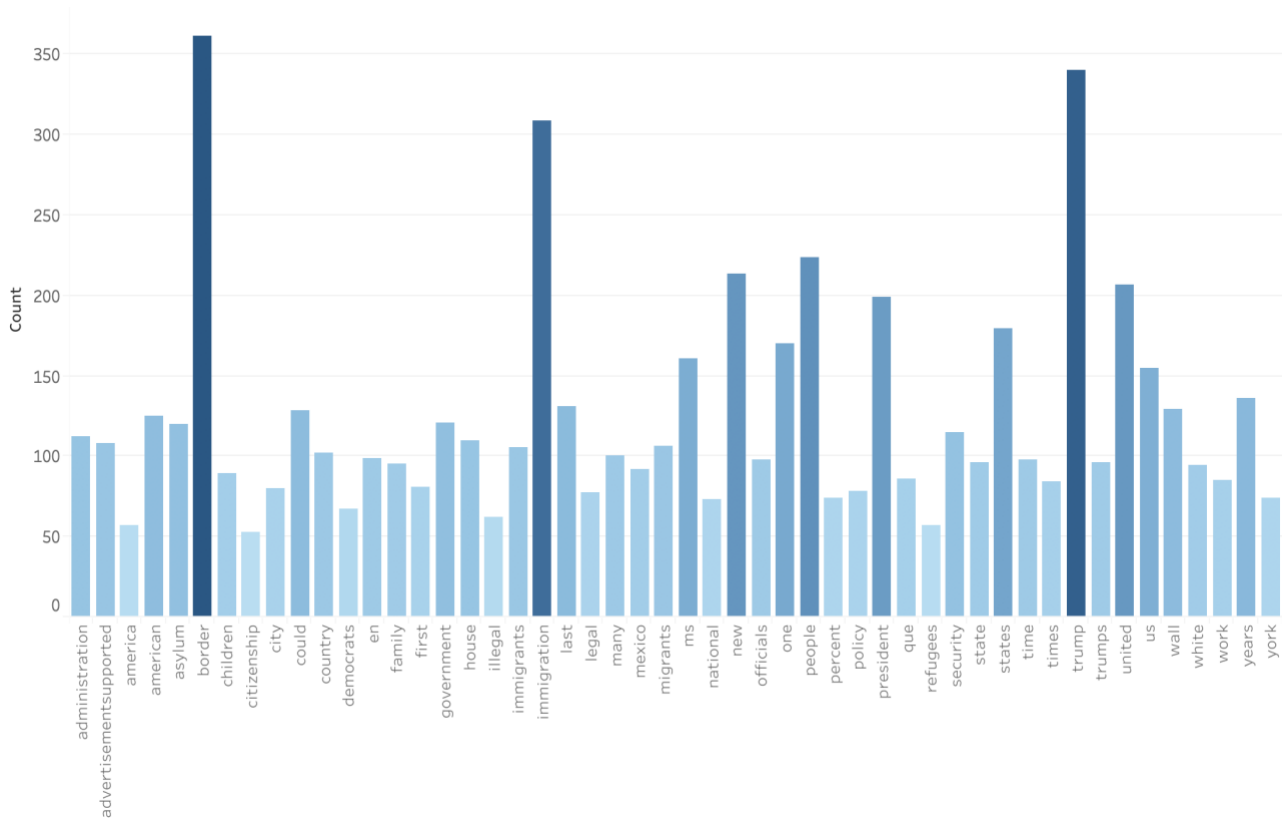
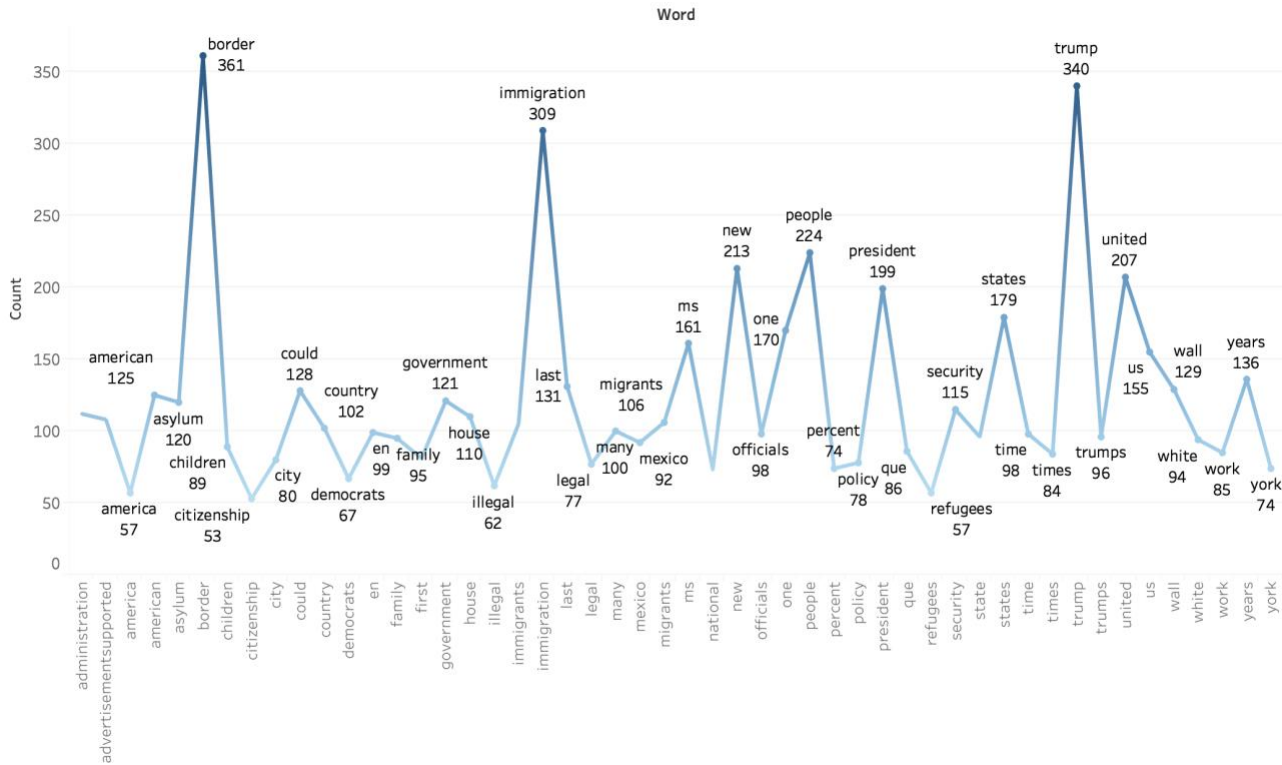


New York Times Top 50 Word Count

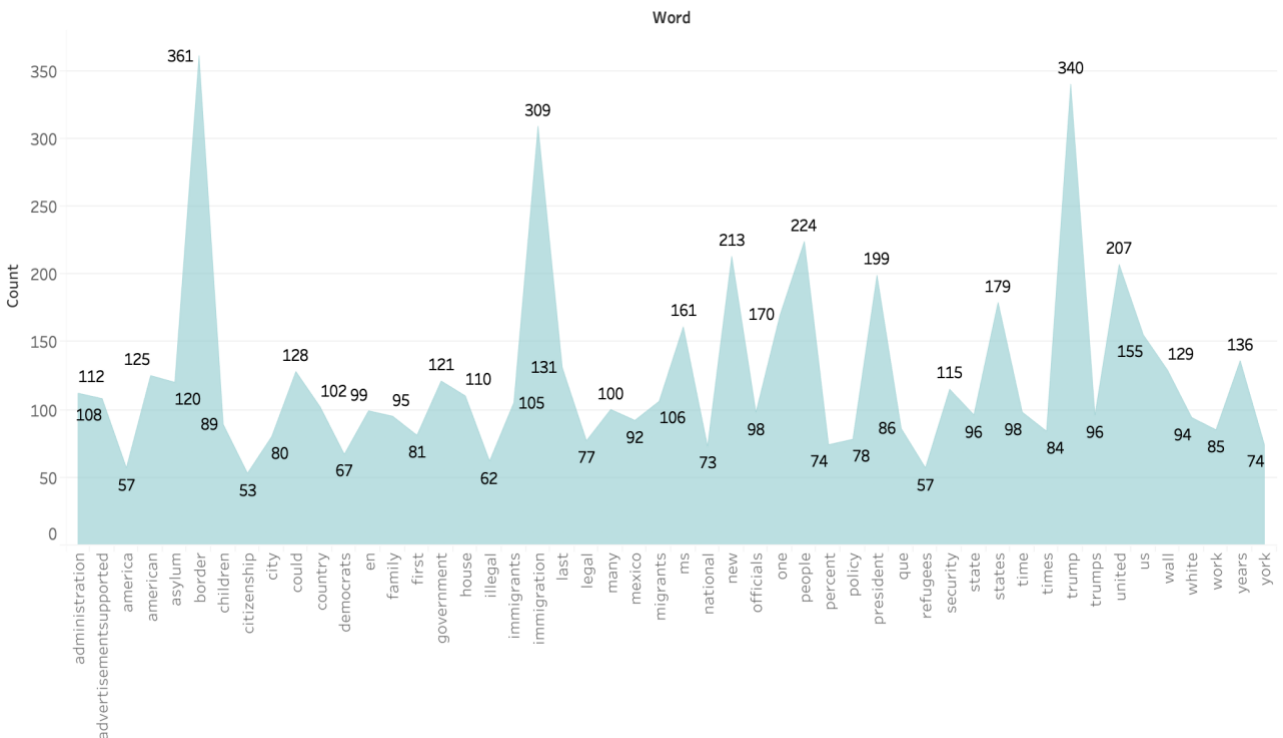
Word Cloud



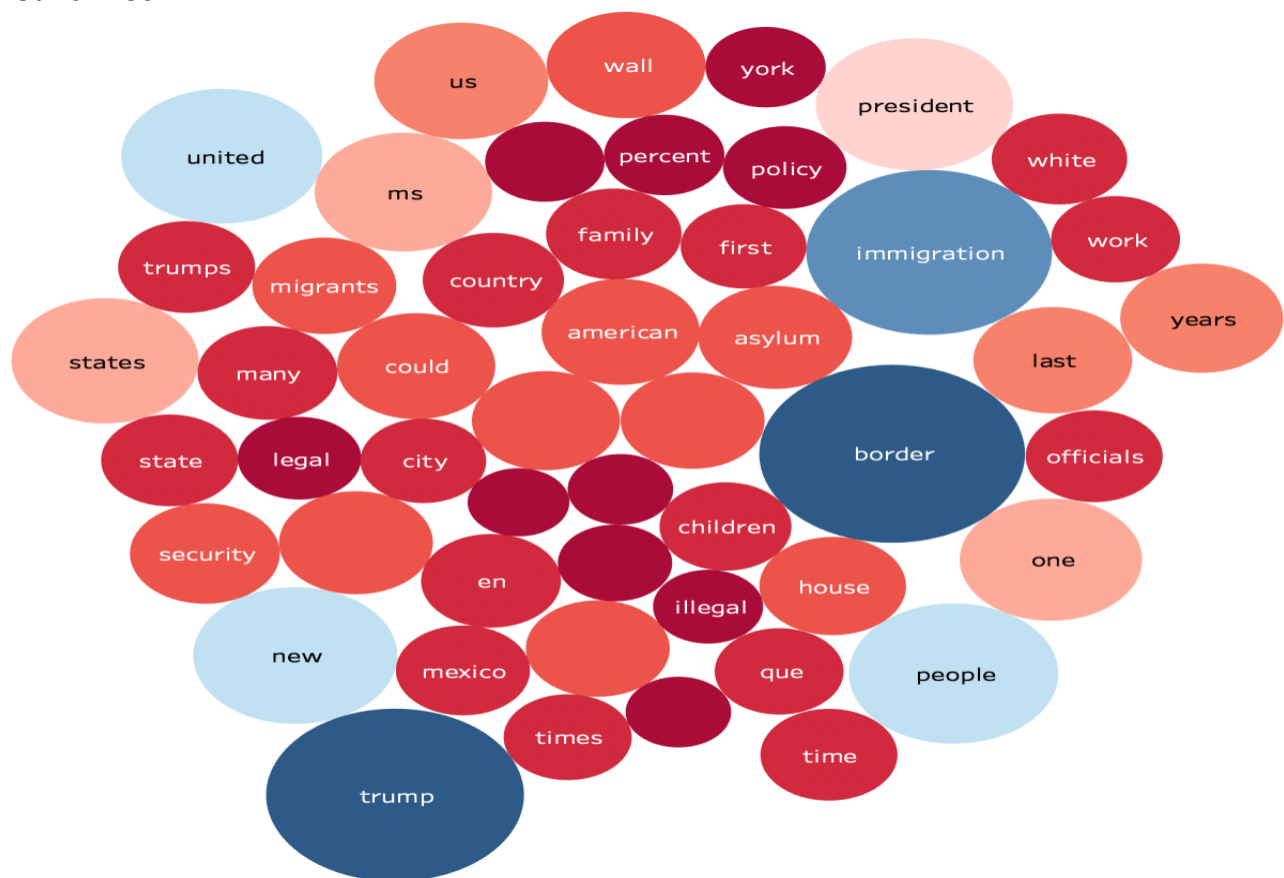
Line, Bar and Area Plot



TEAM MEMBERS: - PRADEEP KUMAR JOSHI, ABHISHEK GOSWAMI

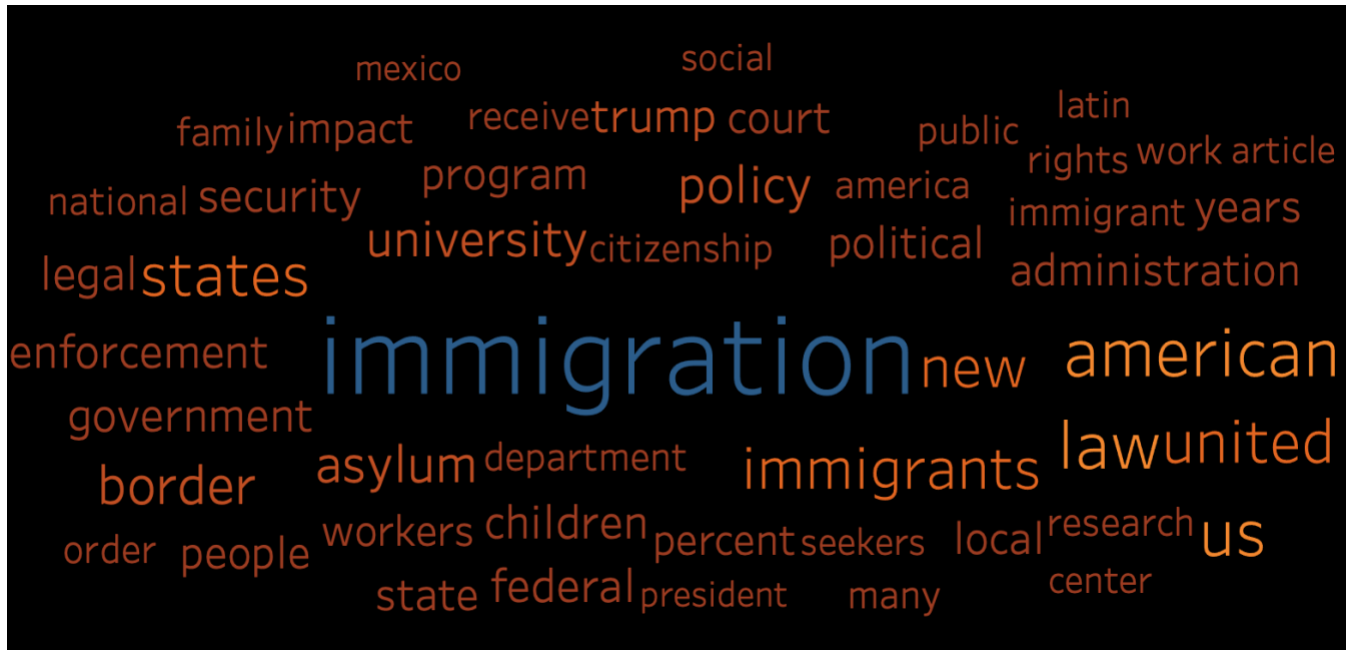


Round Plot

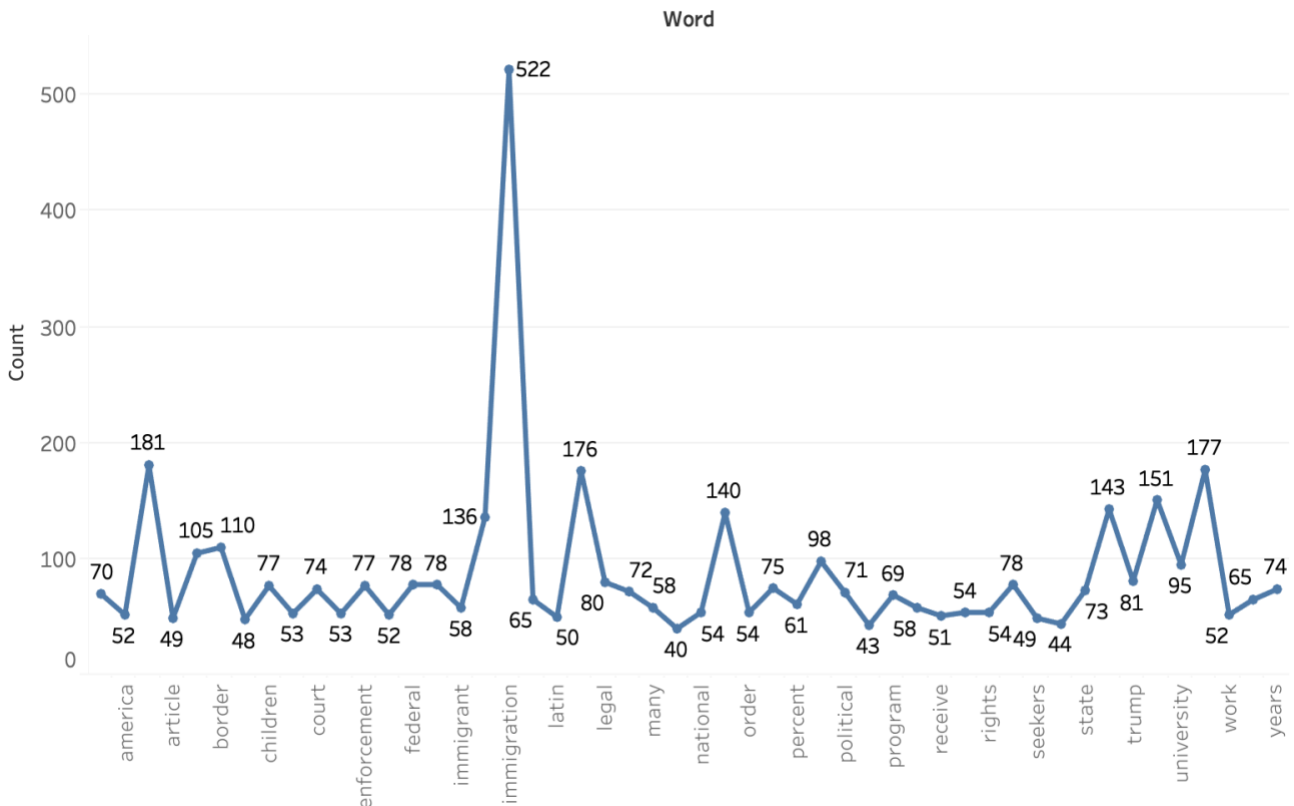


Common Crawl Top 50 Word Count

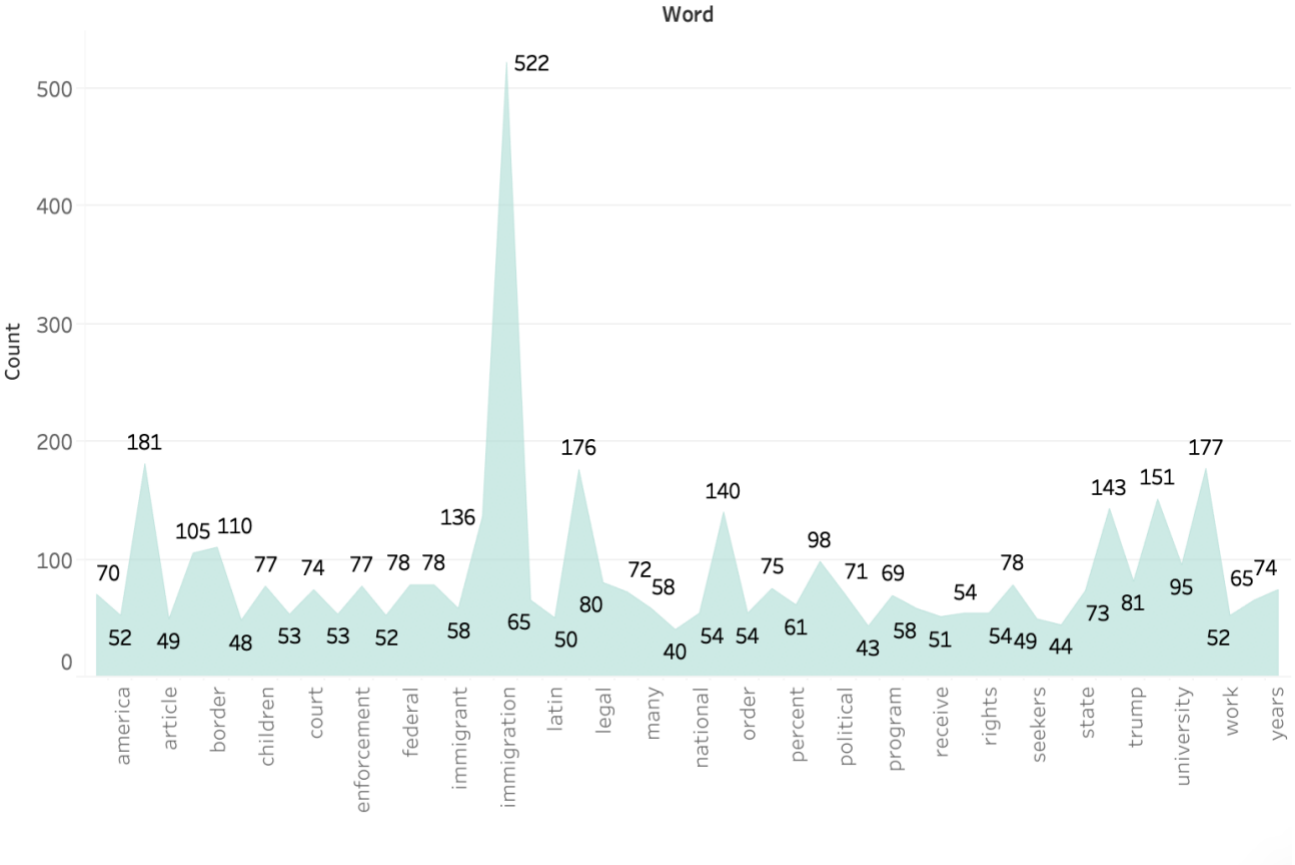
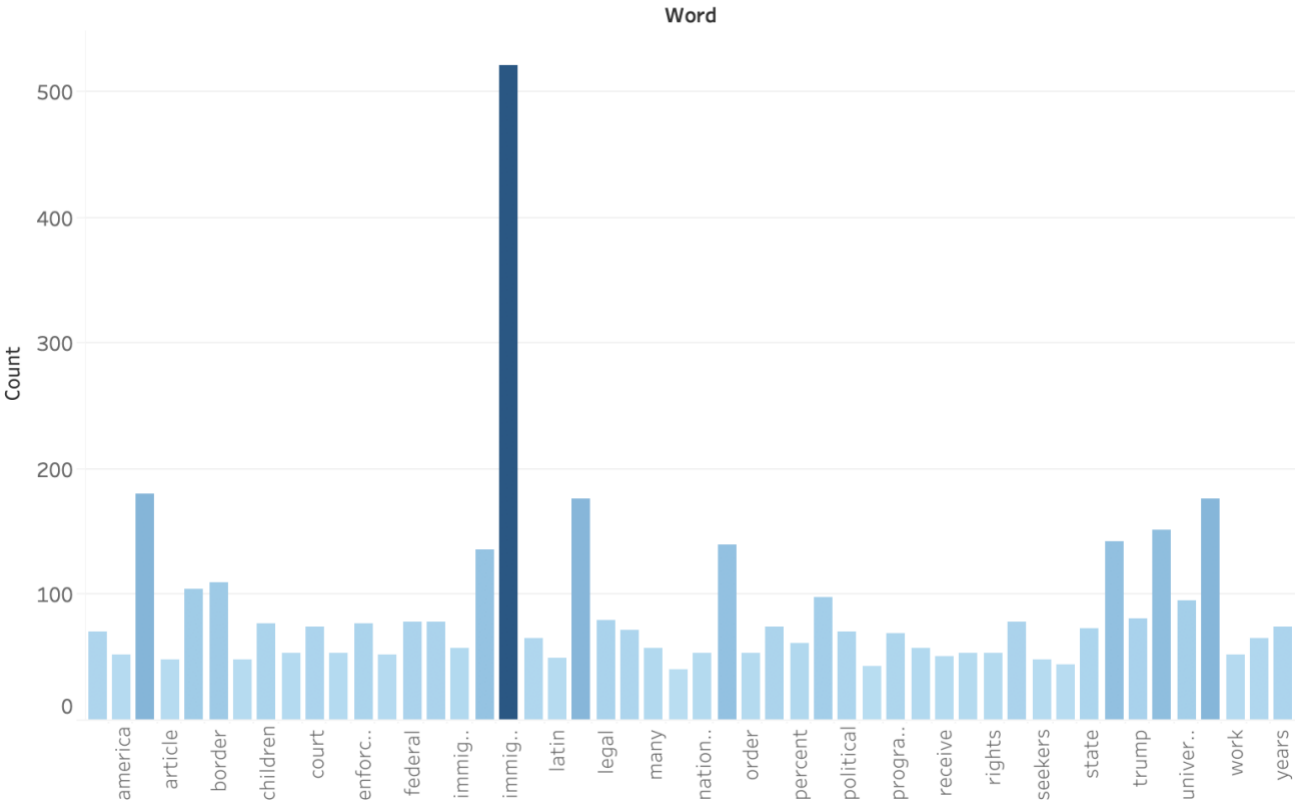
Word Cloud



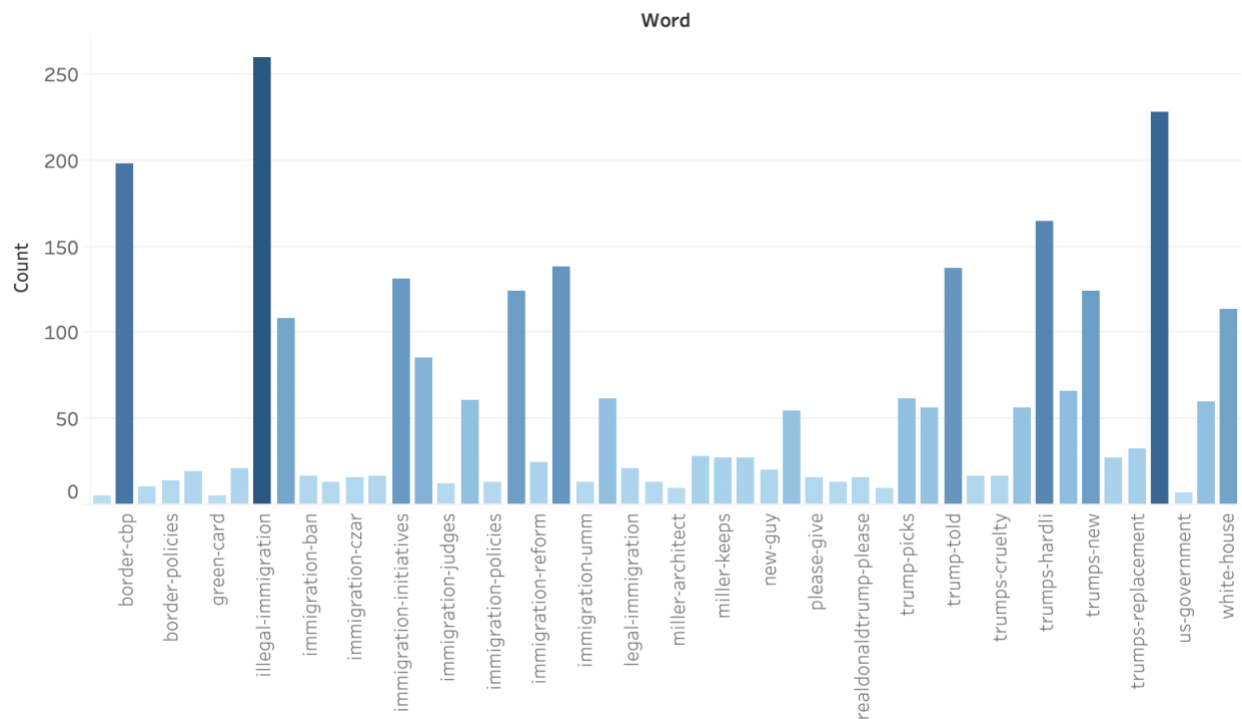
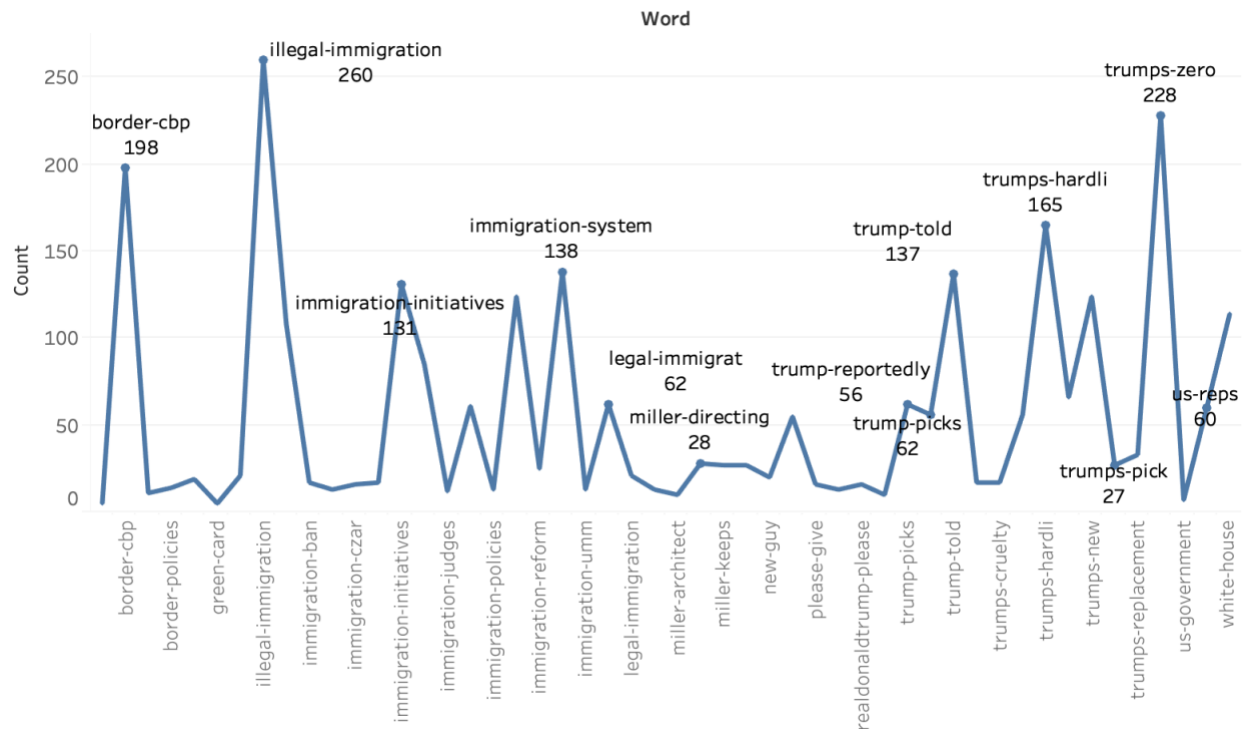
Line, Bar and Area Plot



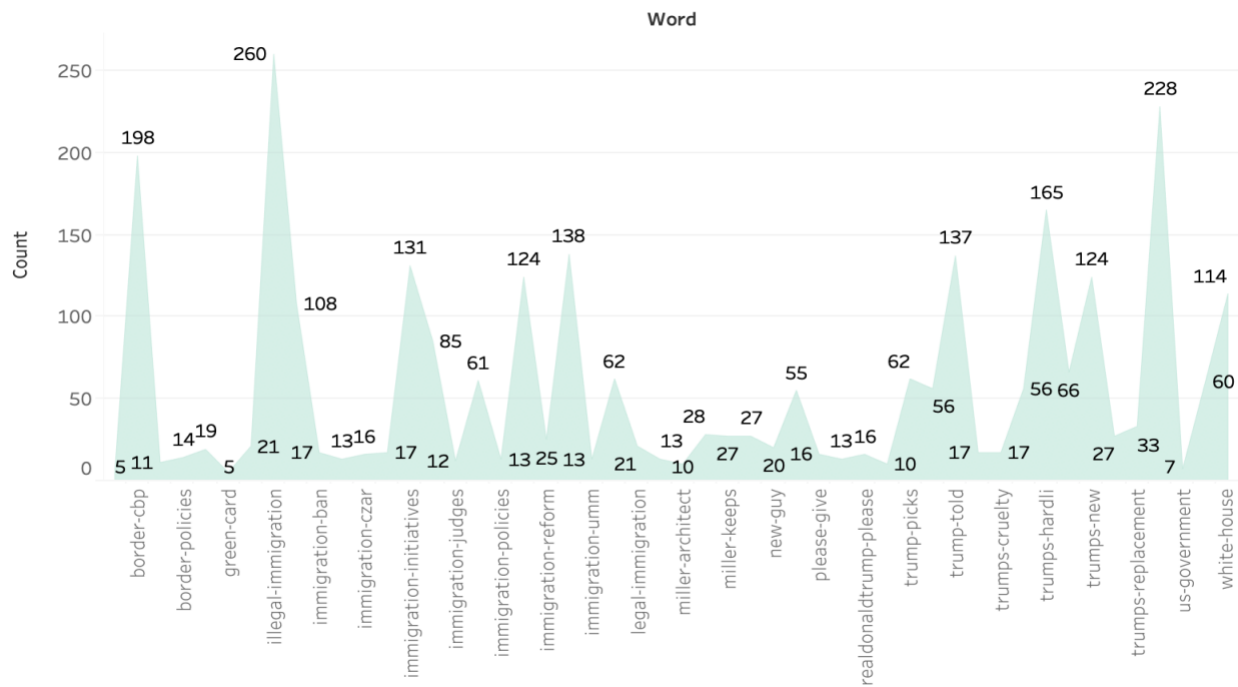
TEAM MEMBERS: - PRADEEP KUMAR JOSHI, ABHISHEK GOSWAMI



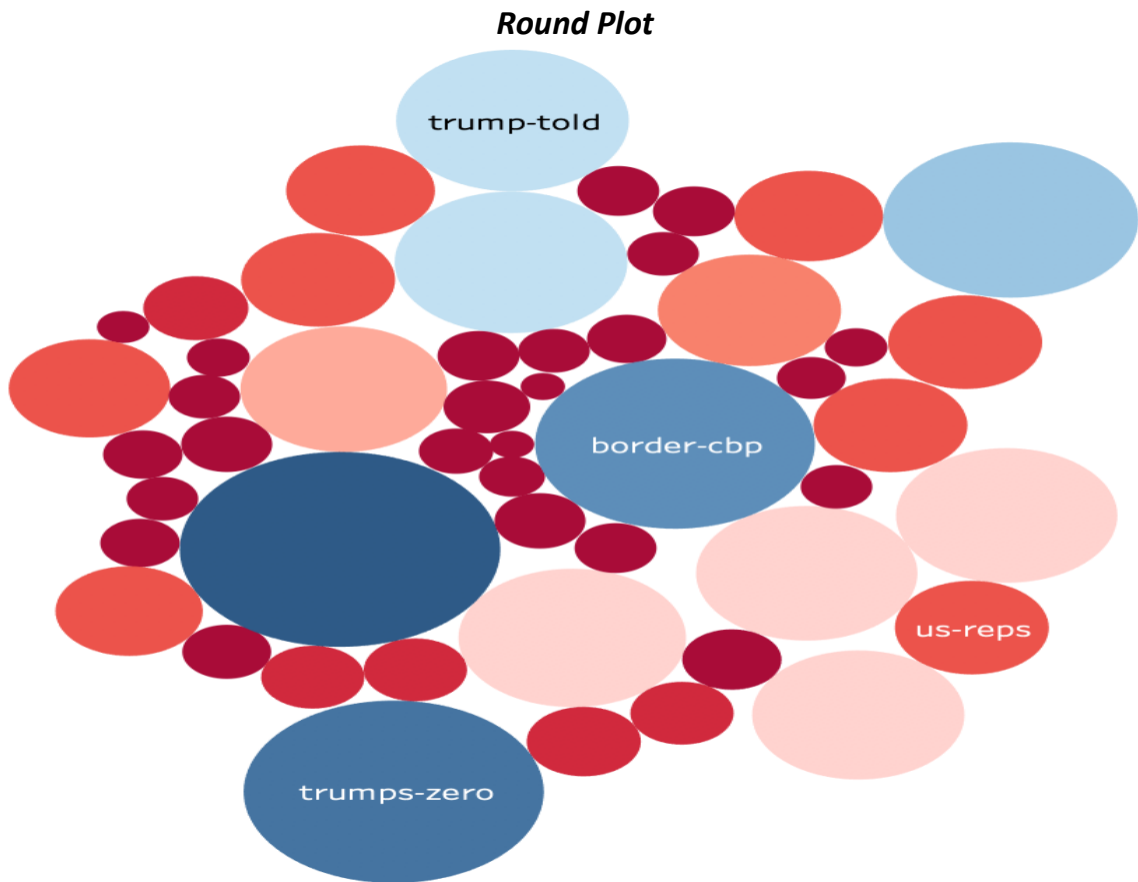
Line, Bar and Area Plot



TEAM MEMBERS: - PRADEEP KUMAR JOSHI, ABHISHEK GOSWAMI



Round Plot

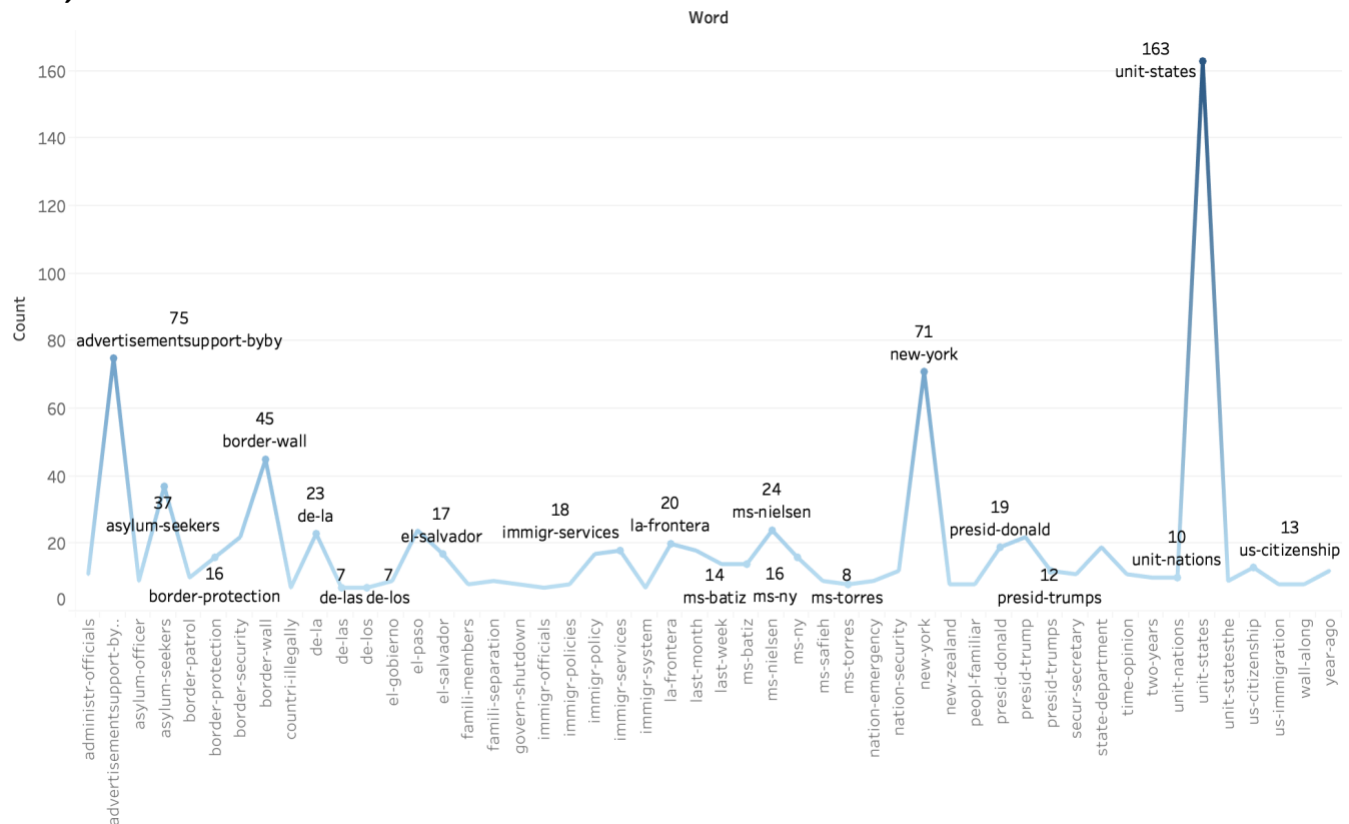


New York Times Top 50 Word Co-occurrence

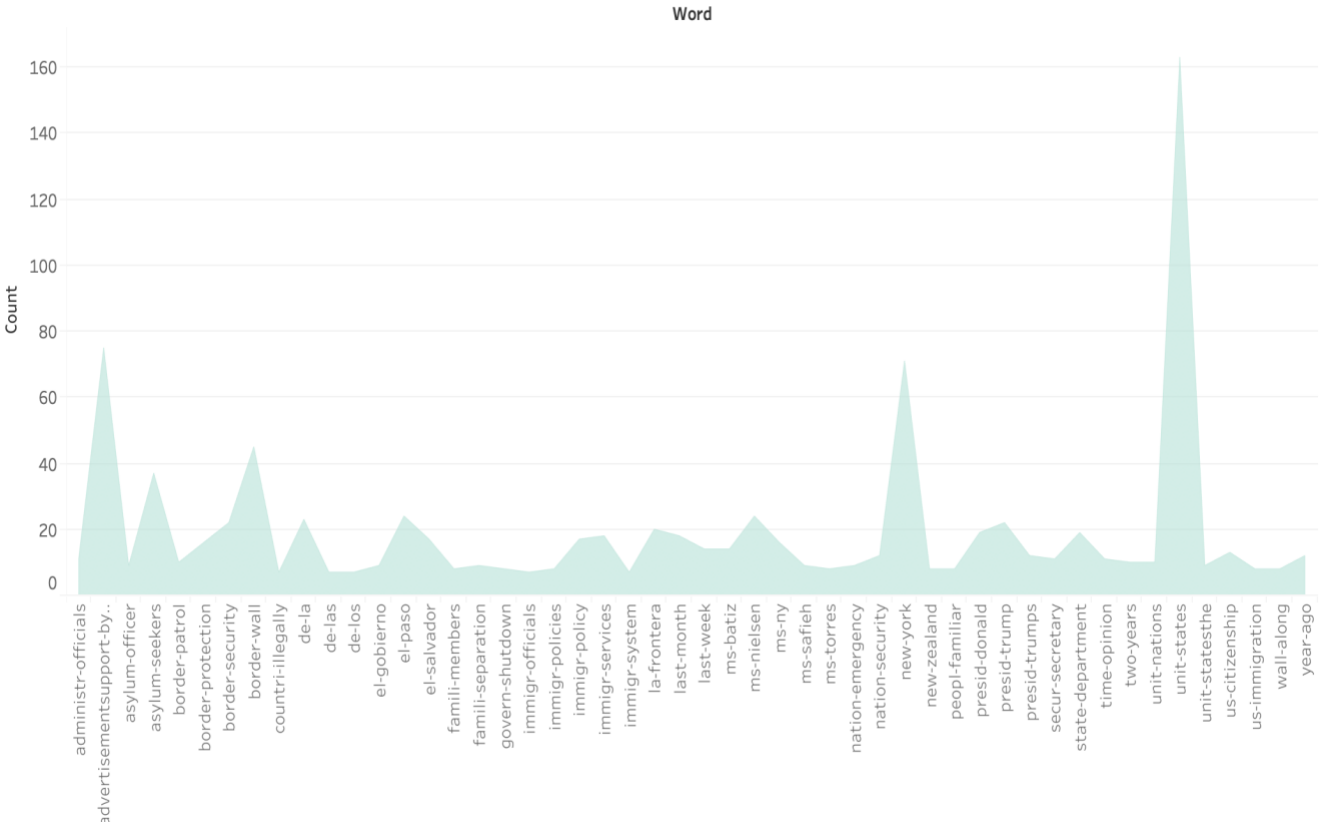
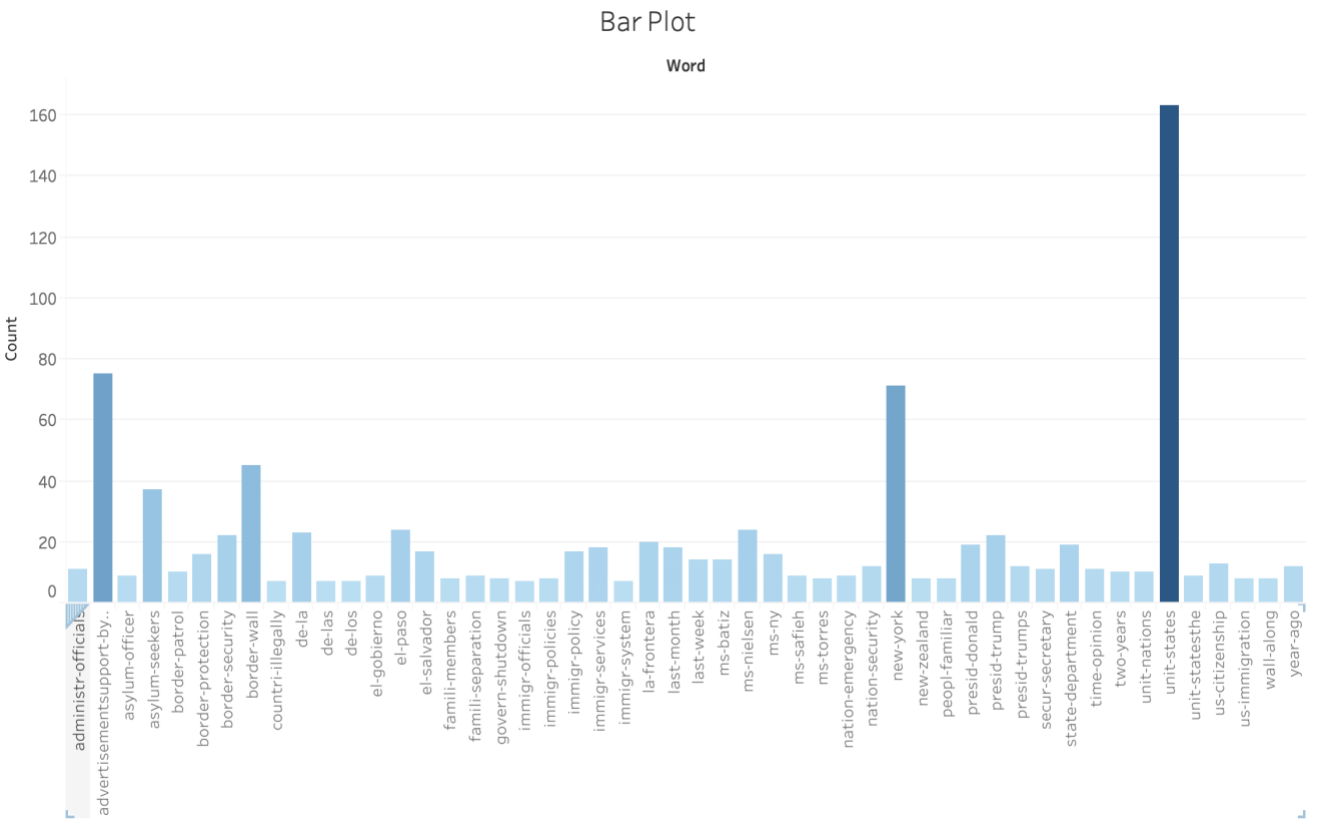
Word Cloud



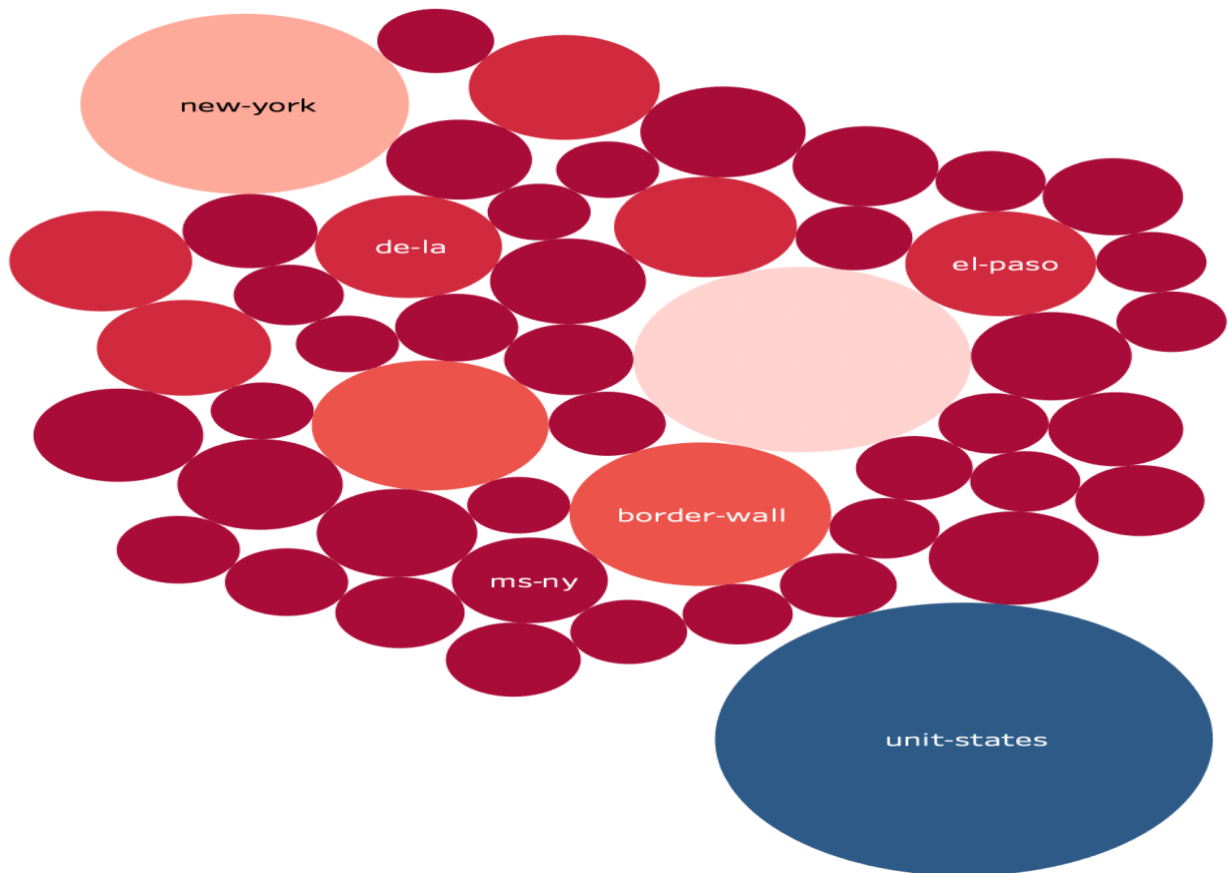
Line, Bar and Area Plot



TEAM MEMBERS: - PRADEEP KUMAR JOSHI, ABHISHEK GOSWAMI

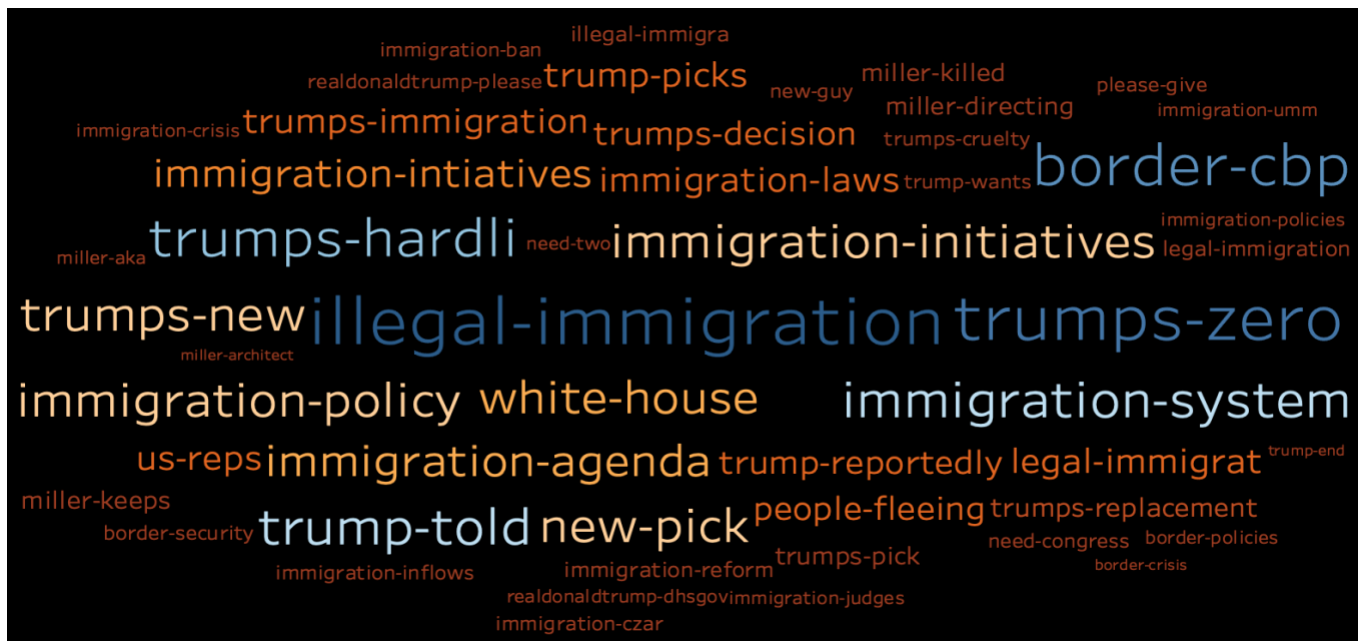


Round Plot

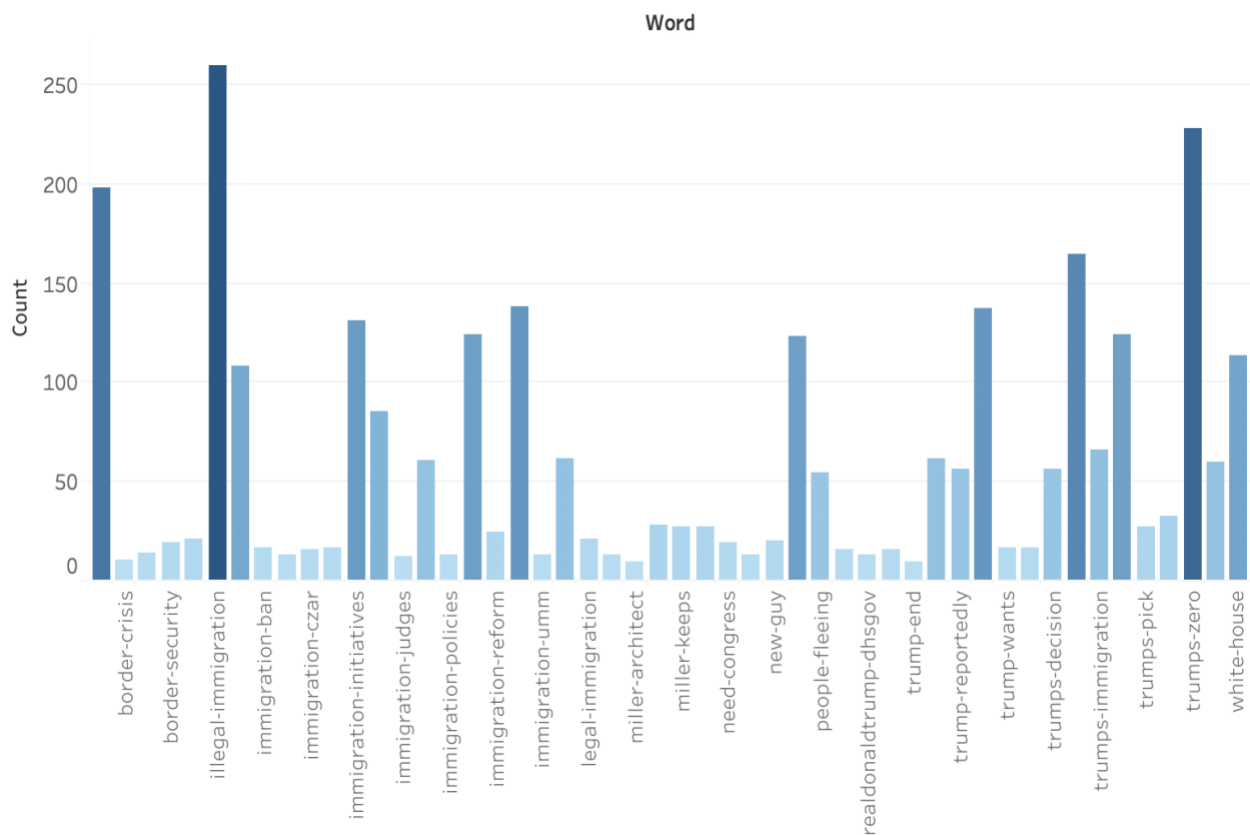
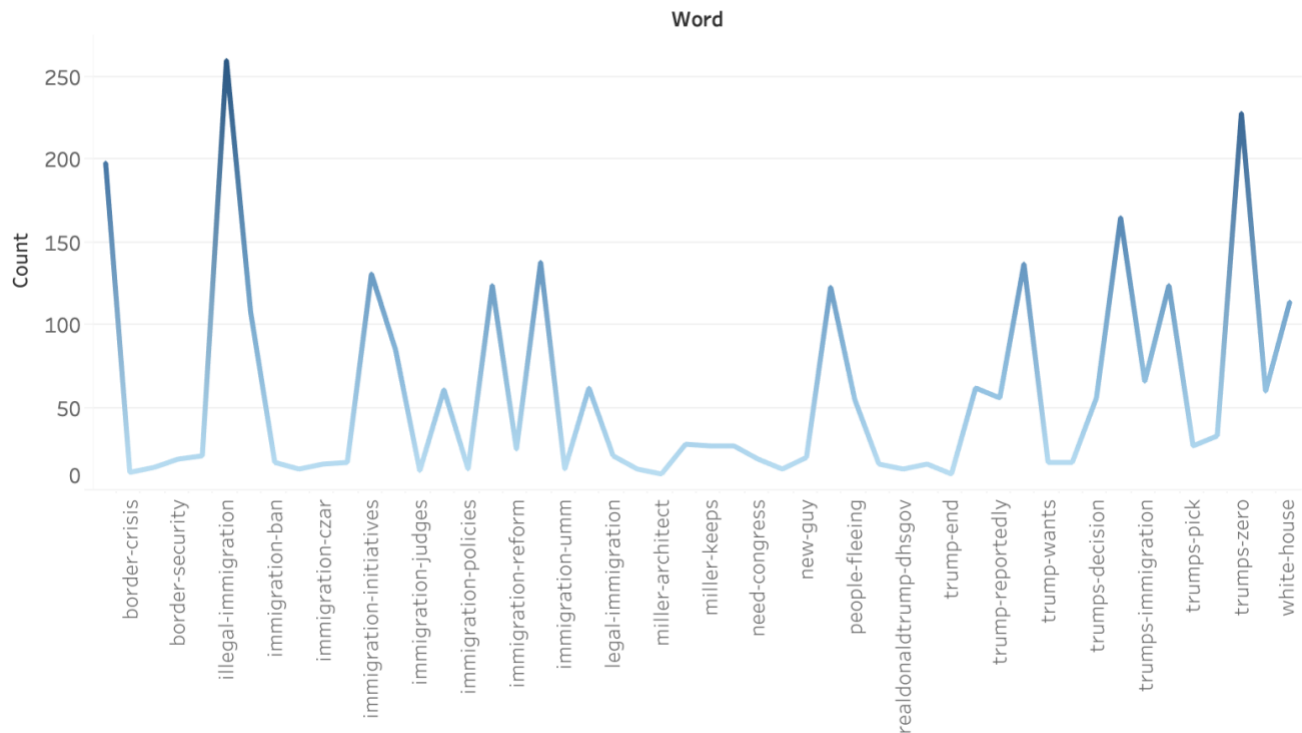


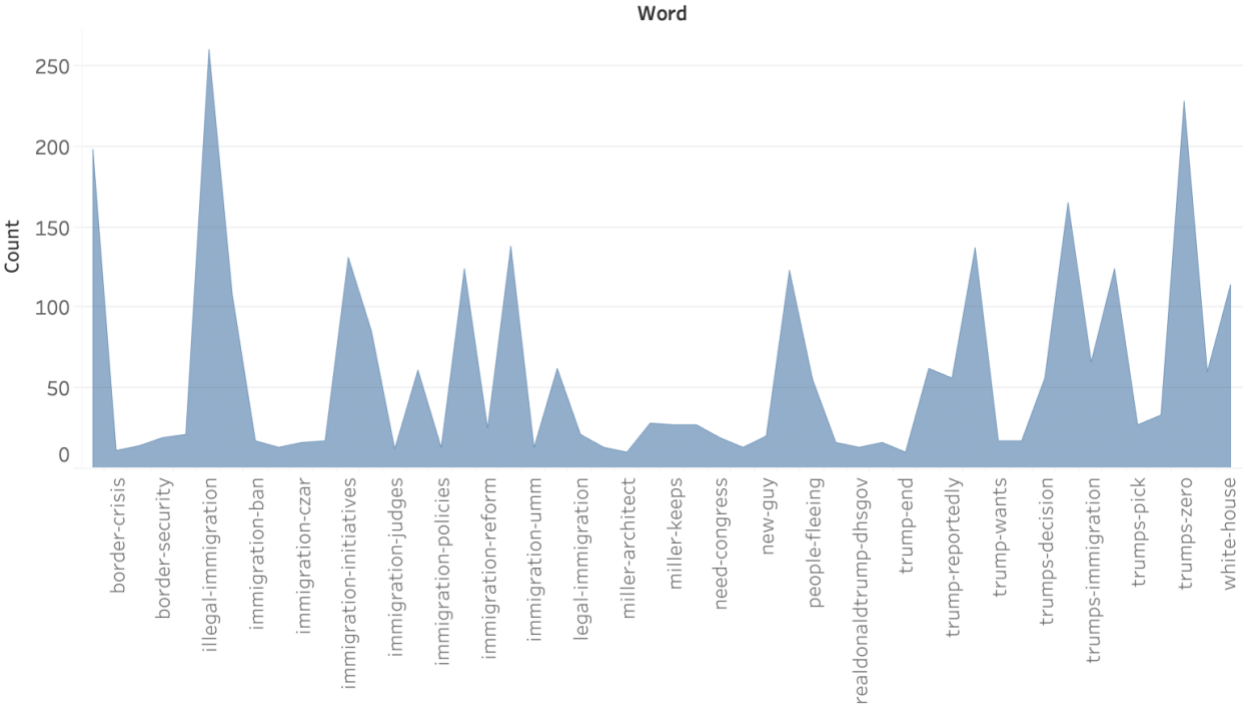
Common Crawl Top 50 Word Co-occurrence

Word Cloud

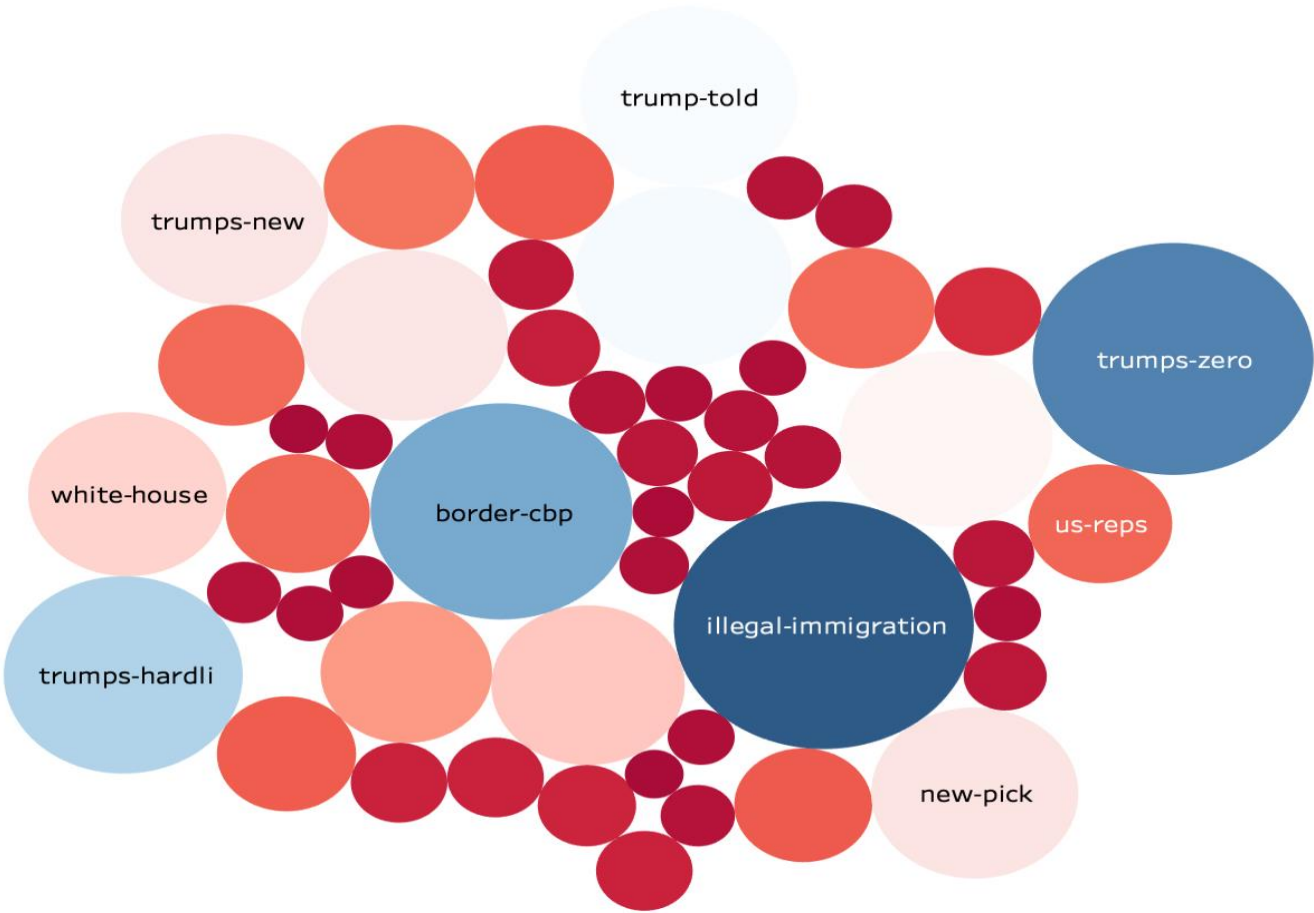


Line, Bar and Area Plot





Round Plot



→ Directory Structure
goswami5Lab2.zip

- **Report.pdf**
- **Video.mp4**
- **README.txt**
- **part1** (folder)
 - **Code** (folder) -> all python and R codes for fetching Twitter, NY Times and Common Crawl Data.
 - **Data** (folder)
 - Twitter (folder) -> all tweets are inside *"tweets.txt"* file
 - NYT (folder) -> all NYT articles are inside *"nytimes.txt"*. Subtopic files are also submitted.
 - Commoncrawl (folder) -> all Common Crawl articles are inside *"common_crawl.txt"*
- **part2** (folder)
 - pg345.txt, mapper.py and reducer.py as provided by the professor.
 - Output files (_SUCCESS and part-00000) and screenshot (Demo_Output) obtained after running the demo.
- **part3** (folder)
 - **Twitter** (folder)
 - Code (folder) -> Contains all the mapper and reducer scripts for word count and word co-occurrence.
 - Images (folder) -> Contains all the *.twbx, image files* for word count and word co-occurrence to create Twitter cloud.
 - Output (folder) -> Contains folders and output files for word count and word co-occurrence for Twitter data.
 - **NYT** (folder)
 - Code (folder) -> Contains all the mapper and reducer scripts for word count and word co-occurrence.
 - Images (folder) -> Contains all the *.twbx, image files* for word count and word co-occurrence to create NYT cloud.
 - Output (folder) -> Contains folders and output files for word count and word co-occurrence for New York Times data.
 - **Commoncrawl** (folder)
 - Code (folder) -> Contains all the mapper and reducer scripts for word count and word co-occurrence.
 - Images (folder) -> Contains all the *.twbx, image files* for word count and word co-occurrence to create Common Crawl cloud.

TEAM MEMBERS: - PRADEEP KUMAR JOSHI, ABHISHEK GOSWAMI

- Output (folder) - > Contains folders and output files for word count and word co-occurrence for Common Crawl data.
- **Webpage** (folder)
 - Interactive webpage file showing the results based on data collected from various sources.
<https://public.tableau.com/profile/abhishek.goswami#!/vizhome/Wordcountandcooccurrencevisualization/DICLAB2Visualization>

CONCLUSION:

The analysis and visualization of word count and word co-occurrence for Twitter, New York Times and Common Crawl is shown in the above plots. Hadoop infrastructure and AWS console was used for running the word-count and word co-occurrence using the appropriate mapper and reducer scripts. The files obtained from the reducer output were then used further in Tableau for visualization.