



IUPUI

SCHOOL OF INFORMATICS AND COMPUTING

Date: 4th December 2023

Election Pulse: Unveiling Voter Sentiments for 2024 US Election Forecasting

GROUP – 19

Hemanth Kumar Reddy Gunnam
Supraja Pericherla
Durga Pravallika Kuchipudi

INTRODUCTION

Aim and Objective:

This project aims to predict the outcome of the 2024 elections using vast amounts of data available online and advanced sentiment analysis techniques. By analyzing public opinion and sentiment from various online sources, we seek to gain insights into voter preferences and trends.

Understanding Public Opinion:

- In the modern political landscape, public opinion is often fragmented and rapidly changing. Our project seeks to capture this complexity through real-time sentiment analysis.
- The significance of this approach lies in its ability to provide a more accurate and nuanced understanding of voter sentiment, beyond what traditional polling methods can offer.

INTRODUCTION

The Role of Sentiment Analysis in Elections:

- Sentiment analysis allows us to process and interpret large-scale online discussions, opinions, and reactions related to political events and candidates.
- This approach is crucial for identifying emerging trends, voter concerns, and overall sentiment towards political entities, which are key indicators for election predictions.

Relevance in the Digital Age:

- In an era where social media and online platforms play a significant role in shaping public opinion, this project's approach is particularly relevant.
- It represents a step towards embracing more technologically advanced and data-driven methods in the field of political forecasting.



DATASET

Reddit Data Mining:

Subreddit Focused: Exploration centered on the 'politics' subreddit, a hub for political discourse.

Search Query: Used 'US elections 2024' to extract relevant discussions, opinions, and predictions from the subreddit.

The screenshot shows the Reddit API developer application creation interface. At the top, there's a navigation bar with links like 'GET NEW REDDIT', 'MY SUBREDDITS', 'HOME', 'POPULAR', 'ALL', 'RANDOM', 'USERS', 'FUNNY', 'ANNOUNCEMENTS', 'REDDIT.COM', 'ASKREDDIT', 'GAMING', 'PICS', 'MOVIES', 'MILDLYINTERESTING', 'TODAYILEARNED', 'WORLDNEWS', 'EXPLAINLIKEMFIVE', 'NEWS', 'VIDEOS', 'DIY', 'AWW', 'TWO', and 'EDIT'. Below the navigation, there are tabs for 'PREFERENCES', 'options', 'apps', 'RSS feeds', 'friends', 'blocked', 'password/email', and 'delete'. The main area is titled 'developed applications' and shows a single entry for 'Testing_api'. The entry includes fields for 'name' (Testing_api), 'description' (Just testing), 'about url' (empty), 'redirect uri' (http://localhost:8080), and 'secret' (redacted). It also lists 'developers' (Amazing-Wedding-7869) and provides options to 'remove' or 'add developer'. Below this, there's a 'create application' section with a note about agreeing to terms and conditions. It has a 'name' field set to 'Testing_api', a 'description' field (Just testing), and a 'redirect uri' field (http://localhost:8080). It also includes a 'web app' radio button (unchecked), an 'installed app' radio button (unchecked), and a 'script' radio button (checked). There are also fields for 'about url' and 'redirect uri', and a 'reCAPTCHA' checkbox at the bottom.

Fig 1: Reddit API

DATASET

Google API Extraction:

Broad Scope: Queries cast a wide net to capture diverse perspectives and information on the subject.

Specific Queries:

- 'US elections 2024' to gather general content related to the upcoming elections.
- '2024 US Presidential Election Candidates' to obtain information on the individuals participating in the race.
- 'US Election 2024 Predictions and Polls' to collect data on forecasts and public opinion polls.

The screenshot shows the Google Cloud Platform API & Services Credentials page. It displays a table for API Keys with one entry: "API key_1" created on Dec 3, 2023. A red arrow points to this entry. Below the table, there's a section for OAuth 2.0 Client IDs and Service Accounts, both of which are currently empty.

Fig 2: Google API

The screenshot shows the Programmable Search Engine creation page. It displays a success message: "Your new search engine has been created". Below it, there's a code snippet for a search engine script, with a red box highlighting the URL and a blue arrow pointing to the "Customize" button.

Fig 3: Google API search engine

Sources used for retrieving the data

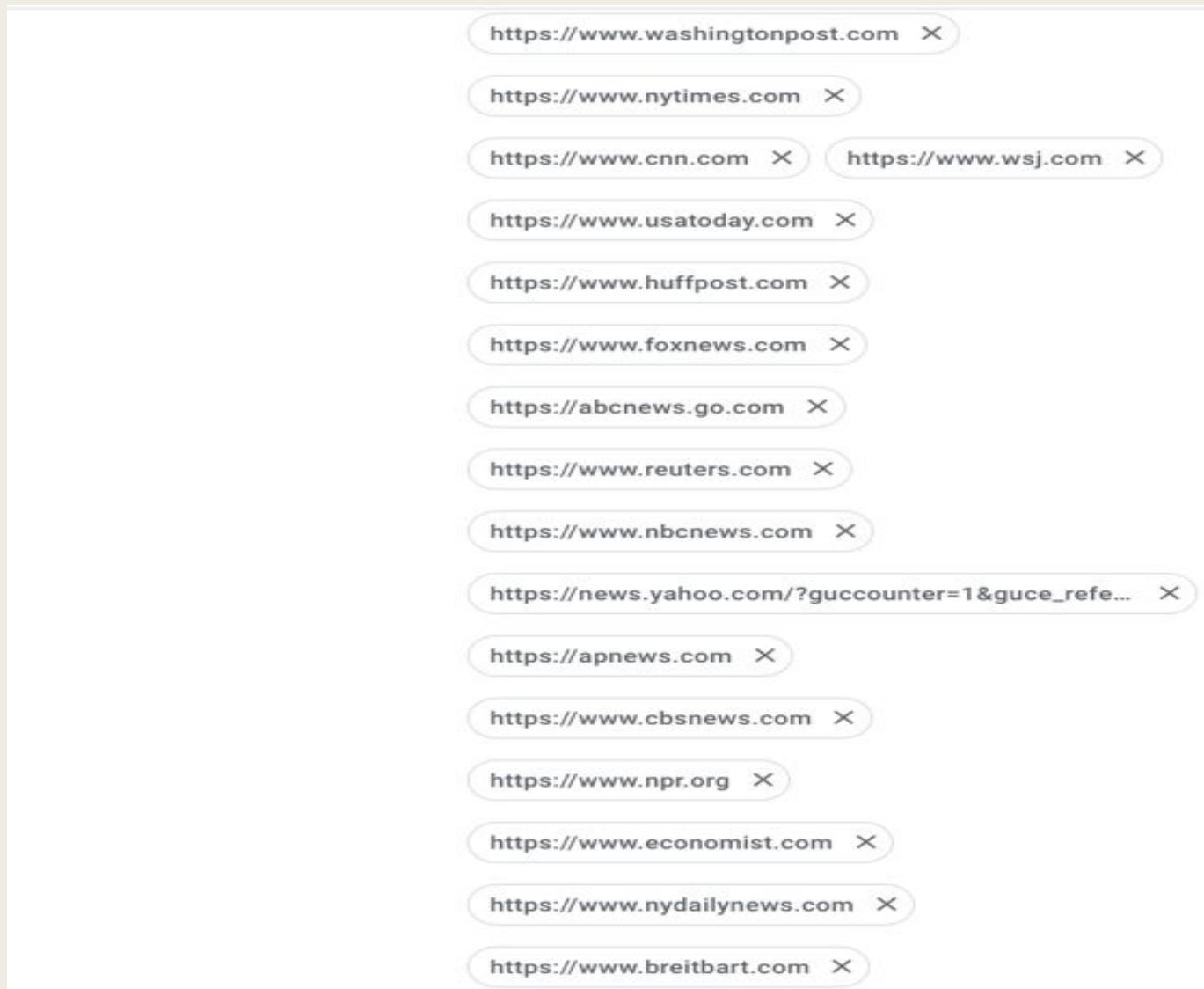


Fig 4: Source links

Comprehensive Methodology

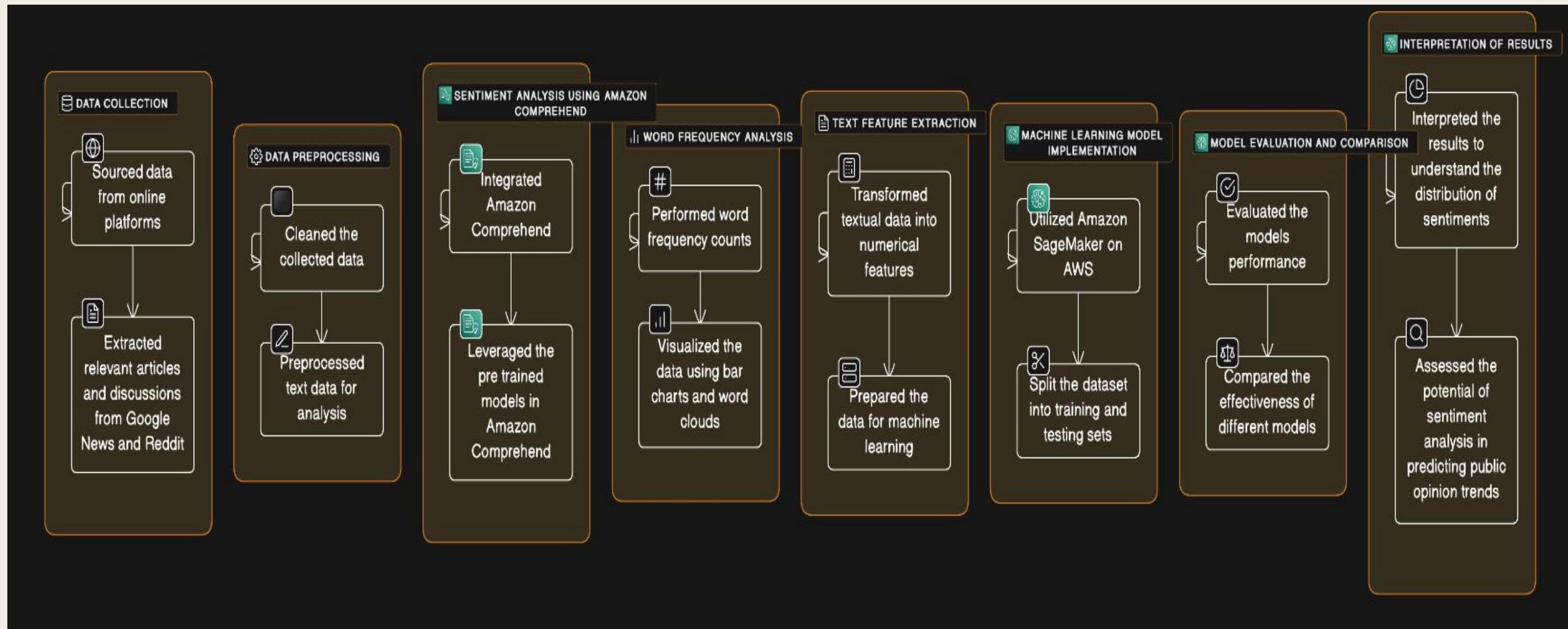


Fig 5: Methodology

Cloud Architecture

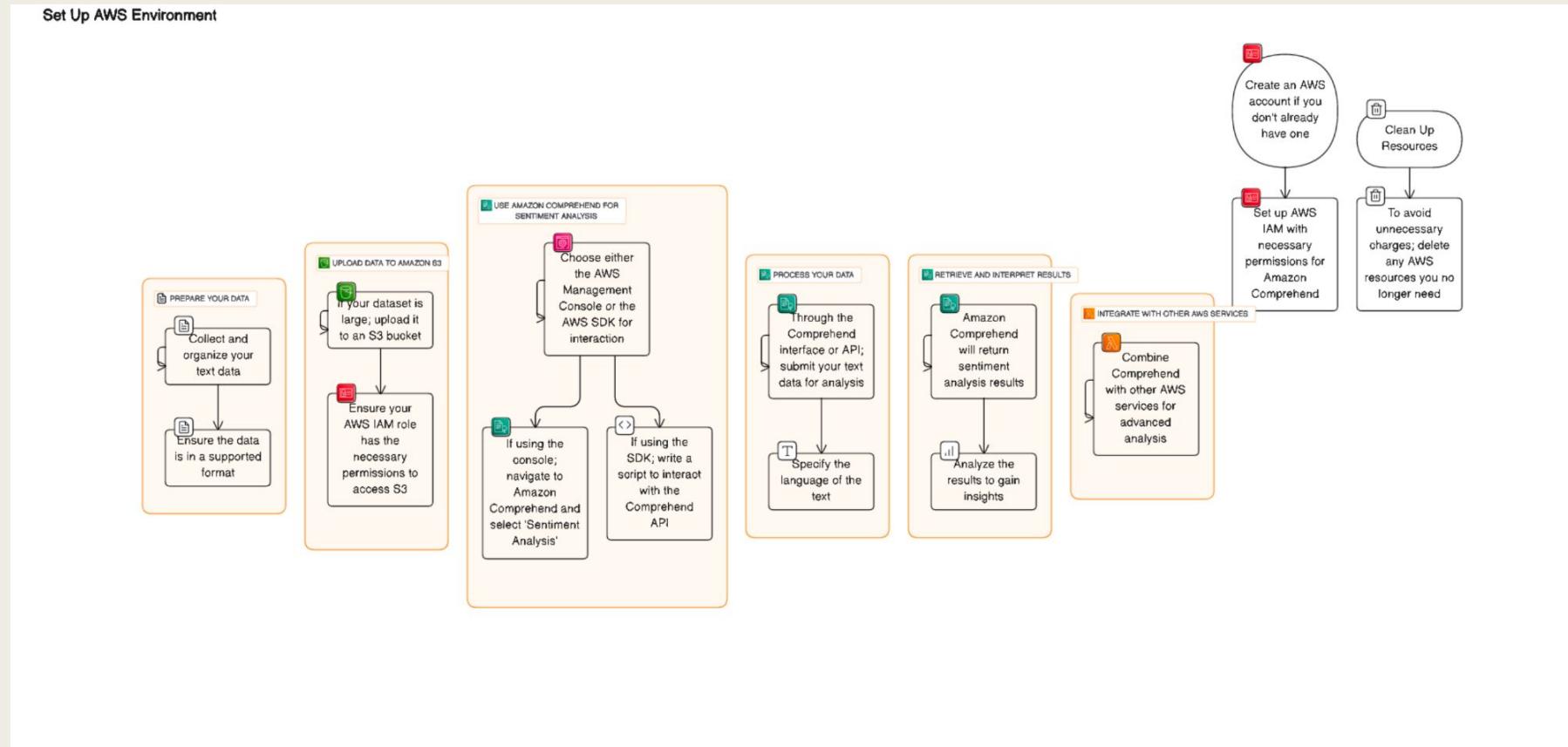


Fig 6: Cloud architecture

Results For google

```
e": "final_cleaned_google_news_final.csv", "Line": 2, "Sentiment": "NEUTRAL", "SentimentScore": {"Mixed": 4.079601239936892e-06, "Negative": 0.02554393745958805, "Neutral": 0.9653922319412231, "Positive": 0.059753268957138} }
e": "final_cleaned_google_news_final.csv", "Line": 4, "Sentiment": "NEUTRAL", "SentimentScore": {"Mixed": 8.190345397451892e-06, "Negative": 0.040790852159261703, "Neutral": 0.937761127948761, "Positive": 43983729183674}
e": "final_cleaned_google_news_final.csv", "Line": 5, "Sentiment": "NEUTRAL", "SentimentScore": {"Mixed": 3.167533577652648e-06, "Negative": 0.016941003501415253, "Neutral": 0.9799470901489258, "Positive": 10875335708260545}
e": "final_cleaned_google_news_final.csv", "Line": 6, "Sentiment": "NEUTRAL", "SentimentScore": {"Mixed": 5.969257017568452e-06, "Negative": 0.08911120891571045, "Neutral": 0.9056525826454163, "Positive": 2303047850728035}
e": "final_cleaned_google_news_final.csv", "Line": 7, "Sentiment": "NEUTRAL", "SentimentScore": {"Mixed": 4.477230049815262e-06, "Negative": 0.03060702607035637, "Neutral": 0.9633992314338684, "Positive": 989262368530035}
e": "final_cleaned_google_news_final.csv", "Line": 15, "Sentiment": "NEUTRAL", "SentimentScore": {"Mixed": 4.092977633263217e-06, "Negative": 0.04457898065447807, "Neutral": 0.9493857026100159, "Positive": 03123614564538}
e": "final_cleaned_google_news_final.csv", "Line": 16, "Sentiment": "NEUTRAL", "SentimentScore": {"Mixed": 3.827872205874883e-06, "Negative": 0.029722079634666443, "Neutral": 0.9663686156272888, "Positive": 905481891706586}
e": "final_cleaned_google_news_final.csv", "Line": 17, "Sentiment": "NEUTRAL", "SentimentScore": {"Mixed": 7.84723033575574e-06, "Negative": 0.050567448139190674, "Neutral": 0.9267773032188416, "Positive": 647405043244362}
e": "final_cleaned_google_news_final.csv", "Line": 18, "Sentiment": "NEUTRAL", "SentimentScore": {"Mixed": 5.062482614448527e-06, "Negative": 0.028772836551070213, "Neutral": 0.9647749066352844, "Positive": 4471573568888056}
e": "final_cleaned_google_news_final.csv", "Line": 19, "Sentiment": "NEUTRAL", "SentimentScore": {"Mixed": 4.793307198269758e-06, "Negative": 0.044693607836961746, "Neutral": 0.9463757872581482, "Positive": 92581976951686}
e": "final_cleaned_google_news_final.csv", "Line": 21, "Sentiment": "NEUTRAL", "SentimentScore": {"Mixed": 3.69299300473358e-06, "Negative": 0.022556137293577194, "Neutral": 0.9746296405792236, "Positive": 905481891706586}
```

Fig 7: Sentiment scores document

Results

For google

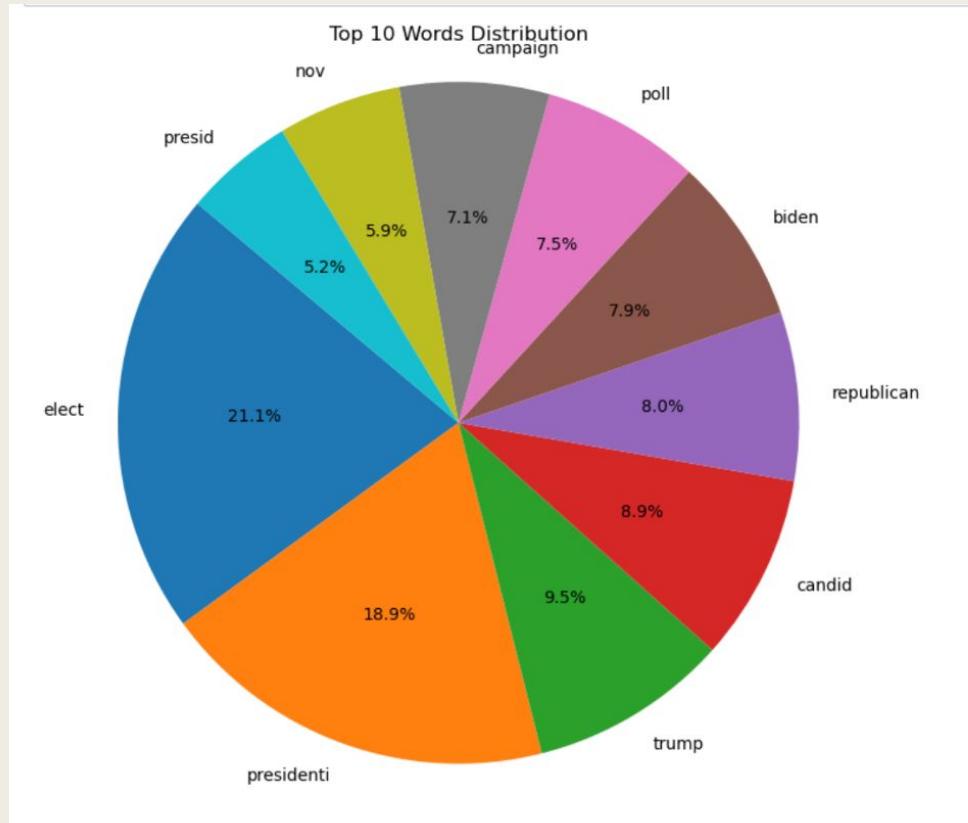


Fig 8: Pie chart of top 10 words

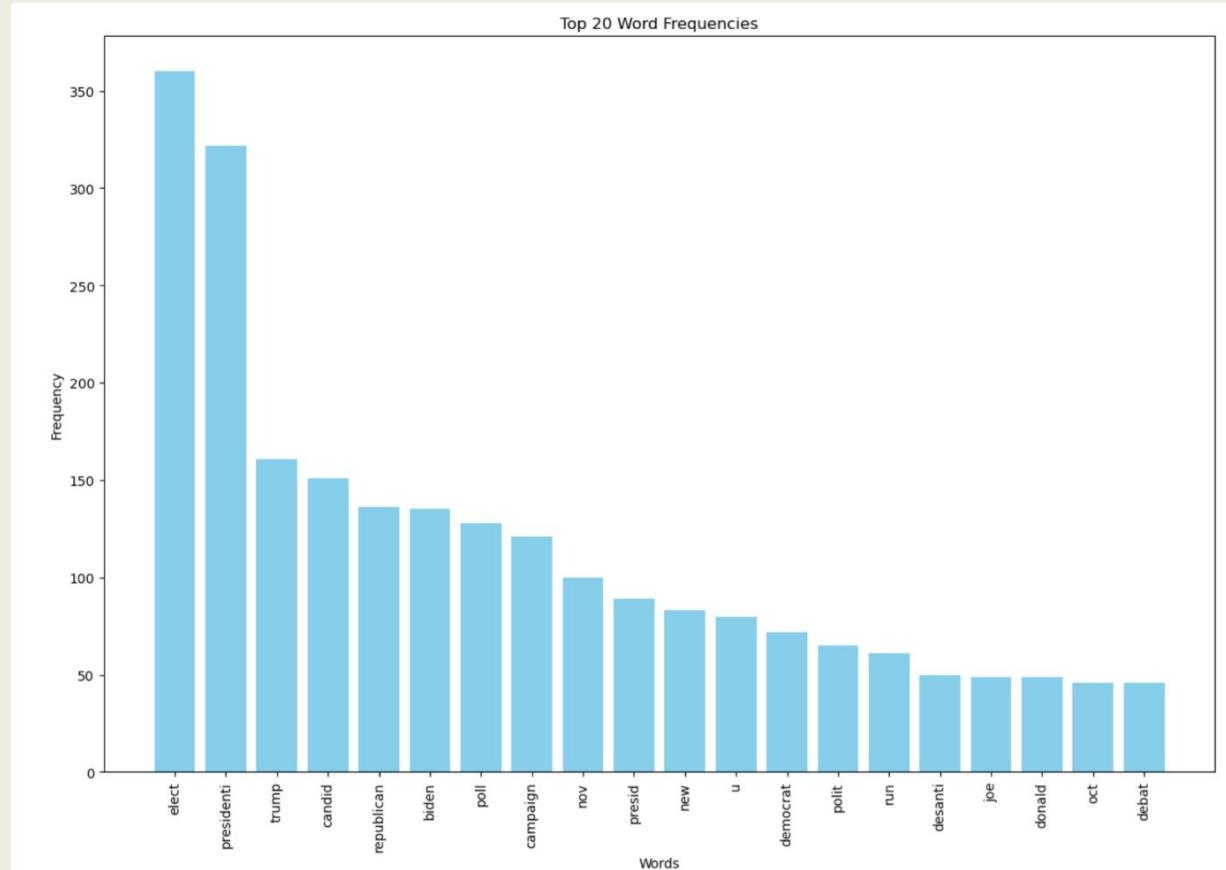


Fig 9: Bar graph showing top 20 word frequencies

Results



Reddit

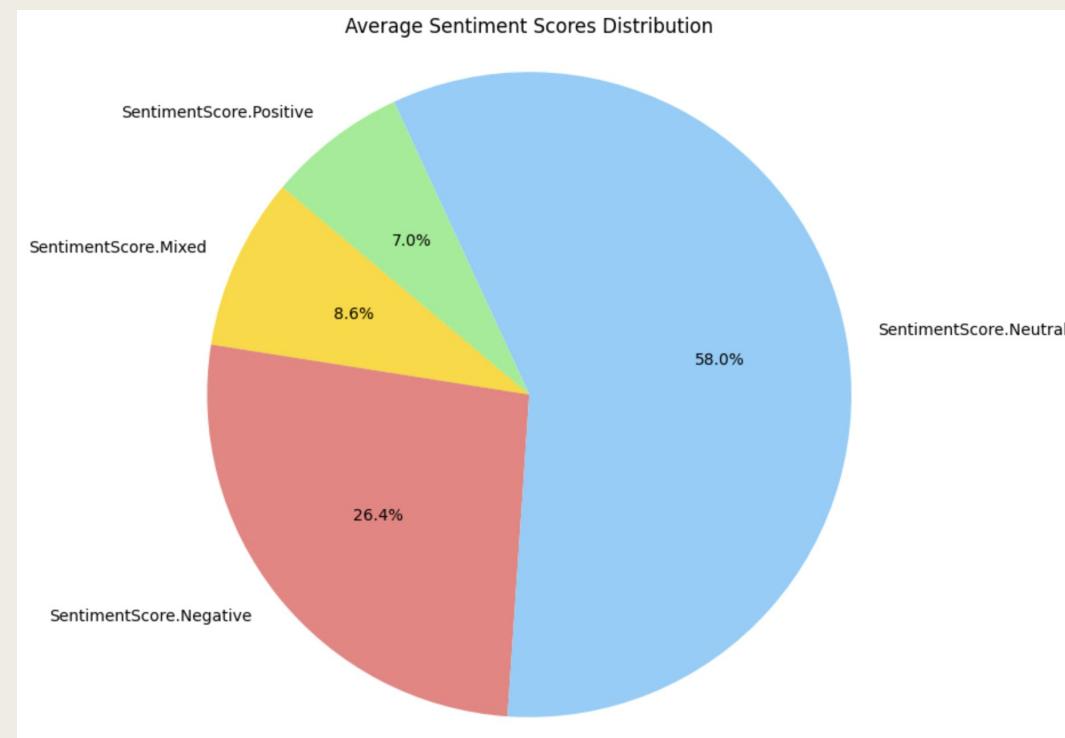


Fig 10: Average sentiment scores distribution from Reddit

Google

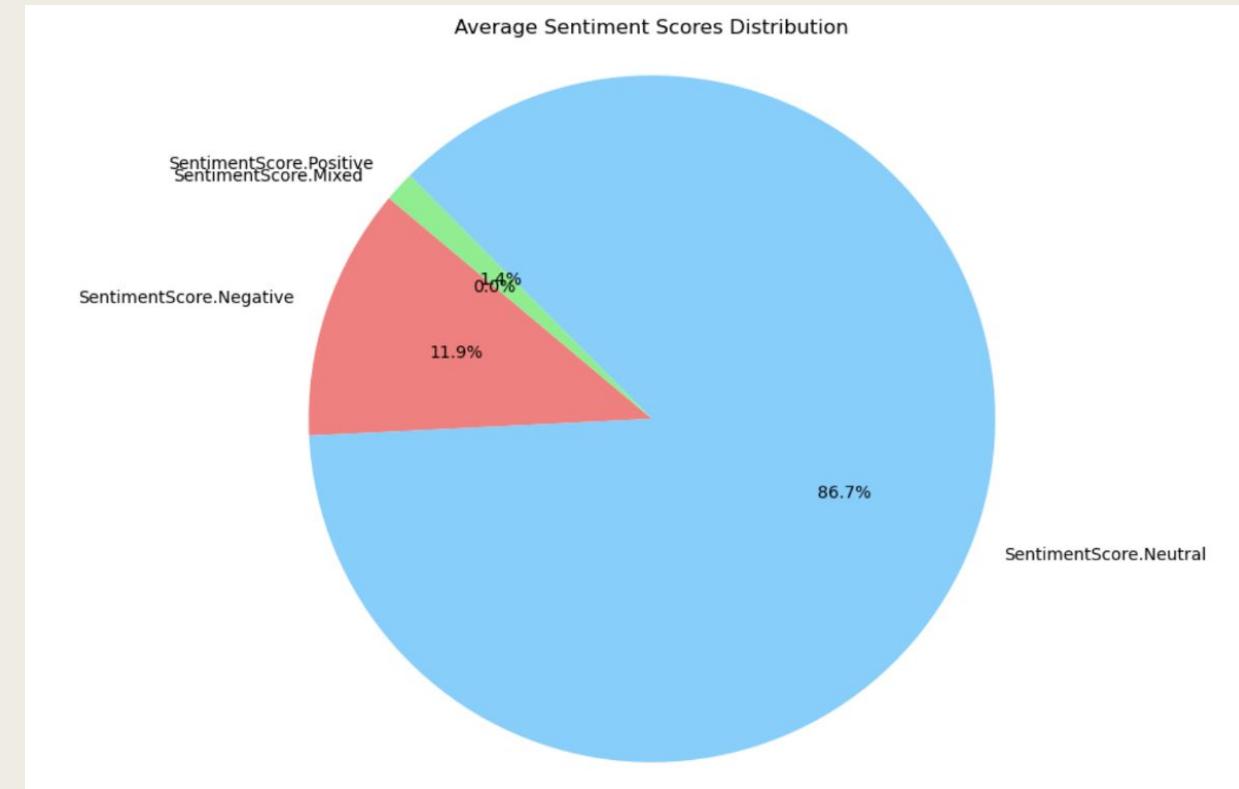


Fig 11: Average sentiment scores distribution from google

Results



IUPUI

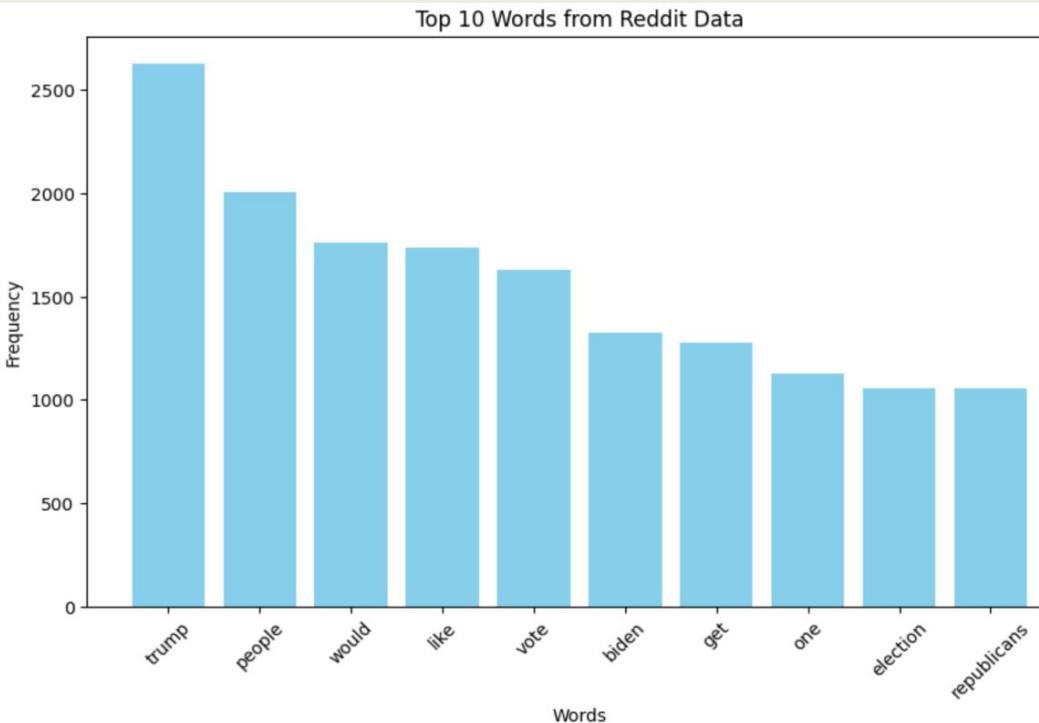
SCHOOL OF INFORMATICS AND COMPUTING

Word Rank and Word frequencies from google and reddit

Reddit:

	Word	Frequency	Rank
0	donald	178	237.0
1	trump	2748	1.0
2	iowa	21	1869.0
3	accuses	3	6586.0
4	biden	1357	6.0

Fig 12: Word rank of Reddit



Google news:

	Word	Frequency	Rank
0	elect	360	1.0
1	presidenti	322	2.0
2	trump	161	3.0
3	candid	151	4.0
4	republican	136	5.0

Fig 13: Word rank of google

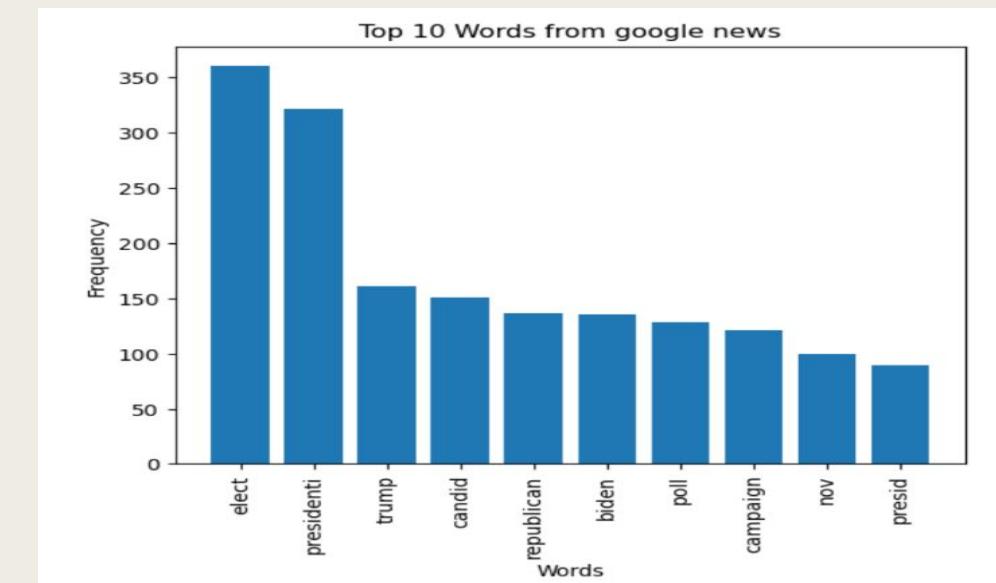


Fig 14: Top 10 words from google

Results - Word Cloud

Google

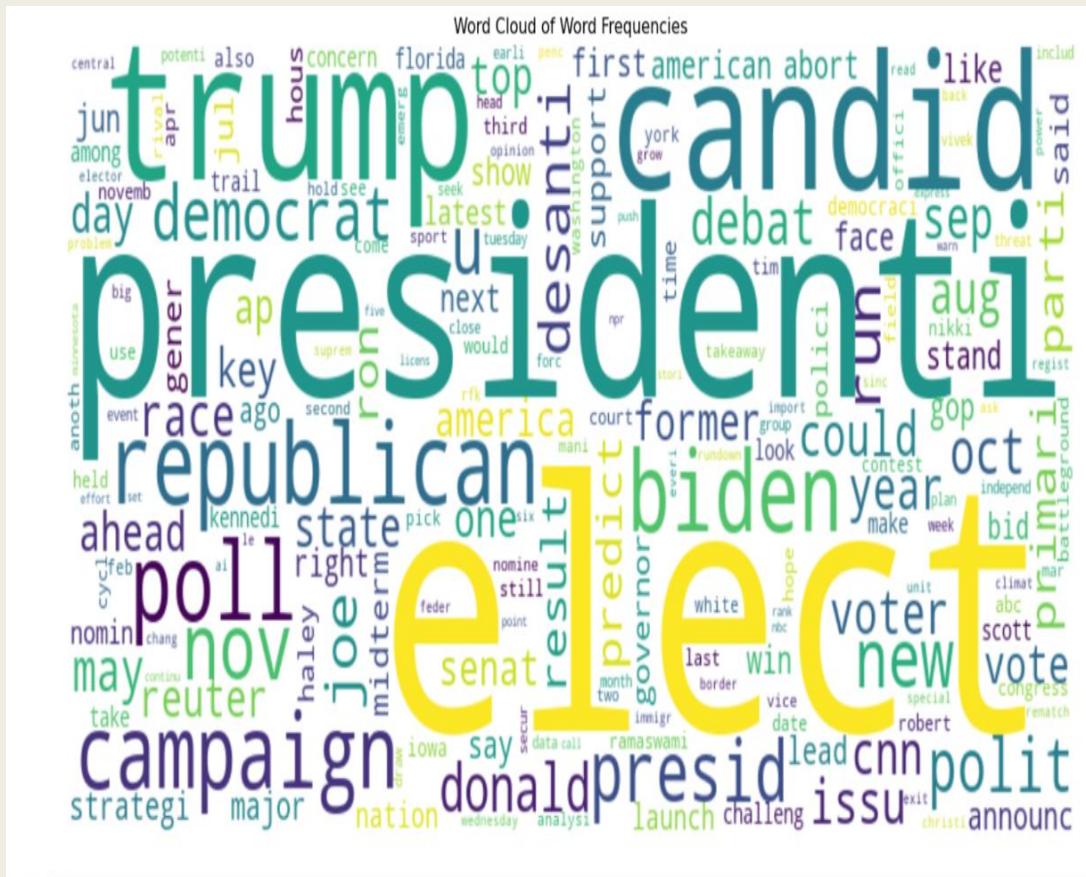


Fig 16: Word cloud from google

Reddit

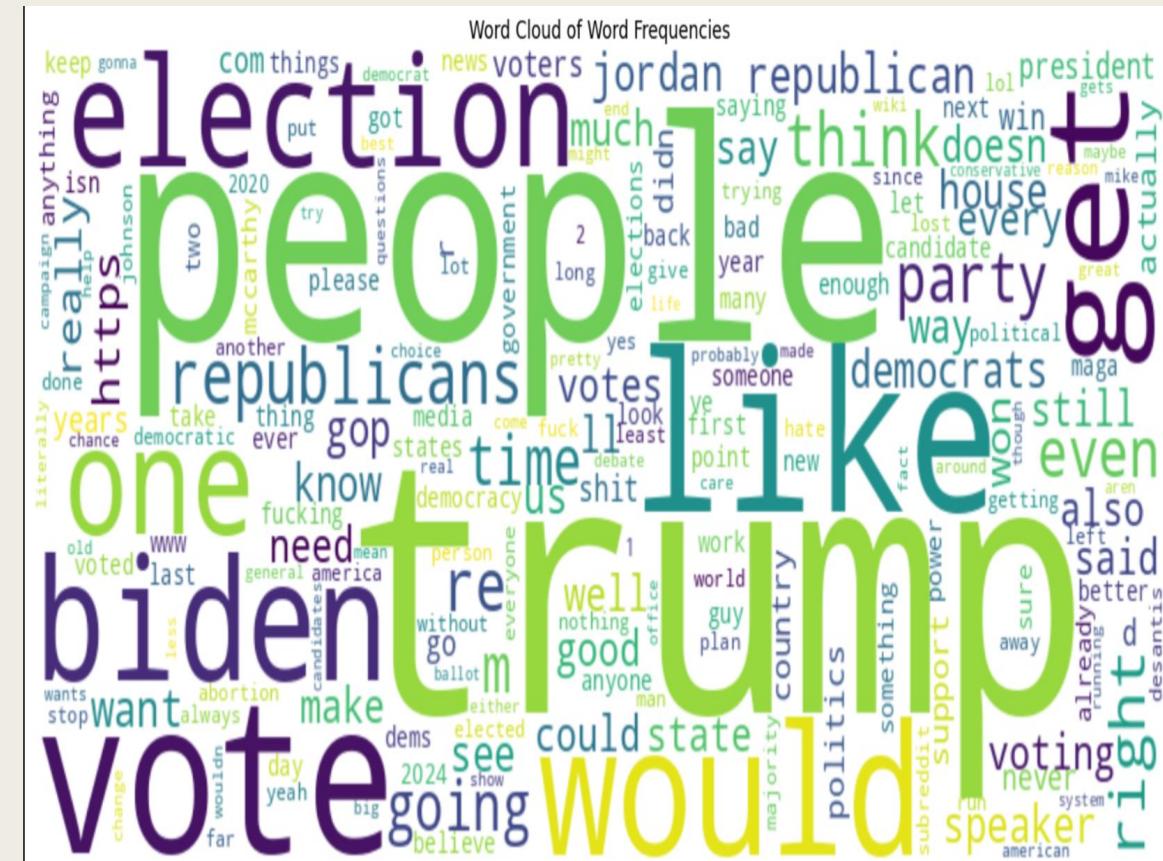
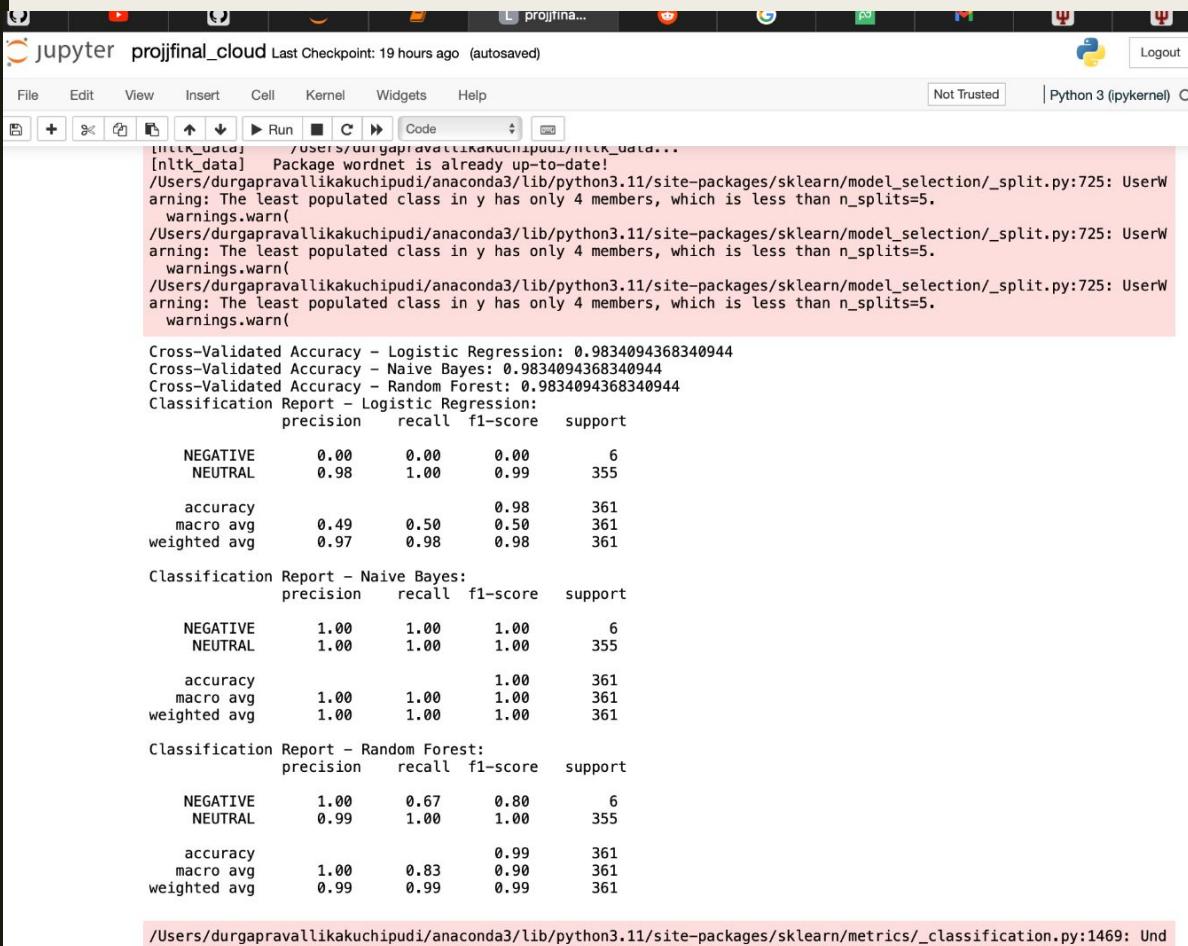


Fig 17: Word cloud from reddit

Results

For google



```

jupyter projffinal_cloud Last Checkpoint: 19 hours ago (autosaved) Logout
File Edit View Insert Cell Kernel Widgets Help Not Trusted Python 3 (ipykernel) O

[nltk_data] Package wordnet is already up-to-date!
[Users/durgapravallikakuchipudi/anaconda3/lib/python3.11/site-packages/sklearn/model_selection/_split.py:725: UserWarning: The least populated class in y has only 4 members, which is less than n_splits=5.
  warnings.warn(
[Users/durgapravallikakuchipudi/anaconda3/lib/python3.11/site-packages/sklearn/model_selection/_split.py:725: UserWarning: The least populated class in y has only 4 members, which is less than n_splits=5.
  warnings.warn(
[Users/durgapravallikakuchipudi/anaconda3/lib/python3.11/site-packages/sklearn/model_selection/_split.py:725: UserWarning: The least populated class in y has only 4 members, which is less than n_splits=5.
  warnings.warn()

Cross-Validated Accuracy - Logistic Regression: 0.9834094368340944
Cross-Validated Accuracy - Naive Bayes: 0.9834094368340944
Cross-Validated Accuracy - Random Forest: 0.9834094368340944
Classification Report - Logistic Regression:
precision recall f1-score support
NEGATIVE 0.00 0.00 0.00 6
NEUTRAL 0.98 1.00 0.99 355

accuracy 0.98 361
macro avg 0.49 0.50 0.50 361
weighted avg 0.97 0.98 0.98 361

Classification Report - Naive Bayes:
precision recall f1-score support
NEGATIVE 1.00 1.00 1.00 6
NEUTRAL 1.00 1.00 1.00 355

accuracy 1.00 361
macro avg 1.00 1.00 1.00 361
weighted avg 1.00 1.00 1.00 361

Classification Report - Random Forest:
precision recall f1-score support
NEGATIVE 1.00 0.67 0.80 6
NEUTRAL 0.99 1.00 1.00 355

accuracy 0.99 0.99 0.99 361
macro avg 1.00 0.83 0.90 361
weighted avg 0.99 0.99 0.99 361

/Users/durgapravallikakuchipudi/anaconda3/lib/python3.11/site-packages/sklearn/metrics/_classification.py:1469: Und

```

Fig 18: Cross-validated accuracy scores and classification reports of each model

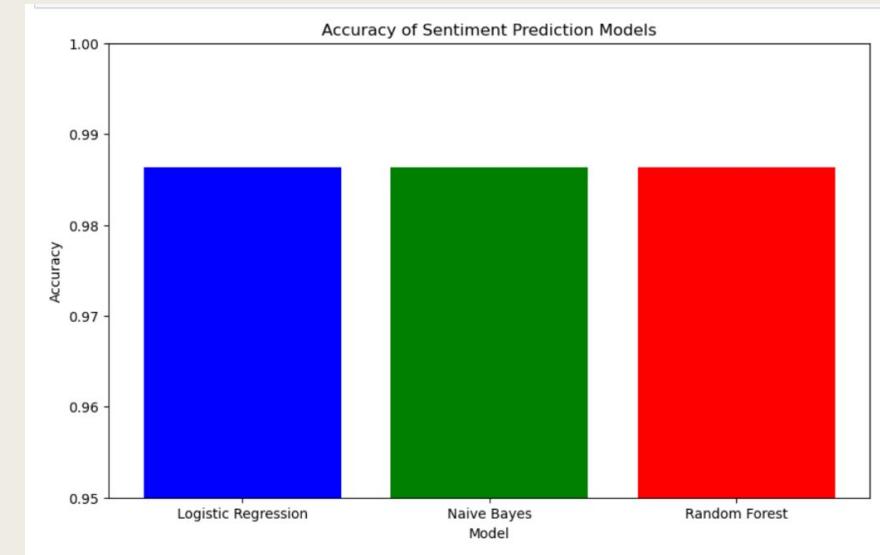


Fig 19: Accuracy score comparisons between different models

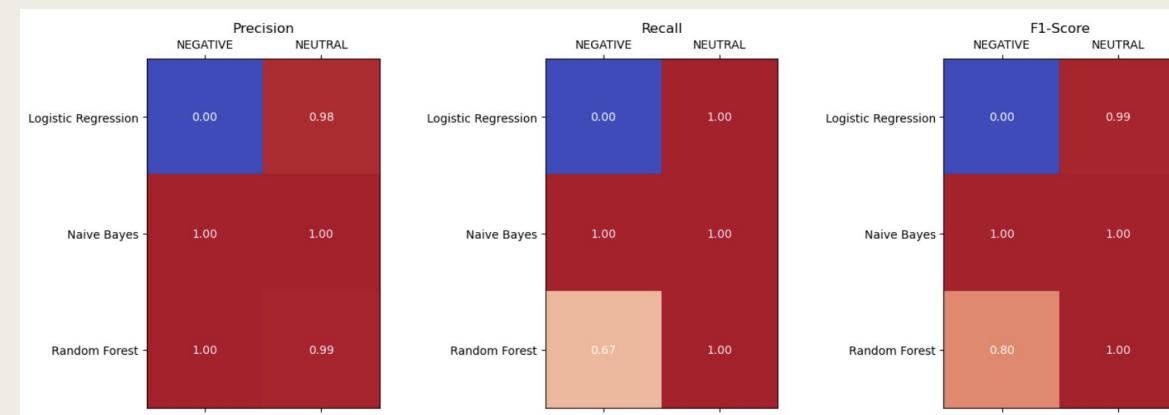


Fig 20: Classification Metrics for Different Models

Results

Reddit:

```
[nltk_data] Downloading package wordnet to /root/nltk_data...
Cross-Validated Accuracy - Logistic Regression: 0.6120182576333739
Cross-Validated Accuracy - Naive Bayes: 0.5945450810324535
Cross-Validated Accuracy - Random Forest: 0.6105921134716603
Classification Report - Logistic Regression:
    precision    recall  f1-score   support
  MIXED       1.00     0.14    0.25      770
  NEGATIVE    1.00     0.09    0.16     3152
  NEUTRAL     0.62     1.00    0.77     6720
  POSITIVE    1.00     0.09    0.17      575

  accuracy                           0.64      11217
  macro avg    0.90     0.33    0.34     11217
  weighted avg  0.77     0.64    0.53     11217

/usr/local/lib/python3.10/dist-packages/sklearn/metrics/_classification.py:1344: UndefinedMetricWarning: Precision and F-score are ill-defined and being set to
      _warn_prfaverage, modifier, msg_start, len(result)
/usr/local/lib/python3.10/dist-packages/sklearn/metrics/_classification.py:1344: UndefinedMetricWarning: Precision and F-score are ill-defined and being set to
      _warn_prfaverage, modifier, msg_start, len(result)
/usr/local/lib/python3.10/dist-packages/sklearn/metrics/_classification.py:1344: UndefinedMetricWarning: Precision and F-score are ill-defined and being set to
      _warn_prfaverage, modifier, msg_start, len(result)
Classification Report - Naive Bayes:
    precision    recall  f1-score   support
  MIXED       1.00     0.02    0.04      770
  NEGATIVE    0.72     0.04    0.08     3152
  NEUTRAL     0.61     1.00    0.76     6720
  POSITIVE    0.00     0.00    0.00      575

  accuracy                           0.61      11217
  macro avg    0.58     0.27    0.22     11217
  weighted avg  0.64     0.61    0.48     11217

Classification Report - Random Forest:
    precision    recall  f1-score   support
  MIXED       0.99     0.14    0.25      770
  NEGATIVE    0.99     0.09    0.16     3152
  NEUTRAL     0.62     1.00    0.77     6720
  POSITIVE    1.00     0.09    0.17      575

  ✓ Connected to Python 3 Google Compute Engine backend
```

Fig 21: Cross-validated accuracy scores and classification reports of each model of Reddit

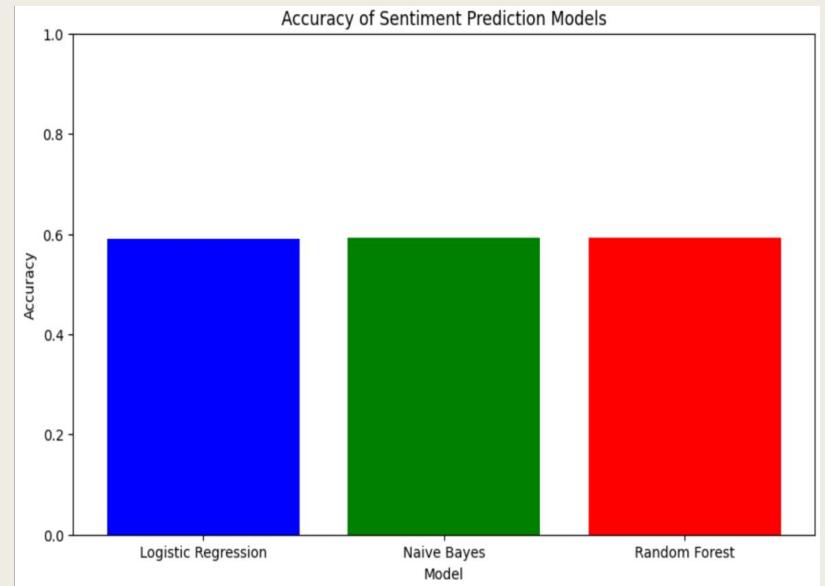


Fig 22: Accuracy score comparisons between different models

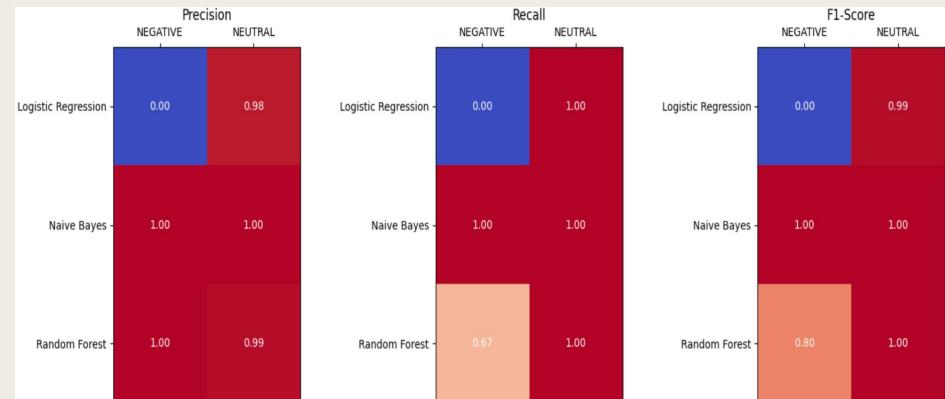


Fig 23: Classification Metrics for Different Models

Limitations

1. **Class Imbalance:** Potential issue in accurately identifying negative sentiments due to overshadowing by neutral sentiments.
2. **Model Overfitting:** Perfect scores of Naive Bayes raise concerns about memorizing training data instead of generalizing.
3. **Data Diversity:** Lack of diverse negative sentiment examples may lead to inadequate model training.
4. **Model Suitability:** Not all models equally suitable for sentiment classes, indicating a need for tailored approaches or ensemble methods.
5. **Evaluation Metrics:** Reliance on accuracy alone is insufficient; consider AUC-ROC, precision-recall curves, or confusion matrices for deeper insights.

Future Work

1. **Enhanced Data Collection:** Incorporate a wider range of real-time data from diverse sources and languages to capture a more global and comprehensive sentiment landscape.
2. **Advanced Analytics:** Employ cutting-edge machine learning techniques and natural language processing algorithms for deeper sentiment contextualization and more accurate prediction models.
3. **Temporal and Demographic Insights:** Implement time-series analysis to monitor sentiment trends and conduct demographic studies to understand sentiment variations across different voter segments.
4. **Interactive Tools:** Develop an interactive dashboard for real-time sentiment tracking and analysis, providing a dynamic tool for various stakeholders.
5. **Ethical Framework:** Strengthen ethical data handling practices, focusing on privacy, consent, and the responsible use of sentiment analysis in political contexts.

Conclusion

- **Accuracy Discrepancy:** High accuracy across Logistic Regression, Naive Bayes, and Random Forest, but performance varies across sentiment categories.
- **Precision, Recall, F1-score Differences:** Metrics reveal significant variations in correctly identifying negative sentiments among the models.
- **Naive Bayes Perfect Scores:** Raises concerns about potential overfitting, as achieving perfect scores may indicate memorization of training data.
- **Logistic Regression Failure:** Fails to identify any negative sentiments, highlighting limitations in its performance for this sentiment class.
- **Random Forest Partial Success:** Shows partial success in identifying negative sentiments, suggesting a potential for improvement or optimization.

Project link

https://github.com/Supraja-p/Cloud_computing