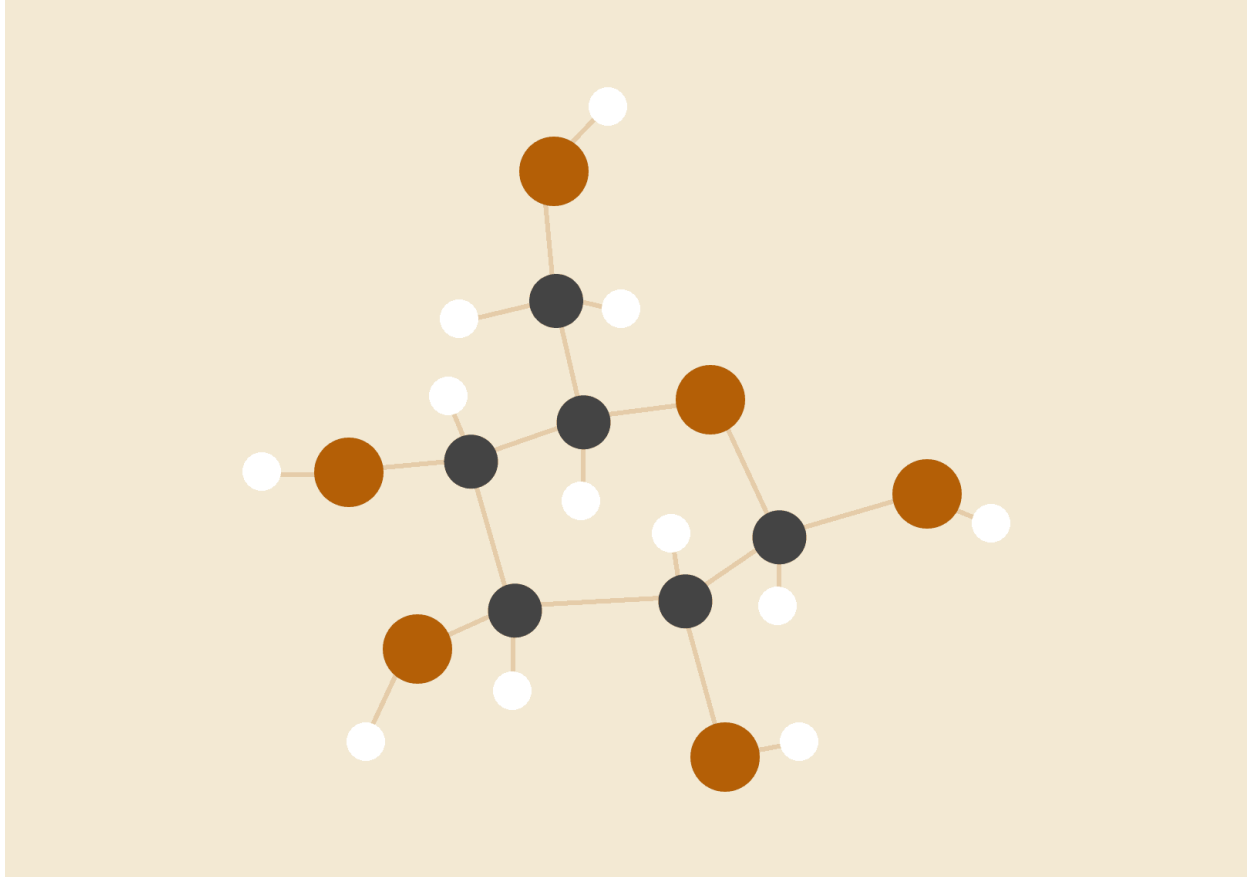


# Sentiment Analysis

*Sam Malik, Sulayman Ahmed, Hasan Ahmed*



Dec 12th, 2023

SI 206

## Goals for The Project

The primary goal of our project was to gain insights into the general public opinion on various topics by leveraging the vast user-generated content available on major social media platforms. To achieve this, we planned to work with three prominent platforms: Reddit, Twitter, and YouTube, utilizing their respective APIs to pull relevant data. The focus was on extracting a rich dataset that included popularity metrics, such as upvotes and likes, as well as timestamps of posts and tweets. This data harvesting was crucial to understanding how certain topics trended and resonated with users over time. Additionally, our project aimed to meticulously visualize the sentiment analysis results and the gathered popularity metrics using Matplotlib, a powerful plotting library in Python. The visualization aspect was intended to provide clear, insightful, and accessible representations of our findings, thereby making the analysis of public sentiment on various topics both comprehensive and intuitive.

## Goals That Were Achieved

In the project, we successfully achieved key goals that involved working with specific APIs from Reddit, Twitter, and YouTube. Utilizing Reddit's PRAW (Python Reddit API Wrapper), Twitter's Tweepy and Twitter API v2, and the Google API Python Client for YouTube, we were able to systematically gather a diverse array of data from these platforms. This data encompassed not only the textual content of posts and tweets but also included essential popularity metrics like upvotes, likes, and timestamps. Such information enabled us to analyze trends and sentiments associated with various topics as they unfolded on these social media platforms. Moreover, through the integration of sentiment analysis techniques, we could interpret and quantify the general public sentiment expressed in these posts and comments. The project effectively gathered and processed data from these platforms, providing a comprehensive view of public opinion and trends across different social media landscapes.

# The Problems We Faced

A large problem we faced was with API access. Since our project focuses on social media and there's a lot of controversy about privacy in this field right now, it was difficult to get API access for some platforms. Additionally, we had to deal with rate limiting for many APIs due to being on the lowest tier.

Another issue was authorization. We couldn't find enough API's that gave us the data information we needed so we had to pay for authorization in bigger APIs like twitter to access enough information to get the 3 APIs working in the same way.

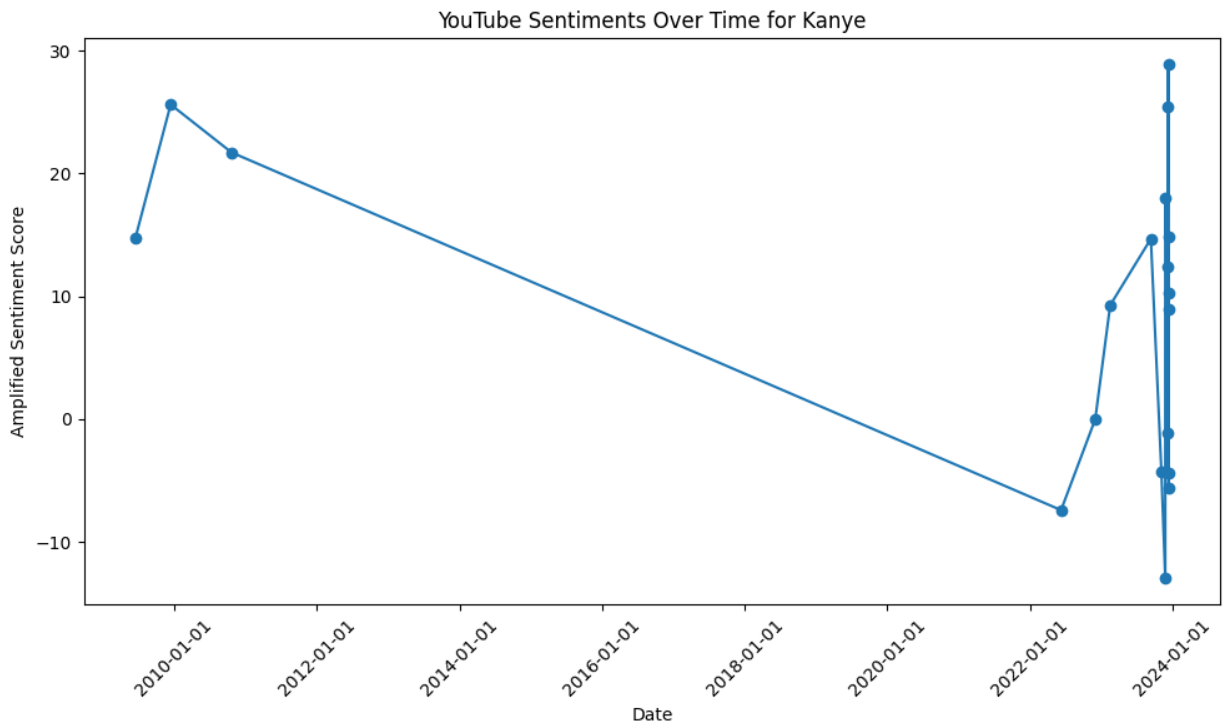
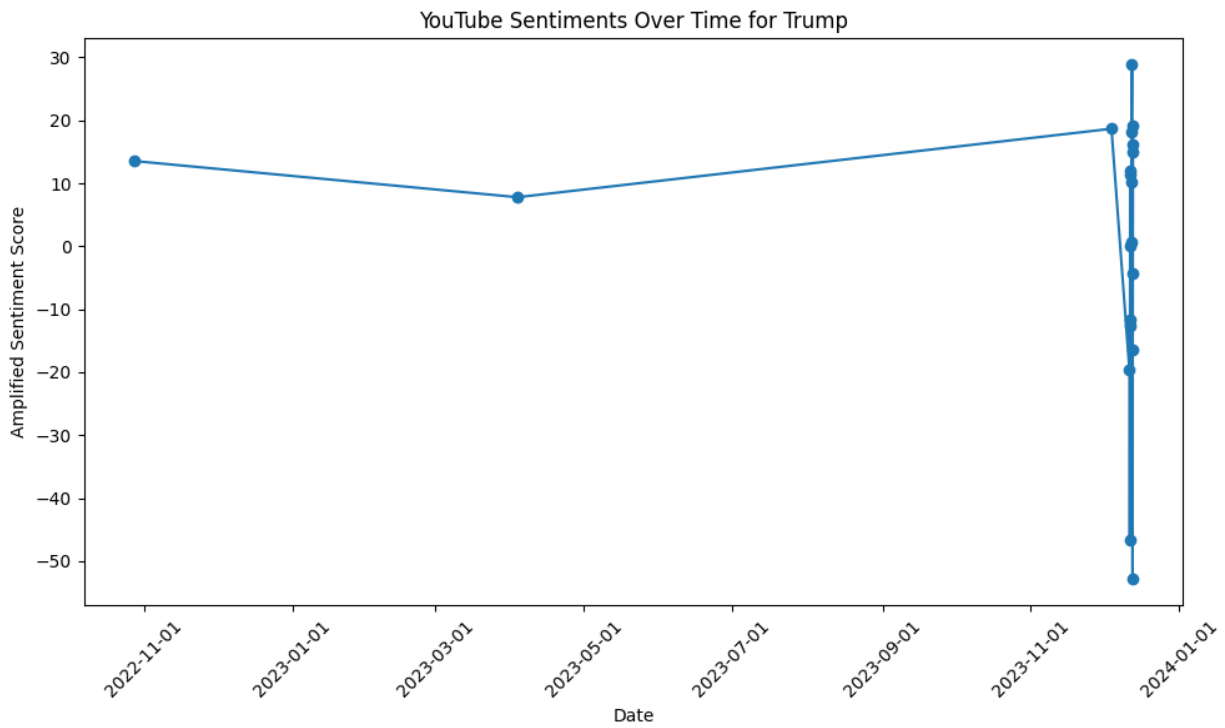
A lot of the data we collected had locked or removed comments sections. We realized this was adding "zero" data that was causing errors in our calculations so we adjusted this by skipping over these data points.

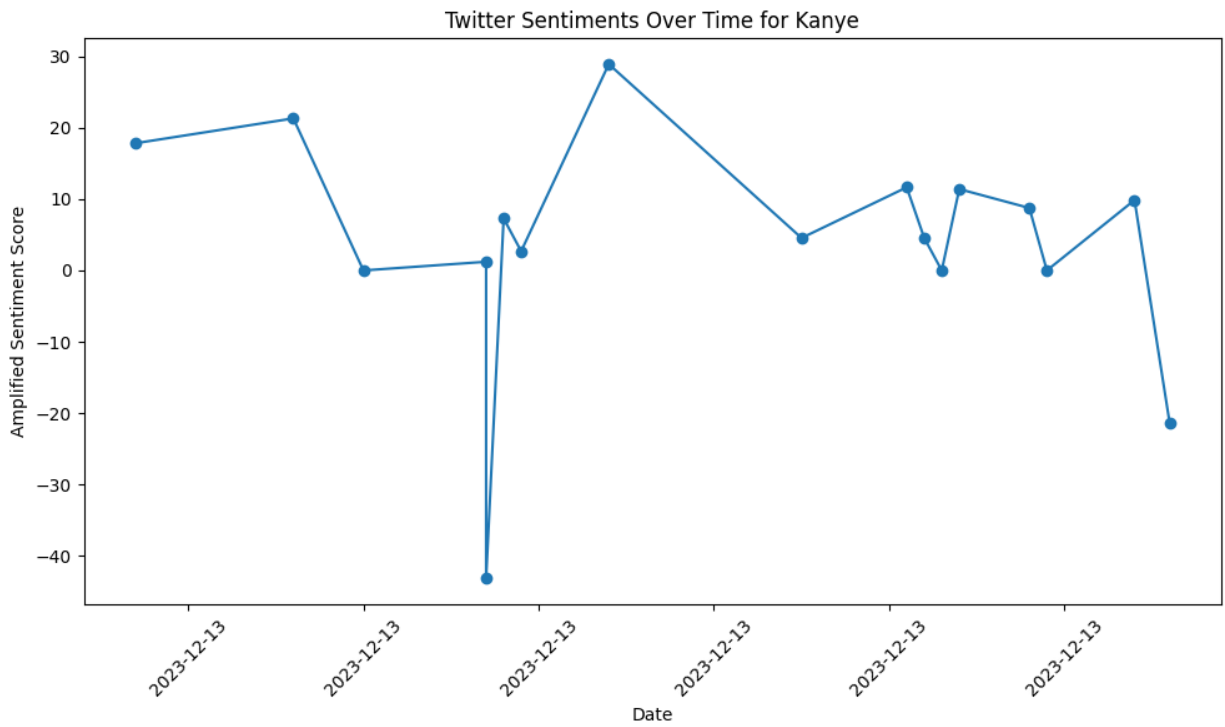
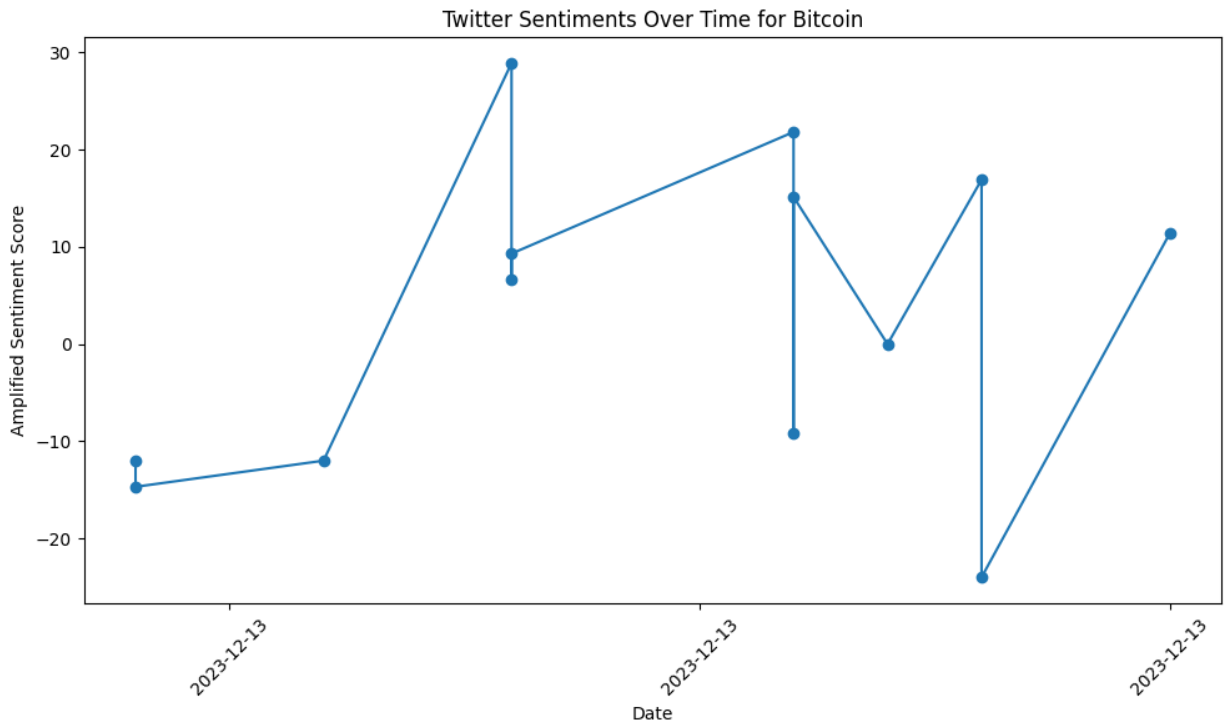
Some of the posts also weren't available in some countries or had safety warnings on twitter so they were locked to a viewer or viewers which we tried to account for in the sentiment analysis.

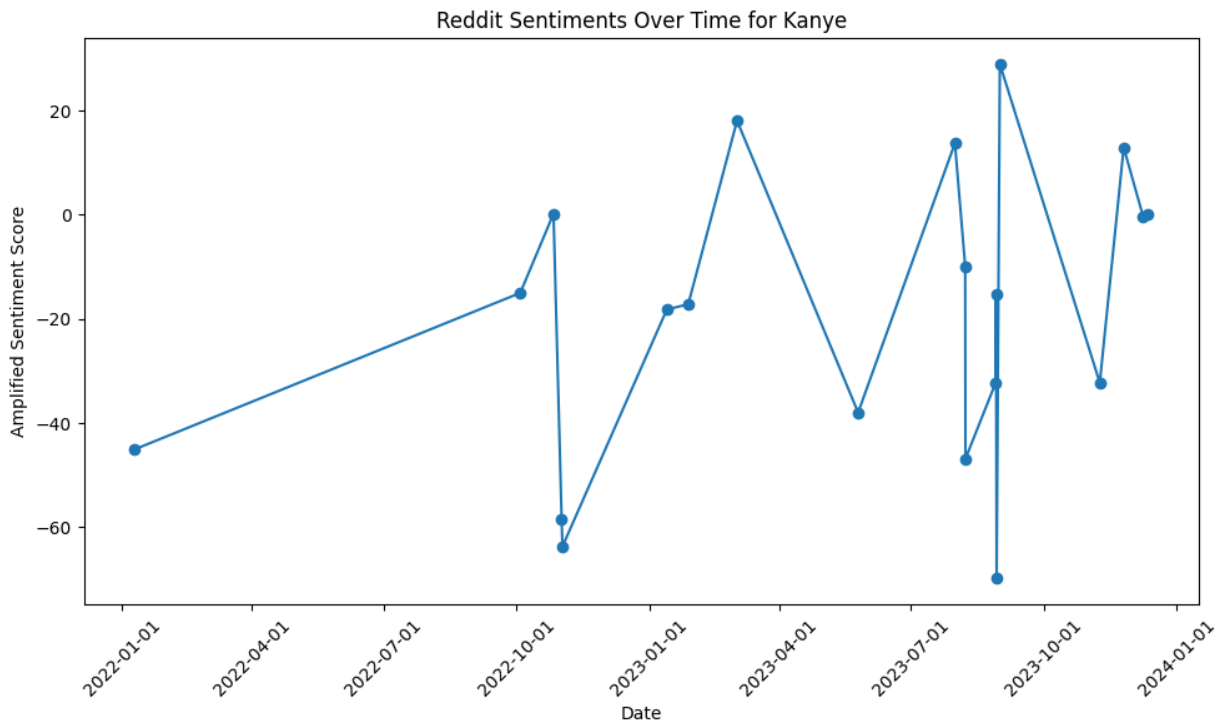
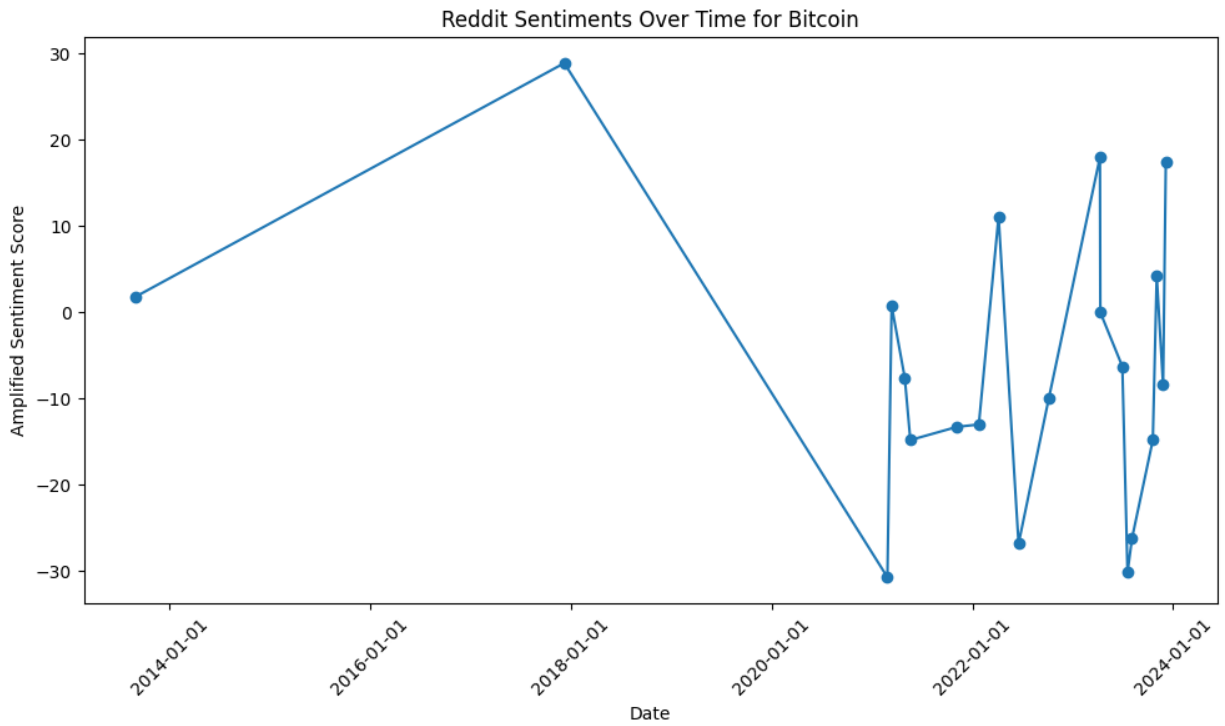
# Calculations from Database

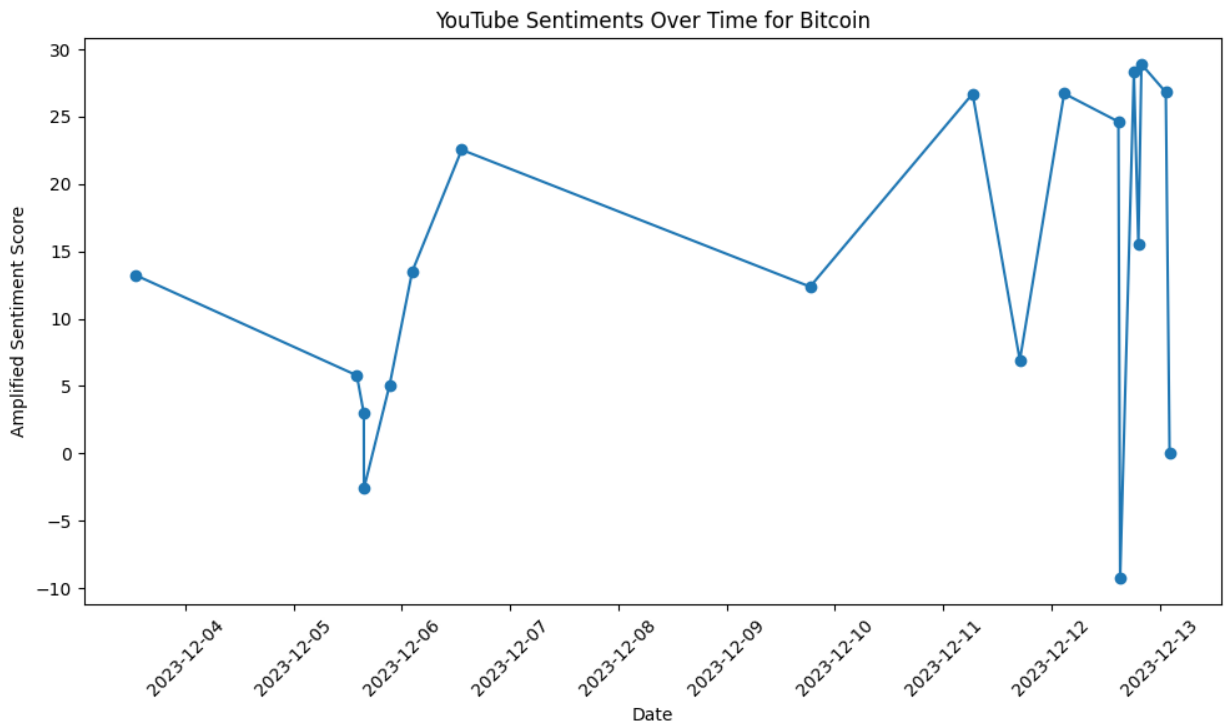
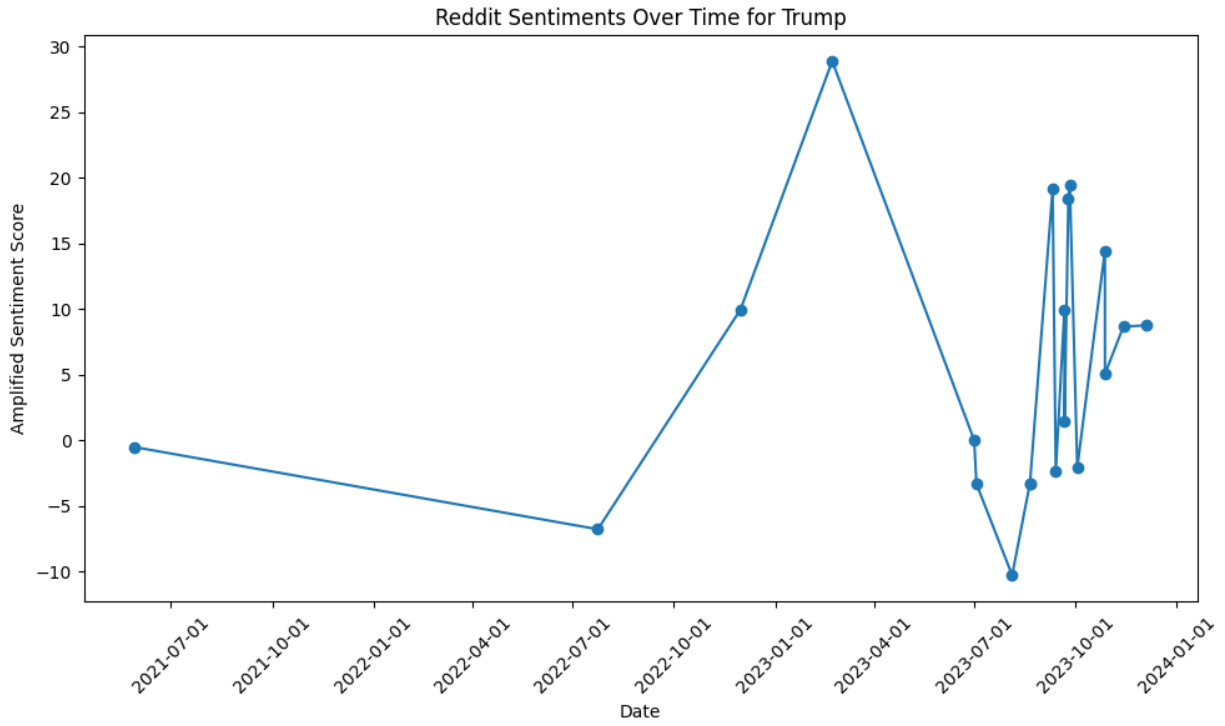
```
-16.34987844169757
14.980843758388818
10.20168756787266
0
sentiments for Trump using Twitter API
-48.8527533411597
-6.113839743334457
-35.072199219481435
18.684882145268266
-2.9807807176138893
10.684882145268266
-28.783826189688858
17.7638126585892
8.283456442719919
-101.85378248479373
3.477231776638836
-9.37946794258585367
1.7748737778184376
-1.8517824847937378
1.4425752935893983
14.632839181524484
28.986482631788784
-36.838891485661134
0
18.684882145268266
-2.9807807176138893
27.852780146995822
-19.84313512398684
0
-89.8565882283488
0
sentiments for Kanye using Reddit API
-45.86658617762159
-15.81524738864558
0.18492697491767888
-38.57289568384528
-63.725512438577425
-18.284729555427783
-17.104137731868864
18.948780739747924
-37.88888784958428
13.724761551979528
-9.888839132784664
-46.86948813175385
-32.29329994158531
-15.187956849114837
-69.79294847588925
28.986482631788784
-32.288484553194736
12.895138746947982
-8.4816824178339247
0
sentiments for Kanye using YouTube API
14.783362878246549
25.617568595258557
21.718757431882958
-7.3988841869882865
0
9.282577951413492
14.689638839482583
-4.381213886617385
-12.948936267355976
17.9669326724821
12.414799388882836
-1.071551872868585
25.407538625531324
28.986482631788784
-4.248893226535214
8.989198615497468
10.294158153845557
-5.53845153807155
14.87788616888483
0
sentiments for Kanye using Twitter API
17.82579997832596
21.36388314799134
0
1.2176993388888882
-43.14525618214888
7.286263857853381
2.691465388898487
```

# Visualizations









## Instructions for Running The Code

The code is all done in the background. All you have to do is run the mainFile.py file. The code will prompt you with asking for inputs for topics you want to test. after Giving 3 topics, the code will run and after querying the APIs and retrieving from the database, will give you the results in text format in output.txt

**Step 1.** Clone repository: \$ git clone

<https://github.com/sammalik111/sentimentAnalysis.git>

**Step 2.** Install packages: \$ pip3 install praw tweepy google-api-python-client textblob

**Step 3.** Run program: \$ python3 mainFile.py

**Step 4.** Give 3 topics when prompted

**Step 5.** Wait for information gathering

**Step 6.** Observe Data

## Documentation for Functions

### grabData.py

- **fetch\_reddit\_data(topic):** Fetches and analyzes data from Reddit. Input is a topic string, and output is a list of dictionaries containing data for each post.
- **fetch\_youtube\_data(topic):** Retrieves and analyzes data from YouTube. Input is a topic string, and output is a list of dictionaries with video data.
- **calculate\_comments\_sentiment(tweet\_data):** Calculates the sentiment score of a tweet. Input is a dictionary containing tweet data, and output is a sentiment score.
- **fetch\_twitter\_data(topic):** Fetches and analyzes data from Twitter. Input is a topic string, and output is a list of tweet data dictionaries.

### mainFile.py

- **calculate\_sentiments\_and\_times(data):** Extracts sentiment scores and times from data. Input is a list of data items, and output is two lists (sentiments and times).
- **main():** Coordinates data gathering. Then pushes the information to the database, and grabs everything from the database it needs to output. calculates all the



information needed to be visualized then calls visualize to display everything.

### **visualize.py**

- **amplify\_variation(scores, base):** Normalizes and amplifies sentiment scores onto a 0-100 score. Input is a list of scores; output is an amplified list of scores.
- **TimePlot(dates, scores, title):** Creates a time-series plot. Inputs are lists of dates and scores, and a title string.
- **histPlot(scores, color, title):** Generates a histogram. Inputs are a list of scores, color for the plot, and a title string. It then outputs the figure and scores to their respective files.
- **barPlot(topics, avg\_sentiments, title):** Creates a bar plot. Inputs are lists of topics and average sentiments, and a title string. It then outputs the figure and scores to their respective files.

### **DatabaseMaker.py**

- **fetch\_data\_from\_database(topic, api\_name):** This function fetches data from the database for a specific topic and API. It takes two parameters, topic (a string) and api\_name (a string), and returns a list of dictionaries containing the fetched data.
- **create\_and\_insert\_table(data, topic, api\_name):** This function creates and inserts data into a table in the database. It takes three parameters, data, topic, and api\_name . It creates a table for the topic and API if it doesn't exist and inserts data into the table while ignoring duplicates.

## **Resource Documentation**

### **Sentiment Analysis Accuracy**

- Date: 12/04/2023
- Issue Description: Challenges in ensuring accuracy of sentiment analysis due to the diverse and ambiguous nature of social media content.
- Location of Resource: Online tutorials, and TextBlob documentation.
- Result: Partially solved. The use of the TextBlob library provided a baseline for

sentiment analysis, but the inherent ambiguities of natural language in social media meant that accuracy was not always optimal.

### **Handling Inconsistent Time Metrics**

- Date: 12/06/2023
- Issue Description: Format of dates was inconsistent across all platforms
- Location of Resource: Data cleaning and preprocessing tutorials, Stack Overflow for Python
- Result: Solved with a lambda function, making sure all data points in visualizations were in order. Implementing data cleaning and preprocessing steps helped in normalizing and sanitizing the data, but some level of data inconsistency remained unavoidable.

### **Data Consistency Across Platforms**

- Date: 12/07/2023
- Issue Description: Difficulty in standardizing and comparing data across different platforms due to varied data structures and metrics.
- Location of Resource: Platform-specific API documentation, data normalization techniques from data science forums and publications.
- Result: Solved to a degree. While each platform's API presented unique data structures, the development of a custom data processing pipeline allowed for the extraction and normalization of key data points (like sentiment scores, timestamps, and popularity metrics) into a consistent format for analysis.

### **Finding a 3rd API That Allowed for Easy Post Retrieval**

- Date: 12/09/2023
- Issue Description: Difficulty in finding free, legal, first-party APIs for post retrieval, particularly for Twitter.
- Location of Resource: Extensive online research, including developer forums and official API documentation.
- Result: Solved. We paid ~\$28 for Basic access to Twitter API v2 (paid prorated amount due to already being partially through the month)