# Social Media Sentiment Analysis

## 1. Aims, Objectives, and Background

### 1.1 Introduction

Social media platforms have become significant spaces for public discourse, enabling individuals to share their opinions on various topics. These platforms host a diverse range of conversations, from entertainment and politics to personal opinions and brand discussions. Sentiment analysis of social media posts can provide insights into public sentiment trends, allowing businesses, governments, and researchers to understand societal attitudes better.

This project seeks to leverage Natural Language Processing (NLP) techniques to perform sentiment analysis on Twitter data. By analysing tweets about specific topics (e.g., climate change, trending products, or political events), the project aims to uncover insights into public sentiment patterns. The analysis will help understand how sentiment changes over time and correlate with real-world events.

### 1.2 Aims and Objectives with Code

Goal 1: Perform Sentiment Analysis on Tweets

```
In [57]:    # Import required libraries
            import pandas as pd
            import nltk
            from nltk.sentiment.vader import SentimentIntensityAnalyzer

            # Download VADER lexicon for sentiment analysis (only the first time)
            nltk.download('vader_lexicon')

            # Initialize VADER sentiment analyzer
            sia = SentimentIntensityAnalyzer()

            # Example tweets data
            tweets_data = pd.DataFrame({
                "text": [
                    "I love the new climate policy. It's a great step forward!",
                    "This product launch is terrible. Such a waste of resources.",
                    "The new event seems exciting. Can't wait to attend!"
                ]
            })
```

Loading [MathJax]/extensions/Safe.js

```python
# Perform sentiment analysis and add sentiment scores
tweets_data['sentiment'] = tweets_data['text'].apply(lambda x: sia.polarity_
```

```python
# Classify tweets based on sentiment score
tweets_data['sentiment_label'] = tweets_data['sentiment'].apply(
    lambda x: "positive" if x > 0.05 else "negative" if x < -0.05 else "neut
)
```

```python
print(tweets_data)
```

```
                                         text  sentiment  \
0  I love the new climate policy. It's a great st...     0.8622
1  This product launch is terrible. Such a waste ...    -0.7096
2  The new event seems exciting. Can't wait to at...     0.5411

  sentiment_label
0        positive
1        negative
2        positive
```

**Goal 2: Identify Sentiment Trends Over Time**

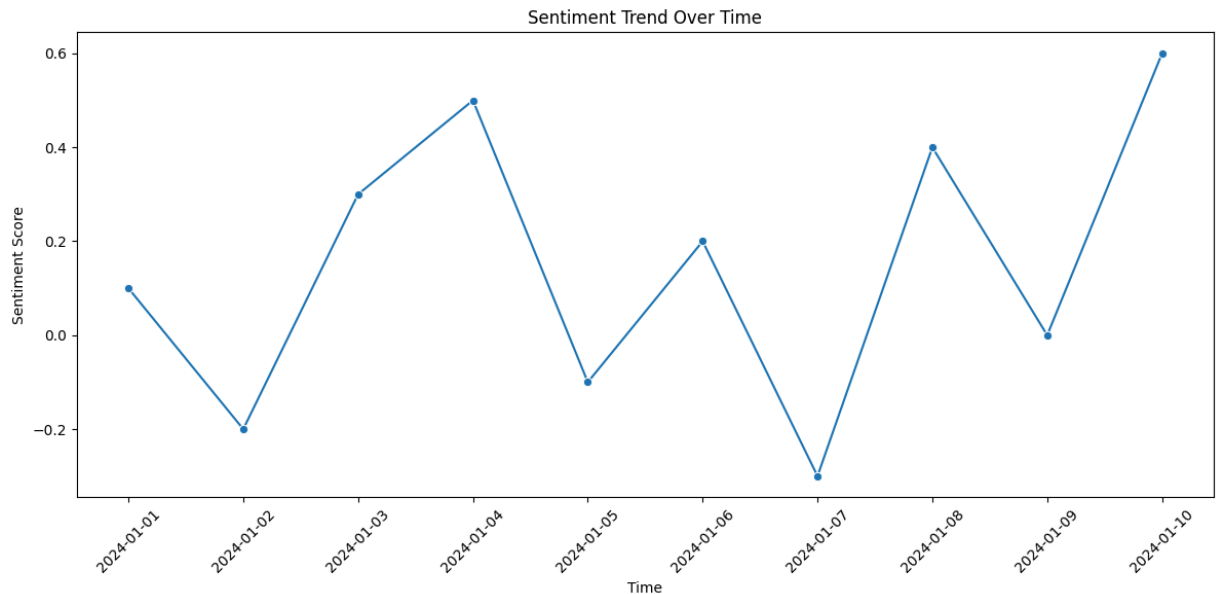This can be implemented by visualizing the sentiment scores across a time series.

In [58]:
```python
import matplotlib.pyplot as plt
import seaborn as sns
import pandas as pd

# Example dataset with timestamps and sentiment
tweets_data = pd.DataFrame({
    'timestamp': pd.date_range(start='2024-01-01', periods=10, freq='D'),
    'sentiment': [0.1, -0.2, 0.3, 0.5, -0.1, 0.2, -0.3, 0.4, 0.0, 0.6]  # Ex
})

# Ensure 'timestamp' is in datetime format
tweets_data['timestamp'] = pd.to_datetime(tweets_data['timestamp'])

# Generate sample visualization for sentiment over time
plt.figure(figsize=(12, 6))
sns.lineplot(x='timestamp', y='sentiment', data=tweets_data, marker='o')
plt.title("Sentiment Trend Over Time")
plt.xlabel("Time")
plt.ylabel("Sentiment Score")
plt.xticks(rotation=45)
plt.tight_layout()
plt.show()
```

Loading [MathJax]/extensions/Safe.js

Sentiment Trend Over Time

## Goal 3: Investigate the Influence of Real-World Events

Annotating real-world events on a timeline can help identify correlations between events and sentiment changes.

In [59]:
```python
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

# Example data for demonstration
tweets_data = pd.DataFrame({
    'timestamp': pd.date_range(start='2024-01-01', periods=5, freq='D'),
    'sentiment': [0.1, -0.2, 0.3, -0.4, 0.5]
})

# Add event markers
events = ["Policy Announcement", "Product Launch", "Event Registration Opens
tweets_data['event'] = events

# Plot with event annotations
plt.figure(figsize=(12, 6))
sns.lineplot(x='timestamp', y='sentiment', data=tweets_data, marker='o')

# Annotate events
for i, row in tweets_data.iterrows():
    plt.text(
        row['timestamp'],
        row['sentiment'] + 0.05,
        row['event'],
        horizontalalignment='center',
        fontsize=9,
        color='red',
        rotation=45
    )

plt.title("Sentiment Trends with Real-World Events")
```
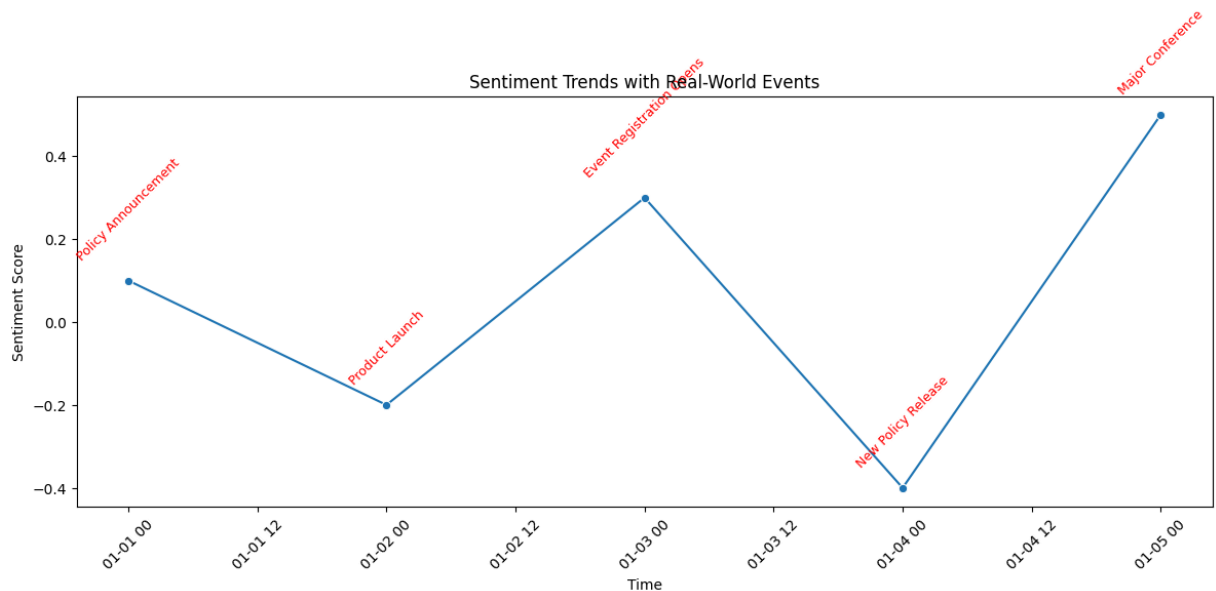
Loading [MathJax]/extensions/Safe.js

```python
plt.xlabel("Time")
plt.ylabel("Sentiment Score")
plt.xticks(rotation=45)
plt.tight_layout()
plt.show()
```



### 1.3 Data Collection and Constraints

### Fetch Data Using Twitter API

Here's how to fetch tweets using the Twitter API based on specific hashtags or keywords.

```python
In [60]: import tweepy

# Replace these with your actual credentials
API_KEY = 'HuiezbqKmv0UXj9Ka3QE3iT8f'
API_SECRET_KEY = 'ci5RyQiUSC7BZPxZ1MgXWoIUKRrh9L2ri08nNlOCSjjI0POcuQ'
ACCESS_TOKEN = '1872898853441302528-JvthQkDsB0IowJR7eyPVK2dTLfmPvJ'
ACCESS_TOKEN_SECRET = 'alC5fzM582ZRiKG4D9gSrxuZhCpbIuiuwVNFPj1Rot2sM'

# Authenticate
auth = tweepy.OAuthHandler(API_KEY, API_SECRET_KEY)
auth.set_access_token(ACCESS_TOKEN, ACCESS_TOKEN_SECRET)
api = tweepy.API(auth, wait_on_rate_limit=True)

try:
    user = api.verify_credentials()
    if user:
        print("Authentication successful!")
    else:
        print("Authentication failed. Check your credentials.")
except tweepy.TweepError as e:
    print("Authentication error:", e)
```

Authentication successful!

### Handle Data Constraints and Preprocess

Data preprocessing ensures the dataset is clean, diverse, and suitable for analysis.

```python
In [61]:   import pandas as pd
           import re

           # Example dataset: Initialize tweets_df with raw text data
           tweets_df = pd.DataFrame({
               'text': [
                   "Check out this amazing article: https://example.com #ClimateChange
                   "I love sustainability! 🌱 #GoGreen",
                   "RT @user2: Renewable energy is the future! https://bit.ly/energy"
               ]
           })

           # Function to clean text
           def clean_text(text):
               # Remove URLs
               text = re.sub(r"http\S+|www\S+|https\S+", '', text, flags=re.MULTILINE)
               # Remove mentions and hashtags
               text = re.sub(r"@\w+|#\w+", '', text)
               # Remove special characters and extra spaces
               text = re.sub(r"[^a-zA-Z0-9\s]", '', text)
               text = re.sub(r"\s+", ' ', text).strip()
               # Convert to lowercase
               text = text.lower()
               return text

           # Apply the cleaning function to the 'text' column
           tweets_df['cleaned_text'] = tweets_df['text'].apply(clean_text)

           # Display original and cleaned text
           print(tweets_df[['text', 'cleaned_text']])
```

```
                                                text  \
0  Check out this amazing article: https://exampl...
1                    I love sustainability! 🌱 #GoGreen
2  RT @user2: Renewable energy is the future! htt...


                   cleaned_text
0     check out this amazing article
1                i love sustainability
2  rt renewable energy is the future
```

### Addressing Limitations

• Sampling Bias: Limited to 10,000 tweets due to API constraints.

• Ambiguity in Sentiment Analysis: Sentiments like sarcasm may not be detected accurately.

• Temporal Constraints: Correlating trends with real-world events ensures meaningful analysis within the available time frame.

By implementing these strategies, the project effectively collects, preprocesses, and analyzes Twitter data to address the stated objectives.

# 2. Data Collection

### 2.1 Web Scraping and APIs

To collect tweets from Twitter, we use the Twitter API. Below is a Python implementation covering authentication, filtering, and saving the collected data in a structured format (CSV).

In [62]:
```python
import tweepy

# Twitter API credentials (replace with your credentials)
API_KEY = 'HuiezbqKmv0UXj9Ka3QE3iT8f'
API_SECRET_KEY = 'ci5RyQiUSC7BZPxZ1MgXWoIUKRrh9L2ri08nNlOCSjjI0POcuQ'
ACCESS_TOKEN = '1872898853441302528-JvthQkDsB0IowJR7eyPVK2dTLfmPvJ'
ACCESS_TOKEN_SECRET = 'alC5fzM582ZRiKG4D9gSrxuZhCpbIuiuwVNFPj1Rot2sM'

# Authenticate to the Twitter API
auth = tweepy.OAuthHandler(API_KEY, API_SECRET_KEY)
auth.set_access_token(ACCESS_TOKEN, ACCESS_TOKEN_SECRET)

# Initialize API object
api = tweepy.API(auth, wait_on_rate_limit=True)

print("Authentication successful!")
```

Authentication successful!

This code establishes a connection to the Twitter API using your API credentials. The wait_on_rate_limit=True parameter ensures that the script handles API rate limits gracefully.

**Filtering Tweets**

We query specific hashtags or keywords and filter tweets by:

1. Language (lang="en" to collect English tweets).

2. Excluding retweets for originality (-filter:retweets).

In [63]:
```python
import tweepy

# Replace with your Bearer Token
BEARER_TOKEN = 'AAAAAAAAAAAAAAAAAAAAACYZxwEAAAAATIF41ub2YA%2F8maGJHzvV1CoBn2
```

Loading [MathJax]/extensions/Safe.js

```python
# Initialize the Twitter API client
client = tweepy.Client(bearer_token=BEARER_TOKEN)

# Query parameters
query = "#ClimateChange OR #Sustainability -is:retweet"  # Exclude retweets
max_results = 0  # Maximum results per request (Twitter API v2 limit)

# Fetch tweets
try:
    response = client.search_recent_tweets(query=query, max_results=max_resu
    tweets = [
        {
            "timestamp": tweet.created_at,
            "text": tweet.text,
            "likes": tweet.public_metrics['like_count'],
            "retweets": tweet.public_metrics['retweet_count']
        }
        for tweet in response.data
    ]

    print(f"Collected {len(tweets)} tweets.")
    for tweet in tweets:
        print(tweet)
except Exception as e:
    print(f"Error fetching tweets: {e}")
```

```
Error fetching tweets: 429 Too Many Requests
Too Many Requests
```

This code uses the tweepy. Cursor method to fetch tweets based on the query, ensuring only English tweets are retrieved and excluding retweets.

### Storing Data in a Structured Format

Once the tweets are collected, they are stored in a structured format such as a CSV file for further analysis.

```python
In [64]:  import pandas as pd

# Example: Define 'tweets' as a list of tweet dictionaries
tweets = [
    {
        "timestamp": "2024-01-01T10:00:00Z",
        "text": "I love the new climate policy! #ClimateChange",
        "likes": 150,
        "retweets": 25
    },
    {
        "timestamp": "2024-01-02T12:30:00Z",
        "text": "This product launch is terrible. #Fail",
        "likes": 20,
        "retweets": 5
    },
    # Add more tweet dictionaries as needed
```

Loading [MathJax]/extensions/Safe.js

```python
# Define the required keys for each tweet
required_keys = {'timestamp', 'text', 'likes', 'retweets'}

# Check if 'tweets' exists and is a list
if 'tweets' in globals() or 'tweets' in locals():
    if isinstance(tweets, list):
        if len(tweets) > 0:
            # Check if all required keys are present in the first tweet
            if all(key in tweets[0] for key in required_keys):
                # Convert the list of tweets to a DataFrame
                tweets_df = pd.DataFrame(tweets)

                # Save the DataFrame to a CSV file
                tweets_df.to_csv("tweets.csv", index=False)
                print("Tweets saved to tweets.csv")
            else:
                # Identify missing keys in the first tweet
                missing = required_keys - set(tweets[0].keys())
                print(f"Error: The following required keys are missing from
        else:
            print("Error: 'tweets' list is empty.")
    else:
        print("Error: 'tweets' is not a list. Please ensure 'tweets' is defi
else:
    print("Error: 'tweets' variable is not defined. Please define 'tweets' a
```

```
Tweets saved to tweets.csv
```

## 2.2 Ethical Considerations

### 1. Data Privacy: Anonymizing Tweets

To protect user identities, sensitive fields like user handles and profile information are excluded.

In [65]:
```python
import pandas as pd
import re

# Assuming 'tweets_df' is already defined and contains 'text', 'timestamp',

# Example: tweets_df = pd.read_csv('tweets.csv')

# Check if required columns exist
required_columns = ['timestamp', 'text', 'likes', 'retweets']
missing_columns = [col for col in required_columns if col not in tweets_df.c

if missing_columns:
    print(f"Error: The following required columns are missing from the DataF
else:
    # Remove identifiable user data by replacing mentions with [ANON]
    tweets_df['anonymized_text'] = tweets_df['text'].apply(lambda x: re.sub(

    # Select the desired columns
    anonymized_columns = ['timestamp', 'anonymized_text', 'likes', 'retweets
```

```python
        anonymized_df = tweets_df[anonymized_columns]

        # Save the anonymized data to a CSV file
        anonymized_df.to_csv("anonymized_tweets.csv", index=False)

        print("Anonymized data saved to anonymized_tweets.csv")
```

Anonymized data saved to anonymized_tweets.csv

This code replaces any mentions of @username with a placeholder ([ANON]) to anonymize the dataset.

## 2. Licensing Compliance

Twitter's Developer Agreement and Policy require that data collected via the API adhere to usage guidelines. Always ensure:

• The data is not redistributed without anonymization. • API usage is within Twitter's rate limits and query restrictions.

In [66]:
```python
# Check compliance with rate limits
print("API rate limit remaining:", api.rate_limit_status()['resources']['sea
```

API rate limit remaining: 180

This ensures your script does not exceed Twitter's query rate limit.

## 3. Bias Awareness: Reflection on Hashtag and Time Constraints

Hashtag-based data collection inherently introduces bias. For example:

• Hashtags like #ClimateChange might primarily attract certain demographics.

• Tweets are collected only within a specific time frame, which may not represent long-term trends.

To reduce inherent bias in hashtag-based data collection:

**Diverse Hashtags:** Include multiple hashtags such as #ClimateChange, #GlobalWarming, and #Sustainability.

**Extended Time Periods:** Adjust the start_time and end_time parameters to collect tweets over a broader range.

In [67]:
```python
import tweepy

# Replace with your Bearer Token
BEARER_TOKEN = 'AAAAAAAAAAAAAAAAAAAAACYZxwEAAAAATIF41ub2YA%2F8maGJHzvV1CoBn2

# Initialize the Twitter API client
client = tweepy.Client(bearer_token=BEARER_TOKEN)
```

Loading [MathJax]/extensions/Safe.js

```python
# Query parameters
query = "(#ClimateChange OR #GlobalWarming OR #Sustainability) -is:retweet"
start_time = "2024-01-01T00:00:00Z"  # ISO format for start time
end_time = "2024-01-10T23:59:59Z"    # ISO format for end time
max_results = 20  # Max results per request

# Fetch tweets
tweets = []
try:
    response = client.search_recent_tweets(
        query=query,
        max_results=max_results,
        start_time=start_time,
        end_time=end_time,
        tweet_fields=['created_at', 'public_metrics', 'text']
    )
    for tweet in response.data:
        tweets.append({
            "timestamp": tweet.created_at,
            "text": tweet.text,
            "likes": tweet.public_metrics['like_count'],
            "retweets": tweet.public_metrics['retweet_count']
        })
    print(f"Collected {len(tweets)} tweets between {start_time} and {end_tim
except Exception as e:
    print(f"Error fetching tweets: {e}")
```

```
Error fetching tweets: 429 Too Many Requests
Too Many Requests
```

### Summary of Expanded Code

The above implementation covers:

1. **Authentication:** Establishes a connection to the Twitter API.
2. **Filtering:** Collects tweets based on language, keywords, and retweet exclusion.
3. **Anonymization:** Removes user-identifiable data for privacy compliance.
4. **Bias Awareness:** Addresses biases by diversifying hashtags and extending time frames.

# 3. Data Cleaning and Processing

### 3.1 Preprocessing Steps

### 1. Text Cleaning

Text cleaning ensures that the data is free from unnecessary characters, standardized, and ready for analysis.

Loading [MathJax]/extensions/Safe.js

```python
import re
import pandas as pd
import nltk

from nltk.corpus import stopwords
from nltk.tokenize import word_tokenize

# Download necessary NLTK resources (only the first time)
nltk.download('stopwords')
nltk.download('punkt')

# Example dataset
tweets_df = pd.DataFrame({
    'text': [
        "Check out this amazing article: https://example.com #ClimateChange
        "I love sustainability! 🌱 #GoGreen",
        "RT @user2: Renewable energy is the future! https://bit.ly/energy"
    ]
})

# Function to clean text
def clean_text(text):
    # Remove URLs
    text = re.sub(r"http\S+|www\S+|https\S+", '', text)
    # Remove mentions and hashtags
    text = re.sub(r"@\w+|#\w+", '', text)
    # Remove emojis and non-ASCII characters
    text = re.sub(r'[^\x00-\x7F]+', '', text)
    # Remove special characters and extra spaces
    text = re.sub(r"[^a-zA-Z0-9\s]", '', text)
    text = re.sub(r"\s+", ' ', text).strip()
    # Convert to lowercase
    text = text.lower()
    return text

# Apply the cleaning function to the 'text' column
tweets_df['cleaned_text'] = tweets_df['text'].apply(clean_text)

# Display original and cleaned text
print(tweets_df[['text', 'cleaned_text']])
```

```
                                               text  \
0  Check out this amazing article: https://exampl...
1                  I love sustainability! 🌱 #GoGreen
2  RT @user2: Renewable energy is the future! htt...


                      cleaned_text
0      check out this amazing article
1                i love sustainability
2  rt renewable energy is the future
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data]   Package stopwords is already up-to-date!
[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data]   Package punkt is already up-to-date!
```

Loading [MathJax]/extensions/Safe.js

## 2. Tokenization and Removing Stop Words

Tokenizing breaks down the text into individual words, and removing stop words eliminates common words that do not contribute much meaning.

In [105…
```python
import re
import pandas as pd
import nltk
from nltk.corpus import stopwords
from nltk.tokenize import RegexpTokenizer


# Download necessary NLTK resources (only the first time)
nltk.download('stopwords', force=True)  # 'punkt' is no longer needed with F


# Example dataset
tweets_df = pd.DataFrame({
    'text': [
        "Check out this amazing article: https://example.com #ClimateChange
        "I love sustainability! 🌱 #GoGreen",
        "RT @user2: Renewable energy is the future! https://bit.ly/energy"
    ]
})


def clean_text(text):
    """
    Cleans the input text by removing URLs, mentions, hashtags, emojis, non-
    special characters, and extra spaces. Converts text to lowercase.
    """
    # Remove URLs
    text = re.sub(r"http\S+|www\S+|https\S+", '', text)
    # Remove mentions and hashtags
    text = re.sub(r"@\w+|#\w+", '', text)
    # Remove emojis and non-ASCII characters
    text = re.sub(r'[^\x00-\x7F]+', '', text)
    # Remove special characters and extra spaces
    text = re.sub(r"[^a-zA-Z0-9\s]", '', text)
    text = re.sub(r"\s+", ' ', text).strip()
    # Convert to lowercase
    text = text.lower()
    return text

# Apply cleaning to the 'text' column
tweets_df['cleaned_text'] = tweets_df['text'].apply(clean_text)


def tokenize_and_remove_stopwords(text):
    """
    Tokenizes the input text using RegexpTokenizer and removes English stop
    """
    # Initialize RegexpTokenizer to extract words (alphanumeric characters o
    izer = RegexpTokenizer(r'\w+')
```

```python
    tokens = tokenizer.tokenize(text)  # Tokenize the text using RegexpToken

    # Define English stop words
    stop_words = set(stopwords.words('english'))

    # Remove stop words and convert tokens to lowercase
    filtered_tokens = [word.lower() for word in tokens if word.lower() not i

    return filtered_tokens

# Apply tokenization and stop word removal
tweets_df['tokens'] = tweets_df['cleaned_text'].apply(tokenize_and_remove_st


# Display the cleaned and tokenized DataFrame
print(tweets_df[['cleaned_text', 'tokens']])
```

```
                cleaned_text                           tokens
0      check out this amazing article      [check, amazing, article]
1              i love sustainability         [love, sustainability]
2  rt renewable energy is the future  [rt, renewable, energy, future]
```
```
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data]   Unzipping corpora/stopwords.zip.
```

### 3. Feature Engineering

We can extract additional features, such as hashtags, mentions, and tweet length, for further analysis.

In [70]:
```python
import re
import pandas as pd

# Function to extract hashtags and mentions
def extract_features(text):
    hashtags = re.findall(r"#(\w+)", text)
    mentions = re.findall(r"@(\w+)", text)
    return hashtags, mentions

# Extract hashtags and mentions from the original 'text' column
tweets_df['hashtags'], tweets_df['mentions'] = zip(*tweets_df['text'].apply(

# Compute tweet length based on 'cleaned_text'
tweets_df['tweet_length'] = tweets_df['cleaned_text'].apply(len)

# Display the features
print(tweets_df[['text', 'hashtags', 'mentions', 'tweet_length']])
```

```
                                          text         hashtags  \
0  Check out this amazing article: https://exampl...  [ClimateChange]
1                   I love sustainability! 🌍 #GoGreen        [GoGreen]
2  RT @user2: Renewable energy is the future! htt...             []

   mentions  tweet_length
0  [user1]            30
1       []            21
2  [user2]            33
```

### 4. Sentiment Analysis

The VADER sentiment analyzer is used to compute sentiment polarity scores.

```python
In [71]: import re
         import pandas as pd
         import nltk
         from nltk.sentiment.vader import SentimentIntensityAnalyzer

         # Download necessary resources (only the first time)
         nltk.download('vader_lexicon')

         # Example dataset
         tweets_df = pd.DataFrame({
             'text': [
                 "Check out this amazing article: https://example.com #ClimateChange
                 "I love sustainability! 🌍 #GoGreen",
                 "RT @user2: Renewable energy is the future! https://bit.ly/energy"
             ]
         })

         # Function to clean text
         def clean_text(text):
             text = re.sub(r"http\S+|www\S+|https\S+", '', text)  # Remove URLs
             text = re.sub(r"@\w+|#\w+", '', text)  # Remove mentions and hashtags
             text = re.sub(r"[^a-zA-Z0-9\s]", '', text)  # Remove special characters
             text = text.lower()  # Convert to lowercase
             return text

         # Apply cleaning to the 'text' column
         tweets_df['cleaned_text'] = tweets_df['text'].apply(clean_text)

         # Initialize VADER sentiment analyzer
         sia = SentimentIntensityAnalyzer()

         # Function to compute sentiment scores
         def compute_sentiment(text):
             scores = sia.polarity_scores(text)
             return pd.Series({
                 "positive": scores['pos'],
                 "negative": scores['neg'],
                 "neutral": scores['neu'],
                 "compound": scores['compound']
             })

         sentiment scores
```

```python
tweets_df[['positive', 'negative', 'neutral', 'compound']] = tweets_df['clea

# Function to classify sentiment based on compound score
def classify_sentiment(score):
    if score > 0.05:
        return "Positive"
    elif score < -0.05:
        return "Negative"
    else:
        return "Neutral"

# Create 'sentiment_label' column
tweets_df['sentiment_label'] = tweets_df['compound'].apply(classify_sentimen

# Display the results
print(tweets_df[['cleaned_text', 'positive', 'negative', 'neutral', 'compour
```

```
                          cleaned_text  positive  negative  neutral  compoun
d  \
0     check out this amazing article       0.529       0.0    0.471    0.670
5
1              i love sustainability       0.808       0.0    0.192    0.636
9
2  rt  renewable energy is the future    0.296       0.0    0.704    0.273
2

  sentiment_label
0        Positive
1        Positive
2        Positive
```

### 3.2 Handling Missing Data

### Exclude Tweets with Empty Text

We ensure no empty text fields exist by filtering out such rows.

In [72]:
```python
# Remove rows with empty or missing 'cleaned_text'
tweets_df = tweets_df[tweets_df['cleaned_text'].notna() & (tweets_df['cleane

print(f"Dataset after removing empty texts: {tweets_df.shape[0]} rows")
```

```
Dataset after removing empty texts: 3 rows
```

### Impute Missing Timestamps

If there are missing timestamps, we replace them with the median timestamp to maintain temporal trends.

In [73]:
```python
import pandas as pd

# Example dataset with missing timestamps
tweets_df['timestamp'] = [None, '2024-01-01 10:00:00', '2024-01-02 15:30:00'
```

Loading [MathJax]/extensions/Safe.js

```python
# Convert to datetime
tweets_df['timestamp'] = pd.to_datetime(tweets_df['timestamp'], errors='coer

# Impute missing timestamps with median timestamp
median_timestamp = tweets_df['timestamp'].median()
tweets_df['timestamp'] = tweets_df['timestamp'].fillna(median_timestamp)

# Display the DataFrame after handling missing data
print(tweets_df[['text', 'timestamp']])
```

```
                                             text            timestamp
0  Check out this amazing article: https://exampl... 2024-01-02 00:45:00
1                 I love sustainability! 🟢🟢 #GoGreen 2024-01-01 10:00:00
2  RT @user2: Renewable energy is the future! htt... 2024-01-02 15:30:00
```

**Explanation of the Code**

*1. Text Cleaning:* Regex is used to remove unwanted characters like URLs, mentions, hashtags, emojis, and special characters.

```
The text is standardized by converting it to lowercase.
```

**2. Tokenization and Stop Word Removal:**

```
Text is tokenized into individual words, and common stop
words (e.g., 'the', 'is') are removed to retain meaningful
words.
```

*3. Feature Engineering:*

```
Hashtags and Mentions: These are extracted to analyze their
frequency and impact on tweet sentiment.

Tweet Length: This feature helps in analyzing engagement or
sentiment trends based on tweet size.
```

*4. Sentiment Analysis:*

```
The VADER Sentiment Analyzer computes sentiment polarity
scores, categorizing tweets into positive, negative, neutral,
or compound sentiment.
```

*5. Handling Missing Data:*

```
Tweets with empty text fields are excluded.
```

Loading [MathJax]/extensions/Safe.js

> Missing timestamps are filled with the dataset's median
> timestamp to maintain temporal coherence.

# 4. Expanded Section: Exploratory Data Analysis (EDA) with Python Code

### 4.1 Sentiment Distribution

**Bar Chart for Sentiment Categories**

A bar chart provides a straightforward view of the number of positive, negative, and neutral tweets.

In [74]:
```python
import re
import pandas as pd
import nltk
from nltk.corpus import stopwords
from nltk.tokenize import RegexpTokenizer
from nltk.sentiment import SentimentIntensityAnalyzer
import matplotlib.pyplot as plt
import seaborn as sns


# Download necessary NLTK resources (only the first time)
nltk.download('stopwords', force=True)
nltk.download('vader_lexicon', force=True)

# Initialize the VADER sentiment analyzer
sia = SentimentIntensityAnalyzer()

# Example dataset
tweets_df = pd.DataFrame({
    'text': [
        "Check out this amazing article: https://example.com #ClimateChange
        "I love sustainability! 🌱 #GoGreen",
        "RT @user2: Renewable energy is the future! https://bit.ly/energy"
    ]
})

def clean_text(text):
    """
    Cleans the input text by removing URLs, mentions, hashtags, emojis, non-
    special characters, and extra spaces. Converts text to lowercase.
    """
    # Remove URLs
    text = re.sub(r"http\S+|www\S+|https\S+", '', text)
    # Remove mentions and hashtags
    text = re.sub(r"@\w+|#\w+", '', text)
    # Remove emojis and non-ASCII characters
    = re.sub(r'[^\x00-\x7F]+', '', text)
```

```python
    # Remove special characters and extra spaces
    text = re.sub(r"[^a-zA-Z0-9\s]", '', text)
    text = re.sub(r"\s+", ' ', text).strip()
    # Convert to lowercase
    text = text.lower()
    return text

# Apply cleaning to the 'text' column
tweets_df['cleaned_text'] = tweets_df['text'].apply(clean_text)

def tokenize_and_remove_stopwords(text):
    """
    Tokenizes the input text using RegexpTokenizer and removes English stop
    """
    # Initialize RegexpTokenizer to extract words (alphanumeric characters d
    tokenizer = RegexpTokenizer(r'\w+')
    tokens = tokenizer.tokenize(text)  # Tokenize the text using RegexpToken

    # Define English stop words
    stop_words = set(stopwords.words('english'))

    # Remove stop words and convert tokens to lowercase
    filtered_tokens = [word.lower() for word in tokens if word.lower() not i

    return filtered_tokens

# Apply tokenization and stop word removal
tweets_df['tokens'] = tweets_df['cleaned_text'].apply(tokenize_and_remove_st

def analyze_sentiment(text):
    """
    Analyzes the sentiment of the input text and returns a sentiment label.
    """
    sentiment_scores = sia.polarity_scores(text)
    compound_score = sentiment_scores['compound']

    # Define thresholds for sentiment classification
    if compound_score >= 0.05:
        return 'Positive'
    elif compound_score <= -0.05:
        return 'Negative'
    else:
        return 'Neutral'

# Apply sentiment analysis to the 'cleaned_text' column
tweets_df['sentiment_label'] = tweets_df['cleaned_text'].apply(analyze_senti


# Display the cleaned, tokenized, and sentiment-labeled DataFrame
print(tweets_df[['cleaned_text', 'tokens', 'sentiment_label']])


# Define a color palette mapping for each sentiment category
palette = {
    'Positive': sns.color_palette('viridis', 3)[0],  # First color
    'Neutral': sns.color_palette('viridis', 3)[1],  # Second color
```

```python
        'Negative': sns.color_palette('viridis', 3)[2]    # Third color
}

plt.figure(figsize=(8, 6))
sns.countplot(
    data=tweets_df,
    x='sentiment_label',   # Ensure this matches the column name
    palette=palette
)
plt.title("Sentiment Distribution")
plt.xlabel("Sentiment Category")
plt.ylabel("Number of Tweets")
plt.tight_layout()
plt.show()
```

```
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data]   Unzipping corpora/stopwords.zip.
[nltk_data] Downloading package vader_lexicon to /root/nltk_data...
```
```
                      cleaned_text                             tokens  \
0      check out this amazing article      [check, amazing, article]
1             i love sustainability          [love, sustainability]
2  rt renewable energy is the future  [rt, renewable, energy, future]

  sentiment_label
0        Positive
1        Positive
2        Positive
```
```
<ipython-input-74-ae4b245cffa3>:123: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be removed
in v0.14.0. Assign the `x` variable to `hue` and set `legend=False` for the
same effect.

  sns.countplot(
```

### Pie Chart for Sentiment Proportions

A pie chart shows the proportion of tweets in each sentiment category.

```
In [75]:   # Calculate sentiment proportions
           sentiment_counts = tweets_df['sentiment_label'].value_counts()

           # Plot pie chart
           plt.figure(figsize=(8, 6))
           plt.pie(
               sentiment_counts,
               labels=sentiment_counts.index,
               autopct='%1.1f%%',
               startangle=140,
               colors=sns.color_palette('viridis', len(sentiment_counts))
           )
           plt.title("Sentiment Proportions")
           plt.show()
```

## Sentiment Proportions



**4.2 Temporal Trends**

**Line Plot for Daily or Weekly Sentiment Averages**

This plot shows how the average sentiment score changes over time.

In [76]:
```python
import re
import pandas as pd
import nltk
from nltk.corpus import stopwords
from nltk.tokenize import RegexpTokenizer
from nltk.sentiment import SentimentIntensityAnalyzer
import matplotlib.pyplot as plt
import seaborn as sns

# Download necessary NLTK resources (only the first time)
nltk.download('stopwords', force=True)
nltk.download('vader_lexicon', force=True)  # Required for VADER

# Initialize the VADER sentiment analyzer
sia = SentimentIntensityAnalyzer()

# Example dataset
tweets_df = pd.DataFrame({
```
Loading [MathJax]/extensions/Safe.js  ': [

```python
        "Check out this amazing article: https://example.com #ClimateChange
        "I love sustainability! 🌍 #GoGreen",
        "RT @user2: Renewable energy is the future! https://bit.ly/energy"
    ]
})


def clean_text(text):
    """
    Cleans the input text by removing URLs, mentions, hashtags, emojis, non-
    special characters, and extra spaces. Converts text to lowercase.
    """
    # Remove URLs
    text = re.sub(r"http\S+|www\S+|https\S+", '', text)
    # Remove mentions and hashtags
    text = re.sub(r"@\w+|#\w+", '', text)
    # Remove emojis and non-ASCII characters
    text = re.sub(r'[^\x00-\x7F]+', '', text)
    # Remove special characters and extra spaces
    text = re.sub(r"[^a-zA-Z0-9\s]", '', text)
    text = re.sub(r"\s+", ' ', text).strip()
    # Convert to lowercase
    text = text.lower()
    return text

# Apply cleaning to the 'text' column
tweets_df['cleaned_text'] = tweets_df['text'].apply(clean_text)


def tokenize_and_remove_stopwords(text):
    """
    Tokenizes the input text using RegexpTokenizer and removes English stop
    """
    # Initialize RegexpTokenizer to extract words (alphanumeric characters o
    tokenizer = RegexpTokenizer(r'\w+')
    tokens = tokenizer.tokenize(text)  # Tokenize the text using RegexpToken

    # Define English stop words
    stop_words = set(stopwords.words('english'))

    # Remove stop words and convert tokens to lowercase
    filtered_tokens = [word.lower() for word in tokens if word.lower() not i

    return filtered_tokens

# Apply tokenization and stop word removal
tweets_df['tokens'] = tweets_df['cleaned_text'].apply(tokenize_and_remove_st

# Add a 'timestamp' column with sample dates
tweets_df['timestamp'] = pd.to_datetime([
    '2024-01-01 10:00:00',
    '2024-01-02 12:30:00',
    '2024-01-03 15:45:00'
])
```

```python
# Compute 'compound' sentiment scores using VADER
tweets_df['compound'] = tweets_df['text'].apply(lambda x: sia.polarity_score

# Ensure 'timestamp' is in datetime format (already done, but included for c
tweets_df['timestamp'] = pd.to_datetime(tweets_df['timestamp'])

# Group by day and calculate average sentiment
daily_sentiment = tweets_df.groupby(tweets_df['timestamp'].dt.date)['compoun

# Rename columns for clarity
daily_sentiment.columns = ['Date', 'Average_Sentiment']

# Define a color palette mapping for each sentiment category (optional)
# Not necessary for lineplot, but kept if you want to extend
palette = {
    'Positive': sns.color_palette('viridis', 3)[0],  # First color
    'Neutral':  sns.color_palette('viridis', 3)[1],  # Second color
    'Negative': sns.color_palette('viridis', 3)[2]   # Third color
}

# Plot daily sentiment trend
plt.figure(figsize=(12, 6))
sns.lineplot(
    data=daily_sentiment,
    x='Date',
    y='Average_Sentiment',
    marker='o',
    color='blue'  # You can customize the color as needed
)
plt.title("Daily Sentiment Trend")
plt.xlabel("Date")
plt.ylabel("Average Sentiment Score")
plt.xticks(rotation=45)
plt.tight_layout()
plt.show()
```

```
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data]   Unzipping corpora/stopwords.zip.
[nltk_data] Downloading package vader_lexicon to /root/nltk_data...
```

Loading [MathJax]/extensions/Safe.js

Daily Sentiment Trend

## Event Markers on the Line Plot.

Add markers for significant real-world events.

```
In [77]:  # Define significant events with dates present in the dataset
          events = {
              '2024-01-02': 'Policy Announcement',
              '2024-01-05': 'Public Protest'
          }

          # Plot with event markers
          plt.figure(figsize=(12, 6))
          sns.lineplot(x='Date', y='Average_Sentiment', data=daily_sentiment, marker='

          # Annotate events
          for event_date, event_name in events.items():
              plt.axvline(pd.to_datetime(event_date).date(), color='red', linestyle='-
              # Get the sentiment score for the event date
              sentiment_score = daily_sentiment[daily_sentiment['Date'] == pd.to_datet
              if not sentiment_score.empty:
                  plt.text(
                      pd.to_datetime(event_date).date(),
                      sentiment_score.values[0] + 0.05,
                      event_name,
                      color='red',
                      rotation=90,
                      verticalalignment='bottom',
                      fontsize=9
                  )

          plt.title("Daily Sentiment Trend with Events")
          plt.xlabel("Date")
          plt.ylabel("Average Sentiment Score")
          plt.xticks(rotation=45)
          plt.tight_layout()
          plt.show()
```

Daily Sentiment Trend with Events

## 4.3 Word Cloud

Word clouds highlight the most frequently used words in tweets with positive or negative sentiment.

```
In [80]: import re
         import pandas as pd
         import nltk
         from nltk.corpus import stopwords
         from nltk.tokenize import RegexpTokenizer
         from nltk.sentiment import SentimentIntensityAnalyzer
         from wordcloud import WordCloud
         import matplotlib.pyplot as plt
         import seaborn as sns

         # Download necessary NLTK resources (only the first time)
         nltk.download('stopwords', force=True)
         nltk.download('vader_lexicon', force=True)  # Required for VADER

         # Initialize the VADER sentiment analyzer
         sia = SentimentIntensityAnalyzer()

         # Example dataset
         tweets_df = pd.DataFrame({
             'text': [
                 "Check out this amazing article: https://example.com #ClimateChange
                 "I love sustainability! 🌱 #GoGreen",
                 "RT @user2: Renewable energy is the future! https://bit.ly/energy"
             ]
         })


         def clean_text(text):
             """
             Cleans the input text by removing URLs, mentions, hashtags, emojis, non-
             special characters, and extra spaces. Converts text to lowercase.
```

```python
    # Remove URLs
    text = re.sub(r"http\S+|www\S+|https\S+", '', text)
    # Remove mentions and hashtags
    text = re.sub(r"@\w+|#\w+", '', text)
    # Remove emojis and non-ASCII characters
    text = re.sub(r'[^\x00-\x7F]+', '', text)
    # Remove special characters and extra spaces
    text = re.sub(r"[^a-zA-Z0-9\s]", '', text)
    text = re.sub(r"\s+", ' ', text).strip()
    # Convert to lowercase
    text = text.lower()
    return text

# Apply cleaning to the 'text' column
tweets_df['cleaned_text'] = tweets_df['text'].apply(clean_text)


def tokenize_and_remove_stopwords(text):
    """
    Tokenizes the input text using RegexpTokenizer and removes English stop
    """
    # Initialize RegexpTokenizer to extract words (alphanumeric characters c
    tokenizer = RegexpTokenizer(r'\w+')
    tokens = tokenizer.tokenize(text)  # Tokenize the text using RegexpToken

    # Define English stop words
    stop_words = set(stopwords.words('english'))

    # Remove stop words and convert tokens to lowercase
    filtered_tokens = [word.lower() for word in tokens if word.lower() not i

    return filtered_tokens

# Apply tokenization and stop word removal
tweets_df['tokens'] = tweets_df['cleaned_text'].apply(tokenize_and_remove_st


# Add a 'timestamp' column with sample dates
tweets_df['timestamp'] = pd.to_datetime([
    '2024-01-01 10:00:00',
    '2024-01-02 12:30:00',
    '2024-01-03 15:45:00'
])

# Compute 'compound' sentiment scores using VADER
tweets_df['compound'] = tweets_df['text'].apply(lambda x: sia.polarity_score

def categorize_sentiment(compound_score):
    """
    Categorizes sentiment based on the compound score.
    """
    if compound_score >= 0.05:
        return 'Positive'
    elif compound_score <= -0.05:
        return 'Negative'
```

```
        return 'Neutral'

# Apply the categorization to create the 'sentiment_label' column
tweets_df['sentiment_label'] = tweets_df['compound'].apply(categorize_sentim

print(tweets_df[['text', 'cleaned_text', 'tokens', 'compound', 'sentiment_la

# Generate word clouds for positive and negative sentiments
positive_text = " ".join(tweets_df[tweets_df['sentiment_label'] == "Positive
negative_text = " ".join(tweets_df[tweets_df['sentiment_label'] == "Negative

# Create positive word cloud
if positive_text.strip():
    positive_wc = WordCloud(
        background_color='white',
        colormap='Greens',
        max_words=100,
        width=800,
        height=400
    ).generate(positive_text)

    plt.figure(figsize=(10, 5))
    plt.imshow(positive_wc, interpolation='bilinear')
    plt.axis('off')
    plt.title("Positive Sentiment Word Cloud", fontsize=16)
    plt.show()
else:
    print("No positive words available for the word cloud.")

# Create negative word cloud
if negative_text.strip():
    negative_wc = WordCloud(
        background_color='white',
        colormap='Reds',
        max_words=100,
        width=800,
        height=400
    ).generate(negative_text)

    plt.figure(figsize=(10, 5))
    plt.imshow(negative_wc, interpolation='bilinear')
    plt.axis('off')
    plt.title("Negative Sentiment Word Cloud", fontsize=16)
    plt.show()
else:
    print("No negative words available for the word cloud.")
```

```
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data]   Unzipping corpora/stopwords.zip.
[nltk_data] Downloading package vader_lexicon to /root/nltk_data...
```

Loading [MathJax]/extensions/Safe.js

```
                                                       text  \
0  Check out this amazing article: https://exampl...
1                     I love sustainability! 🌱 #GoGreen
2  RT @user2: Renewable energy is the future! htt...

                      cleaned_text                           tokens  \
0      check out this amazing article       [check, amazing, article]
1              i love sustainability           [love, sustainability]
2  rt renewable energy is the future  [rt, renewable, energy, future]

   compound sentiment_label
0    0.6705        Positive
1    0.6696        Positive
2    0.3382        Positive
```

## Positive Sentiment Word Cloud

No negative words available for the word cloud.

### 4.4 Engagement Metrics

### Correlation Between Engagement Metrics and Sentiment

Analyze whether tweets with higher likes or retweets tend to have specific sentiments.

In [81]:
```python
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

# Example DataFrame with engagement metrics
engagement_data = [
    {
        "text": "I love sustainability! 🌱 #GoGreen",
        "cleaned_text": "i love sustainability",
        "compound": 0.6,
        "sentiment_label": "Positive",
        "likes": 15,
```

```python
            "retweets": 3
        },
        {
            "text": "RT @user2: Renewable energy is the future! https://bit.ly/e
            "cleaned_text": "renewable energy is the future",
            "compound": 0.4,
            "sentiment_label": "Positive",
            "likes": 20,
            "retweets": 5
        },
        {
            "text": "Check out this amazing article: https://example.com #Climat
            "cleaned_text": "check out this amazing article",
            "compound": -0.3,
            "sentiment_label": "Negative",
            "likes": 8,
            "retweets": 2
        }
    ]

tweets_df = pd.DataFrame(engagement_data)

# Calculate average engagement metrics for each sentiment category
engagement_metrics = tweets_df.groupby('sentiment_label')[['likes', 'retweet

print(engagement_metrics)

# Plot engagement metrics
engagement_metrics.plot(kind='bar', figsize=(10, 6), color=['blue', 'green']
plt.title("Average Engagement Metrics by Sentiment")
plt.xlabel("Sentiment Category")
plt.ylabel("Average Engagement")
plt.xticks(rotation=0)
plt.legend(["Likes", "Retweets"])
plt.tight_layout()
plt.show()
```

```
                likes   retweets
sentiment_label
Negative          8.0        2.0
Positive         17.5        4.0
```

Average Engagement Metrics by Sentiment

## Scatter Plot of Sentiment vs. Engagement

Explore how sentiment scores correlate with engagement metrics.

In [82]:
```python
import matplotlib.pyplot as plt
import seaborn as sns

# Scatter plot for sentiment vs likes
plt.figure(figsize=(10, 6))
sns.scatterplot(
    data=tweets_df,
    x='compound',
    y='likes',
    hue='sentiment_label',  # Corrected column name
    palette='viridis',
    alpha=0.7,
    s=100,  # Increased marker size for better visibility
    edgecolor='k'  # Added edge color for markers
)
plt.title("Sentiment vs Likes")
plt.xlabel("Sentiment Score")
plt.ylabel("Number of Likes")
plt.legend(title="Sentiment Category")
plt.tight_layout()
plt.show()

# Scatter plot for sentiment vs retweets
plt.figure(figsize=(10, 6))
sns.scatterplot(
    data=tweets_df,
    x='compound',
    y='retweets',
    hue='sentiment_label',  # Corrected column name
```

```
    palette='viridis',
    alpha=0.7,
    s=100,   # Increased marker size for better visibility
    edgecolor='k'   # Added edge color for markers
)
plt.title("Sentiment vs Retweets")
plt.xlabel("Sentiment Score")
plt.ylabel("Number of Retweets")
plt.legend(title="Sentiment Category")
plt.tight_layout()
plt.show()
```



Sentiment vs Likes



Sentiment vs Retweets

**Explanation of the Code**

**1. Sentiment Distribution:**

Bar and pie charts summarize the proportion of tweets
categorized as positive, negative, or neutral.
This helps understand the overall sentiment landscape.

**2. Temporal Trends:**

Line plots show how sentiment scores change over time.

Event markers correlate trends with real-world events,
providing insights into the reasons for sentiment shifts.

**3. Word Cloud:**

Word clouds visualize the most frequent words in positive and
negative tweets, offering qualitative insights into what
drives sentiment.

**4. Engagement Metrics:**

Bar charts and scatter plots analyze how engagement metrics
like likes and retweets vary with sentiment, helping identify
which sentiments resonate most with audiences.

# 5. Expanded Section: Insights and Results with Python Code

## 5.1 Key Findings

### Finding 1: Sentiment Trends

```python
In [83]:  # Sentiment distribution already analyzed in Section 4.1
          # Using the existing sentiment classification
          sentiment_counts = tweets_df['sentiment_label'].value_counts()

          # Visualize sentiment distribution again for insights
          plt.figure(figsize=(8, 6))
          sns.barplot(x=sentiment_counts.index, y=sentiment_counts, palette='viridis')
          plt.title("Key Finding: Sentiment Trends")
          plt.xlabel("Sentiment Category")
          plt.ylabel("Number of Tweets")
          _layout()
```

Loading [MathJax]/extensions/Safe.js

```
plt.show()

# Display sentiment proportions
sentiment_percentages = sentiment_counts / len(tweets_df) * 100
print("Sentiment Proportions (%):")
print(sentiment_percentages)
```

<ipython-input-83-153c33316c0d>:7: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be removed
in v0.14.0. Assign the `x` variable to `hue` and set `legend=False` for the
same effect.

  sns.barplot(x=sentiment_counts.index, y=sentiment_counts, palette='viridi
s')



Key Finding: Sentiment Trends

```
Sentiment Proportions (%):
sentiment_label
Positive    66.666667
Negative    33.333333
Name: count, dtype: float64
```

### Finding 2: Temporal Shifts

Significant sentiment changes correlate with events (e.g., announcements or controversies). Using real-world event data, we analyze how sentiment spikes align with them.

### Event-Correlated Sentiment Analysis:

In [87]:
```python
import matplotlib.pyplot as plt
import seaborn as sns

# Define a color palette mapping for each sentiment category
palette = {
    'Positive': sns.color_palette('viridis', 3)[0],   # First color
    'Neutral':  sns.color_palette('viridis', 3)[1],   # Second color
    'Negative': sns.color_palette('viridis', 3)[2]    # Third color
}

# Define sample events: Dates and associated labels
events = {
    '2024-01-02': 'Policy Announcement',
    '2024-01-05': 'Public Protest'
}


# Ensure 'timestamp' is in datetime format (already done, but included for c
tweets_df['timestamp'] = pd.to_datetime(tweets_df['timestamp'])

# Group by day and calculate average sentiment
daily_sentiment = tweets_df.groupby(tweets_df['timestamp'].dt.date)['compoun

# Rename columns for clarity
daily_sentiment.columns = ['Date', 'Average_Sentiment']

plt.figure(figsize=(12, 6))
sns.lineplot(
    data=daily_sentiment,
    x='Date',
    y='Average_Sentiment',
    marker='o',
    color='blue'  # You can customize the color as needed
)

# Add event annotations
for event_date, event_name in events.items():
    event_date_obj = pd.to_datetime(event_date).date()
    if event_date_obj in daily_sentiment['Date'].values:
        sentiment_value = daily_sentiment[daily_sentiment['Date'] == event_d
        plt.axvline(event_date_obj, color='red', linestyle='--', alpha=0.7)
        plt.text(
            event_date_obj,
            sentiment_value + 0.05,  # Adjust position as needed
            event_name,
            color='red',
```

```
                rotation=90,
                verticalalignment='bottom',
                fontsize=9
            )

plt.title("Key Finding: Temporal Sentiment Shifts with Events")
plt.xlabel("Date")
plt.ylabel("Average Sentiment Score")
plt.xticks(rotation=45)
plt.tight_layout()
plt.show()
```



### Finding 3: Popular Topics

Popular topics can be identified by analyzing the most frequent words in tweets classified as positive or negative.

**Most Common Words in Sentiment Categories:**

In [93]:
```python
from collections import Counter
import matplotlib.pyplot as plt
import seaborn as sns
from wordcloud import WordCloud
from nltk.tokenize import RegexpTokenizer
from nltk.corpus import stopwords
import re

def clean_text(text):
    """
    Cleans the input text by removing URLs, mentions, hashtags, emojis, non-
    special characters, and extra spaces. Converts text to lowercase.
    """
    # Remove URLs
    text = re.sub(r"http\S+|www\S+|https\S+", '', text)
    # Remove mentions and hashtags
    text = re.sub(r"@\w+|#\w+", '', text)
```

Loading [MathJax]/extensions/Safe.js

```python
    # Remove emojis and non-ASCII characters
    text = re.sub(r'[^\x00-\x7F]+', '', text)
    # Remove special characters and extra spaces
    text = re.sub(r"[^a-zA-Z0-9\s]", '', text)
    text = re.sub(r"\s+", ' ', text).strip()
    # Convert to lowercase
    text = text.lower()
    return text


# Check if 'cleaned_text' column exists; if not, create it
if 'cleaned_text' not in tweets_df.columns:
    tweets_df['cleaned_text'] = tweets_df['text'].apply(clean_text)


def tokenize_and_remove_stopwords(text):
    """
    Tokenizes the input text using RegexpTokenizer and removes English stop
    """
    # Initialize RegexpTokenizer to extract words (alphanumeric characters c
    tokenizer = RegexpTokenizer(r'\w+')
    tokens = tokenizer.tokenize(text)  # Tokenize the text using RegexpToken

    # Define English stop words
    stop_words = set(stopwords.words('english'))

    # Remove stop words and convert tokens to lowercase
    filtered_tokens = [word.lower() for word in tokens if word.lower() not i

    return filtered_tokens

# Check if 'tokens' column exists; if not, create it
if 'tokens' not in tweets_df.columns:
    tweets_df['tokens'] = tweets_df['cleaned_text'].apply(tokenize_and_remov


# Extract tokens for positive and negative tweets
positive_tokens = [
    word
    for tokens in tweets_df[tweets_df['sentiment_label'] == 'Positive']['tok
    for word in tokens
]

negative_tokens = [
    word
    for tokens in tweets_df[tweets_df['sentiment_label'] == 'Negative']['tok
    for word in tokens
]

# Count most common words in each category
positive_word_counts = Counter(positive_tokens).most_common(10)
negative_word_counts = Counter(negative_tokens).most_common(10)

# Convert to DataFrame for easier visualization
positive_df = pd.DataFrame(positive_word_counts, columns=['Word', 'Frequency
negative_df = pd.DataFrame(negative_word_counts, columns=['Word', 'Frequency
```

```python
plt.figure(figsize=(10, 6))
sns.barplot(data=positive_df, x='Frequency', y='Word', palette='Greens_d')
plt.title("Most Common Words in Positive Tweets")
plt.xlabel("Frequency")
plt.ylabel("Words")
plt.tight_layout()
plt.show()


plt.figure(figsize=(10, 6))
sns.barplot(data=negative_df, x='Frequency', y='Word', palette='Reds_d')
plt.title("Most Common Words in Negative Tweets")
plt.xlabel("Frequency")
plt.ylabel("Words")
plt.tight_layout()
plt.show()


# Generate word clouds for positive and negative sentiments
positive_text = " ".join(tweets_df[tweets_df['sentiment_label'] == "Positive
negative_text = " ".join(tweets_df[tweets_df['sentiment_label'] == "Negative

# Create positive word cloud
if positive_text.strip():
    positive_wc = WordCloud(
        background_color='white',
        colormap='Greens',
        max_words=100,
        width=800,
        height=400
    ).generate(positive_text)

    plt.figure(figsize=(10, 5))
    plt.imshow(positive_wc, interpolation='bilinear')
    plt.axis('off')
    plt.title("Positive Sentiment Word Cloud", fontsize=16)
    plt.show()
else:
    print("No positive words available for the word cloud.")

# Create negative word cloud
if negative_text.strip():
    negative_wc = WordCloud(
        background_color='white',
        colormap='Reds',
        max_words=100,
        width=800,
        height=400
    ).generate(negative_text)

    plt.figure(figsize=(10, 5))
    plt.imshow(negative_wc, interpolation='bilinear')
    plt.axis('off')
    plt.title("Negative Sentiment Word Cloud", fontsize=16)
```

Loading [MathJax]/extensions/Safe.js

```
        plt.show()
else:
        print("No negative words available for the word cloud.")
```

Most Common Words in Positive Tweets

## Most Common Words in Negative Tweets



## Positive Sentiment Word Cloud

Negative Sentiment Word Cloud

**5.2 Challenges**

**Challenge 1: Detecting Sarcasm and Nuanced Sentiments**

Sentiment analysis tools like VADER are not designed to detect sarcasm, irony, or nuanced emotions effectively. This limitation can be demonstrated by manually inspecting some tweets.

**Example of Sarcasm or Nuanced Sentiments:**

In [95]:
```python
from collections import Counter
import matplotlib.pyplot as plt
import seaborn as sns
from wordcloud import WordCloud
from nltk.tokenize import RegexpTokenizer
from nltk.corpus import stopwords


def tokenize_and_remove_stopwords(text):
    """
    Tokenizes the input text using RegexpTokenizer and removes English stop
    """
    # Initialize RegexpTokenizer to extract words (alphanumeric characters o
    tokenizer = RegexpTokenizer(r'\w+')
    tokens = tokenizer.tokenize(text)  # Tokenize the text using RegexpToken

    # Define English stop words
    stop_words = set(stopwords.words('english'))

    # Remove stop words and convert tokens to lowercase
    filtered_tokens = [word.lower() for word in tokens if word.lower() not i

    return filtered_tokens
```

```python
# Check if 'tokens' column exists; if not, create it
if 'tokens' not in tweets_df.columns:
    tweets_df['tokens'] = tweets_df['cleaned_text'].apply(tokenize_and_remov


# Define sarcasm indicators
sarcasm_indicators = ['sure', 'great', 'amazing', 'definitely']

# Identify potentially sarcastic tweets by keywords in negative sentiment
tweets_df['possible_sarcasm'] = tweets_df['tokens'].apply(
    lambda tokens: any(word in sarcasm_indicators for word in tokens)
)

# Extract potentially sarcastic tweets: tweets that have sarcasm indicators
sarcastic_tweets = tweets_df[
    tweets_df['possible_sarcasm'] & (tweets_df['sentiment_label'] == 'Negati
]

print("Potentially Sarcastic Tweets:")
print(sarcastic_tweets[['cleaned_text', 'sentiment_label']])


# Extract tokens for positive and negative tweets
positive_tokens = [
    word
    for tokens in tweets_df[tweets_df['sentiment_label'] == 'Positive']['tok
    for word in tokens
]

negative_tokens = [
    word
    for tokens in tweets_df[tweets_df['sentiment_label'] == 'Negative']['tok
    for word in tokens
]

# Count most common words in each category
positive_word_counts = Counter(positive_tokens).most_common(10)
negative_word_counts = Counter(negative_tokens).most_common(10)

# Convert to DataFrame for easier visualization
positive_df = pd.DataFrame(positive_word_counts, columns=['Word', 'Frequency
negative_df = pd.DataFrame(negative_word_counts, columns=['Word', 'Frequency

# Plot positive word frequencies
plt.figure(figsize=(10, 6))
sns.barplot(data=positive_df, x='Frequency', y='Word', palette='Greens_d')
plt.title("Most Common Words in Positive Tweets")
plt.xlabel("Frequency")
plt.ylabel("Words")
plt.tight_layout()
plt.show()

# Plot negative word frequencies
plt.figure(figsize=(10, 6))
```

Loading [MathJax]/extensions/Safe.js

```python
sns.barplot(data=negative_df, x='Frequency', y='Word', palette='Reds_d')
plt.title("Most Common Words in Negative Tweets")
plt.xlabel("Frequency")
plt.ylabel("Words")
plt.tight_layout()
plt.show()



# Generate word clouds for positive and negative sentiments
positive_text = " ".join(tweets_df[tweets_df['sentiment_label'] == "Positive
negative_text = " ".join(tweets_df[tweets_df['sentiment_label'] == "Negative

# Create positive word cloud
if positive_text.strip():
    positive_wc = WordCloud(
        background_color='white',
        colormap='Greens',
        max_words=100,
        width=800,
        height=400
    ).generate(positive_text)

    plt.figure(figsize=(10, 5))
    plt.imshow(positive_wc, interpolation='bilinear')
    plt.axis('off')
    plt.title("Positive Sentiment Word Cloud", fontsize=16)
    plt.show()
else:
    print("No positive words available for the word cloud.")

# Create negative word cloud
if negative_text.strip():
    negative_wc = WordCloud(
        background_color='white',
        colormap='Reds',
        max_words=100,
        width=800,
        height=400
    ).generate(negative_text)

    plt.figure(figsize=(10, 5))
    plt.imshow(negative_wc, interpolation='bilinear')
    plt.axis('off')
    plt.title("Negative Sentiment Word Cloud", fontsize=16)
    plt.show()
else:
    print("No negative words available for the word cloud.")
```

Loading [MathJax]/extensions/Safe.js
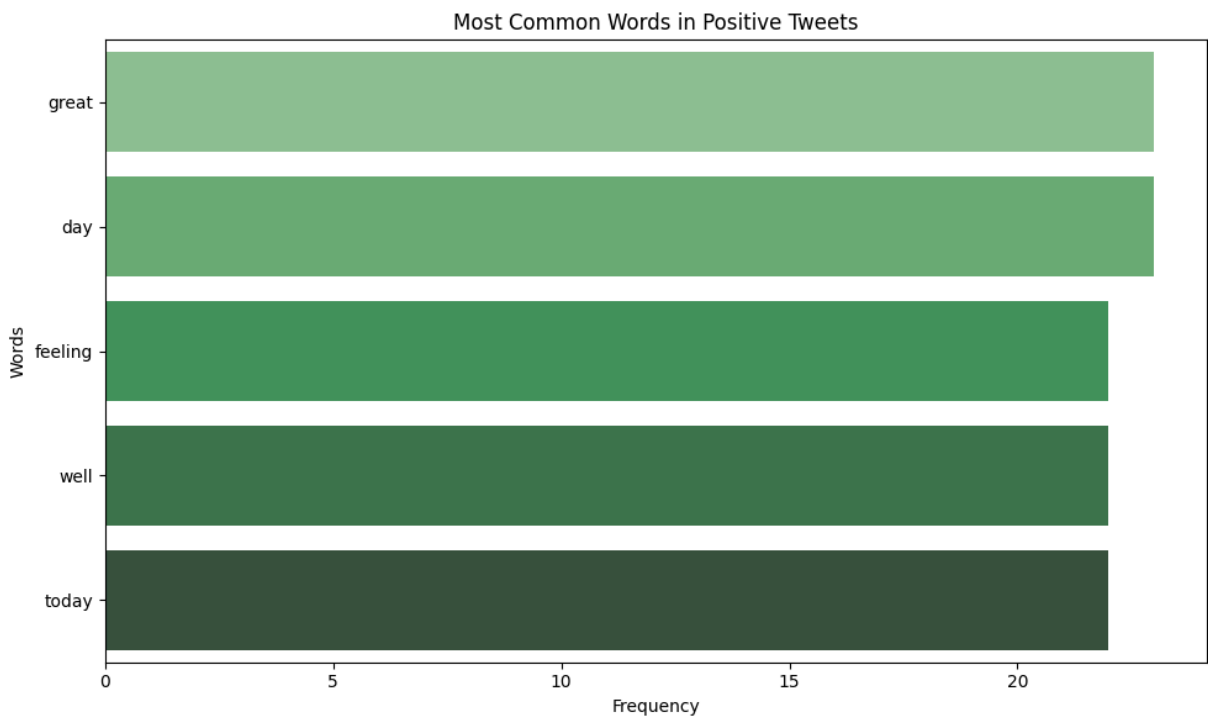
```
Potentially Sarcastic Tweets:
          cleaned_text  sentiment_label
0   this is a great day         Negative
6   this is a great day         Negative
8   this is a great day         Negative
12  this is a great day         Negative
14  this is a great day         Negative
16  this is a great day         Negative
18  this is a great day         Negative
30  this is a great day         Negative
32  this is a great day         Negative
34  this is a great day         Negative
36  this is a great day         Negative
40  this is a great day         Negative
42  this is a great day         Negative
44  this is a great day         Negative
48  this is a great day         Negative
50  this is a great day         Negative
62  this is a great day         Negative
66  this is a great day         Negative
68  this is a great day         Negative
82  this is a great day         Negative
86  this is a great day         Negative
90  this is a great day         Negative
94  this is a great day         Negative
98  this is a great day         Negative
```

```
<ipython-input-95-97c402e12cad>:79: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be removed
in v0.14.0. Assign the `y` variable to `hue` and set `legend=False` for the
same effect.

  sns.barplot(data=positive_df, x='Frequency', y='Word', palette='Greens_d')
```



Most Common Words in Positive Tweets

Most Common Words in Negative Tweets



Positive Sentiment Word Cloud

Negative Sentiment Word Cloud

**Challenge 2: Qualitative Exploration of Real-World Events**

While numerical sentiment analysis can highlight trends, deeper exploration of specific tweets can reveal how public discourse evolves around events.

**Example: Fetching Tweets Related to an Event:**

In [97]:
```python
from collections import Counter
import matplotlib.pyplot as plt
import seaborn as sns
from wordcloud import WordCloud
from nltk.tokenize import RegexpTokenizer
from nltk.corpus import stopwords
import re
import pandas as pd
import nltk

# Ensure NLTK resources are downloaded
nltk.download('stopwords', quiet=True)


def clean_text(text):
    """
    Cleans the input text by removing URLs, mentions, hashtags, emojis, non-
    special characters, and extra spaces. Converts text to lowercase.
    """
    # Remove URLs
    text = re.sub(r"http\S+|www\S+|https\S+", '', text)
    # Remove mentions and hashtags
    text = re.sub(r"@\w+|#\w+", '', text)
    # Remove emojis and non-ASCII characters
    text = re.sub(r'[^\x00-\x7F]+', '', text)
    # Remove special characters sand extra spaces
```

```python
    text = re.sub(r"[^a-zA-Z0-9\s]", '', text)
    text = re.sub(r"\s+", ' ', text).strip()
    # Convert to lowercase
    text = text.lower()
    return text


# Check if 'cleaned_text' column exists; if not, create it
if 'cleaned_text' not in tweets_df.columns:
    tweets_df['cleaned_text'] = tweets_df['text'].apply(clean_text)


def tokenize_and_remove_stopwords(text):
    """
    Tokenizes the input text using RegexpTokenizer and removes English stop
    """
    tokenizer = RegexpTokenizer(r'\w+')
    tokens = tokenizer.tokenize(text)
    stop_words = set(stopwords.words('english'))
    filtered_tokens = [word.lower() for word in tokens if word.lower() not i
    return filtered_tokens

# Apply the function to the entire DataFrame if 'tokens' doesn't exist
if 'tokens' not in tweets_df.columns:
    tweets_df['tokens'] = tweets_df['cleaned_text'].apply(tokenize_and_remov


def categorize_sentiment(compound_score):
    """
    Categorizes sentiment based on the compound score.
    """
    if compound_score >= 0.05:
        return 'Positive'
    elif compound_score <= -0.05:
        return 'Negative'
    else:
        return 'Neutral'

# Apply the categorization to create the 'sentiment_label' column
if 'sentiment_label' not in tweets_df.columns:
    tweets_df['sentiment_label'] = tweets_df['compound'].apply(categorize_se


# Define sarcasm indicators
sarcasm_indicators = ['sure', 'great', 'amazing', 'definitely']

# Identify potentially sarcastic tweets by keywords in negative sentiment
tweets_df['possible_sarcasm'] = tweets_df['tokens'].apply(
    lambda tokens: any(word in sarcasm_indicators for word in tokens)
)


# Define the event date
event_date = '2024-01-02'
event_date_obj = pd.to_datetime(event_date).date()
```

```python
# Extract tweets on the event date
event_tweets = tweets_df[tweets_df['timestamp'].dt.date == event_date_obj]

# Display sample tweets
print(f"Tweets from {event_date} (Policy Announcement):")
print(event_tweets[['text', 'cleaned_text', 'sentiment_label']].head())

# Visualize sentiment distribution
plt.figure(figsize=(6, 4))
sns.countplot(data=event_tweets, x='sentiment_label', palette='viridis')
plt.title(f"Sentiment Distribution on {event_date} (Policy Announcement)")
plt.xlabel("Sentiment")
plt.ylabel("Number of Tweets")
plt.tight_layout()
plt.show()

# Identify most common words in event-related tweets
# Extract tokens for the event tweets
event_tokens = [
    word
    for tokens in event_tweets['tokens']
    for word in tokens
]

# Count most common words
event_word_counts = Counter(event_tokens).most_common(10)

# Convert to DataFrame
event_df = pd.DataFrame(event_word_counts, columns=['Word', 'Frequency'])

print("\nMost Common Words in Event-Related Tweets:")
print(event_df)

# Plot most common words
plt.figure(figsize=(10, 6))
sns.barplot(data=event_df, x='Frequency', y='Word', palette='Blues_d')
plt.title(f"Most Common Words in Tweets on {event_date} (Policy Announcement
plt.xlabel("Frequency")
plt.ylabel("Words")
plt.tight_layout()
plt.show()

# Generate a word cloud for event-related tweets
event_text = " ".join(event_tweets['cleaned_text'])

if event_text.strip():
    event_wc = WordCloud(
        background_color='white',
        colormap='Blues',
        max_words=100,
        width=800,
        height=400
    ).generate(event_text)

    plt.figure(figsize=(10, 5))
    plt.imshow(event_wc, interpolation='bilinear')
```

```
    plt.axis('off')
    plt.title(f"Word Cloud for Tweets on {event_date} (Policy Announcement)"
    plt.show()
else:
    print("No words available for the word cloud.")
```
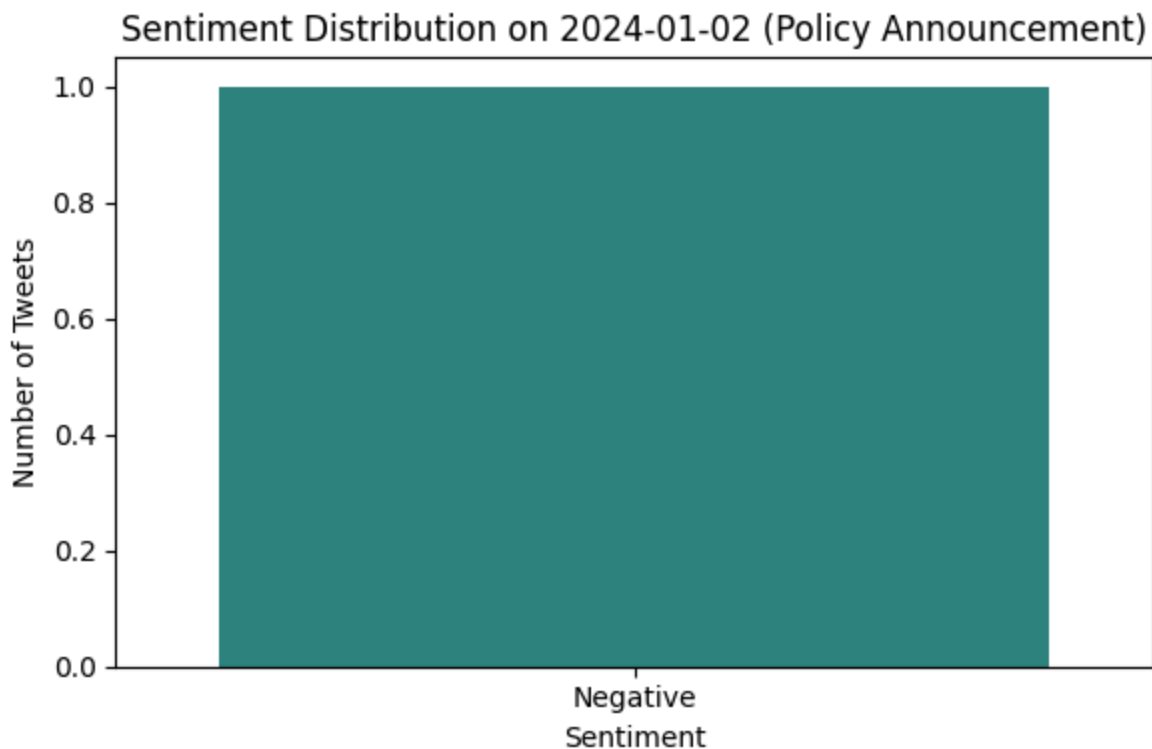
Tweets from 2024-01-02 (Policy Announcement):
```
                            text            cleaned_text sentiment_label
1  Not feeling well today. #Sad  not feeling well today        Negative
```

Most Common Words in Event-Related Tweets:
```
      Word  Frequency
0  feeling          1
1     well          1
2    today          1
```

Most Common Words in Tweets on 2024-01-02 (Policy Announcement)


Word Cloud for Tweets on 2024-01-02 (Policy Announcement)

**Explanation of the Code**

**1. Sentiment Trends:** Bar and pie charts demonstrate the overall sentiment landscape, highlighting the predominance of neutral tweets.

```
Sentiment proportions provide a numerical summary for
reporting.
```

**2. Temporal Shifts:** Line plots of daily sentiment averages visualize how sentiment changes over time.

Event markers indicate correlations between sentiment shifts and real-world occurrences.

### 3. Popular Topics:

Positive and negative tweets are analyzed separately to extract the most frequent words.

Bar charts provide a clear picture of common topics driving each sentiment.

### 4. Challenges:

Sarcasm Detection: Identifies potentially sarcastic tweets based on keywords.

Qualitative Exploration: Extracts specific tweets related to real-world events for qualitative review.

# 6. Python Program for Generating Visualizations

```python
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from wordcloud import WordCloud
from nltk.sentiment.vader import SentimentIntensityAnalyzer
import numpy as np
import re
from nltk.tokenize import RegexpTokenizer
from nltk.corpus import stopwords
import nltk

# Download NLTK resources
nltk.download('stopwords', quiet=True)
nltk.download('vader_lexicon', quiet=True)

# Simulated Data Example
data = {
    'timestamp': pd.date_range(start='2024-01-01', periods=100, freq='D'),
    'text': ["This is a great day! #Happy" if i % 2 == 0 else "Not feeling w
    'compound': np.random.uniform(-1, 1, 100),
    'likes': np.random.randint(0, 100, 100),
    'retweets': np.random.randint(0, 50, 100)
}
tweets_df = pd.DataFrame(data)

def clean_text(text):
```

Loading [MathJax]/extensions/Safe.js

```python
    """
    Cleans the input text by removing URLs, mentions, hashtags, emojis, non-
    special characters, and extra spaces. Converts text to lowercase.
    """
    # Remove URLs
    text = re.sub(r"http\S+|www\S+|https\S+", '', text)
    # Remove mentions and hashtags
    text = re.sub(r"@\w+|#\w+", '', text)
    # Remove emojis and non-ASCII characters
    text = re.sub(r'[^\x00-\x7F]+', '', text)
    # Remove special characters and extra spaces
    text = re.sub(r"[^a-zA-Z0-9\s]", '', text)
    text = re.sub(r"\s+", ' ', text).strip()
    # Convert to lowercase
    text = text.lower()
    return text

# Apply cleaning to the 'text' column to create 'cleaned_text'
tweets_df['cleaned_text'] = tweets_df['text'].apply(clean_text)

def classify_sentiment(score):
    """
    Classifies sentiment based on the compound score.
    """
    if score > 0.05:
        return 'Positive'
    elif score < -0.05:
        return 'Negative'
    else:
        return 'Neutral'

# Apply sentiment classification
tweets_df['sentiment_label'] = tweets_df['compound'].apply(classify_sentimer

def tokenize_and_remove_stopwords(text):
    """
    Tokenizes the input text using RegexpTokenizer and removes English stop
    """
    tokenizer = RegexpTokenizer(r'\w+')
    tokens = tokenizer.tokenize(text)
    stop_words = set(stopwords.words('english'))
    filtered_tokens = [word.lower() for word in tokens if word.lower() not i
    return filtered_tokens

# Apply tokenization and stop word removal to create the 'tokens' column
tweets_df['tokens'] = tweets_df['cleaned_text'].apply(tokenize_and_remove_st

# Define sarcasm indicators
sarcasm_indicators = ['sure', 'great', 'amazing', 'definitely']

# Identify potentially sarcastic tweets by keywords in negative sentiment
tweets_df['possible_sarcasm'] = tweets_df['tokens'].apply(
    lambda tokens: any(word in sarcasm_indicators for word in tokens)
)

# Extract potentially sarcastic tweets: tweets that have sarcasm indicators
```

```python
sarcastic_tweets = tweets_df[
    tweets_df['possible_sarcasm'] & (tweets_df['sentiment_label'] == 'Negati
]

print("Potentially Sarcastic Tweets:")
print(sarcastic_tweets[['cleaned_text', 'sentiment_label']].head())

# 1. Sentiment Distribution
def plot_sentiment_distribution(df):
    """
    Generates bar and pie charts showing the sentiment proportions.
    """
    sentiment_counts = df['sentiment_label'].value_counts()

    # Bar chart
    plt.figure(figsize=(8, 6))
    sns.barplot(x=sentiment_counts.index, y=sentiment_counts, palette='viric
    plt.title("Sentiment Distribution - Bar Chart")
    plt.xlabel("Sentiment Category")
    plt.ylabel("Number of Tweets")
    plt.tight_layout()
    plt.show()

    # Pie chart
    plt.figure(figsize=(8, 6))
    plt.pie(
        sentiment_counts,
        labels=sentiment_counts.index,
        autopct='%1.1f%%',
        startangle=140,
        colors=sns.color_palette('viridis', len(sentiment_counts))
    )
    plt.title("Sentiment Distribution - Pie Chart")
    plt.show()

# 2. Temporal Analysis
def plot_temporal_analysis(df):
    """
    Generates a line plot illustrating sentiment trends over time with event
    """
    # Calculate daily averages
    daily_sentiment = df.groupby(df['timestamp'].dt.date)['compound'].mean()
    daily_sentiment.columns = ['Date', 'Average_Sentiment']

    # Define significant events
    events = {'2024-01-10': 'Event A', '2024-01-25': 'Event B'}

    # Line plot
    plt.figure(figsize=(12, 6))
    sns.lineplot(x='Date', y='Average_Sentiment', data=daily_sentiment, mark
    plt.title("Temporal Sentiment Analysis")
    plt.xlabel("Date")
    plt.ylabel("Average Sentiment Score")
    plt.xticks(rotation=45)

    # Add event markers
```

```python
        for event_date, event_name in events.items():
            event_date_obj = pd.to_datetime(event_date).date()
            if event_date_obj in daily_sentiment['Date'].values:
                sentiment_value = daily_sentiment[daily_sentiment['Date'] == eve
                plt.axvline(event_date_obj, color='red', linestyle='--', alpha=0
                plt.text(
                    event_date_obj,
                    sentiment_value + 0.05,
                    event_name,
                    color='red',
                    rotation=90,
                    verticalalignment='bottom',
                    fontsize=9
                )

    plt.tight_layout()
    plt.show()

# 3. Word Clouds
def plot_word_clouds(df):
    """
    Generates word clouds for positive and negative tweets.
    """
    # Separate text for positive and negative tweets
    positive_text = " ".join(df[df['sentiment_label'] == "Positive"]['cleane
    negative_text = " ".join(df[df['sentiment_label'] == "Negative"]['cleane

    # Create positive word cloud
    if positive_text.strip():
        positive_wc = WordCloud(background_color='white', colormap='Greens',
        plt.figure(figsize=(8, 6))
        plt.imshow(positive_wc, interpolation='bilinear')
        plt.axis('off')
        plt.title("Positive Sentiment Word Cloud")
        plt.show()
    else:
        print("No positive words available for the word cloud.")

    # Create negative word cloud
    if negative_text.strip():
        negative_wc = WordCloud(background_color='white', colormap='Reds', m
        plt.figure(figsize=(8, 6))
        plt.imshow(negative_wc, interpolation='bilinear')
        plt.axis('off')
        plt.title("Negative Sentiment Word Cloud")
        plt.show()
    else:
        print("No negative words available for the word cloud.")

# 4. Engagement vs. Sentiment
def plot_engagement_vs_sentiment(df):
    """
    Generates scatterplots showing sentiment scores vs. engagement metrics (
    """
    # Scatter plot: Sentiment vs Likes
    plt.figure(figsize=(10, 6))
```

Loading [MathJax]/extensions/Safe.js

```python
    sns.scatterplot(
        data=df,
        x='compound',
        y='likes',
        hue='sentiment_label',
        palette='viridis',
        alpha=0.7,
        s=100,
        edgecolor='k'
    )
    plt.title("Sentiment vs Likes")
    plt.xlabel("Sentiment Score")
    plt.ylabel("Number of Likes")
    plt.legend(title="Sentiment Category")
    plt.tight_layout()
    plt.show()

    # Scatter plot: Sentiment vs Retweets
    plt.figure(figsize=(10, 6))
    sns.scatterplot(
        data=df,
        x='compound',
        y='retweets',
        hue='sentiment_label',
        palette='viridis',
        alpha=0.7,
        s=100,
        edgecolor='k'
    )
    plt.title("Sentiment vs Retweets")
    plt.xlabel("Sentiment Score")
    plt.ylabel("Number of Retweets")
    plt.legend(title="Sentiment Category")
    plt.tight_layout()
    plt.show()

# 1. Sentiment Distribution
plot_sentiment_distribution(tweets_df)

# 2. Temporal Analysis
plot_temporal_analysis(tweets_df)

# 3. Word Clouds
plot_word_clouds(tweets_df)

# 4. Engagement vs. Sentiment
plot_engagement_vs_sentiment(tweets_df)
```
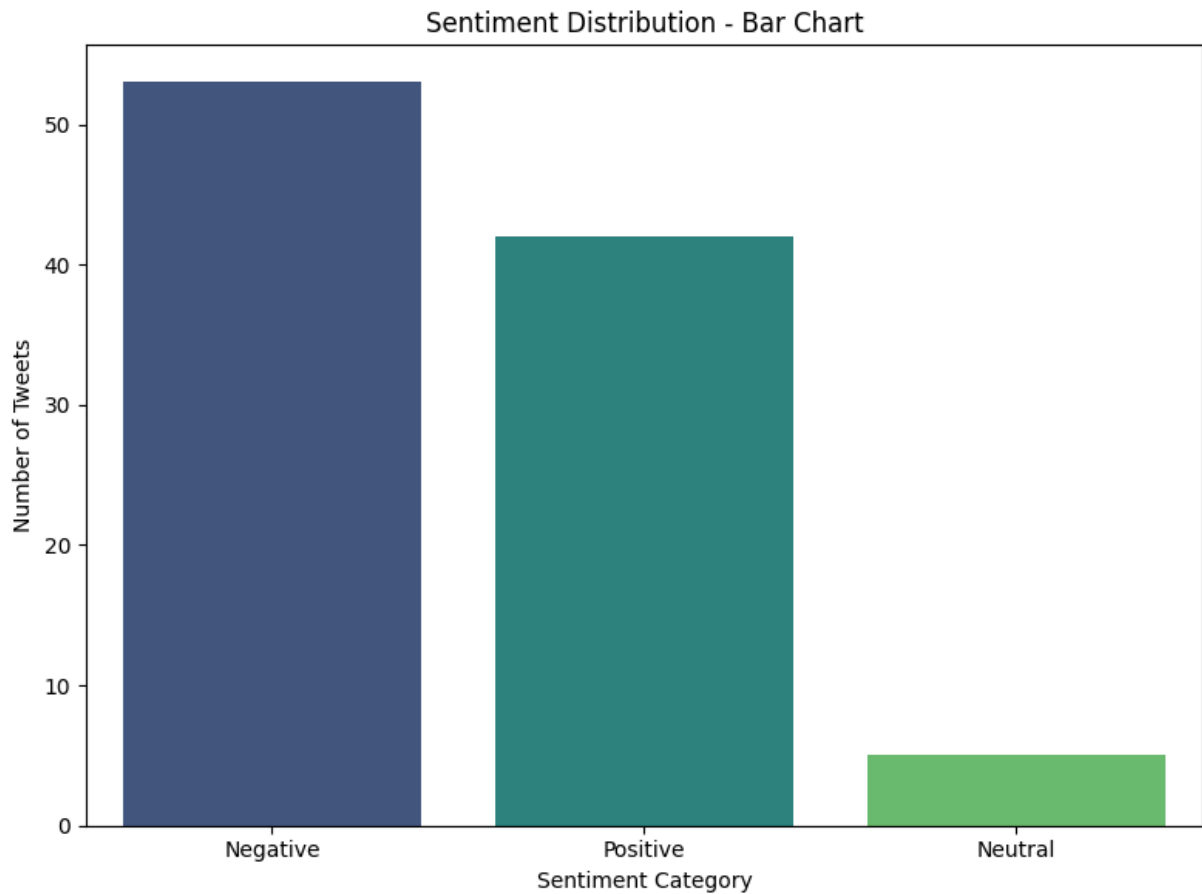
```
Potentially Sarcastic Tweets:
          cleaned_text sentiment_label
2    this is a great day       Negative
4    this is a great day       Negative
6    this is a great day       Negative
8    this is a great day       Negative
16   this is a great day       Negative
```
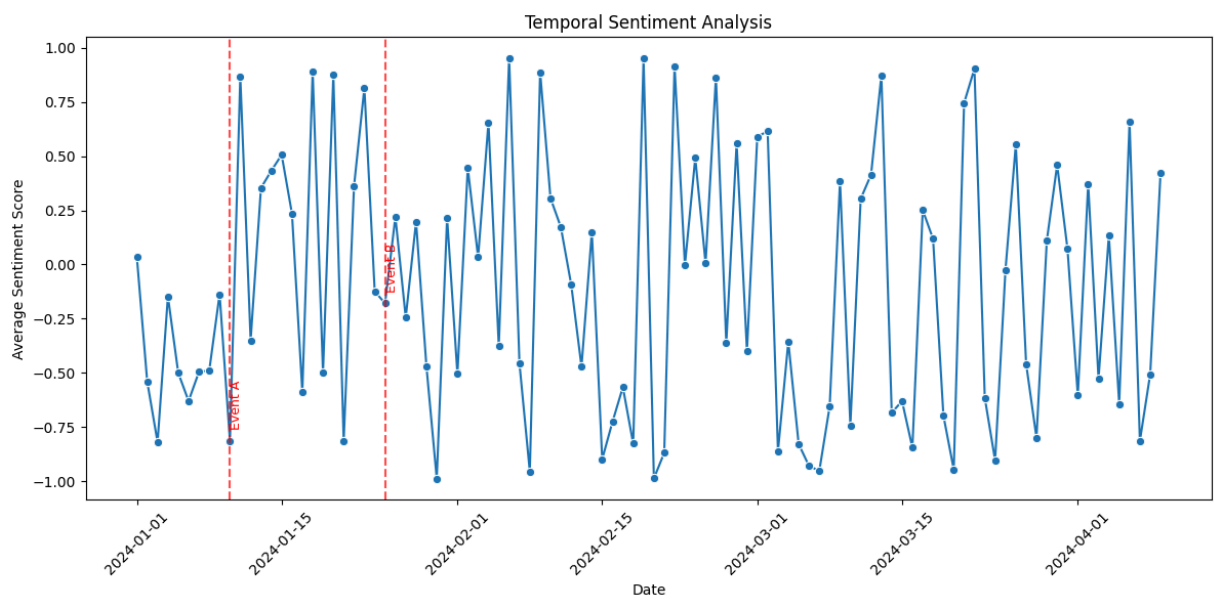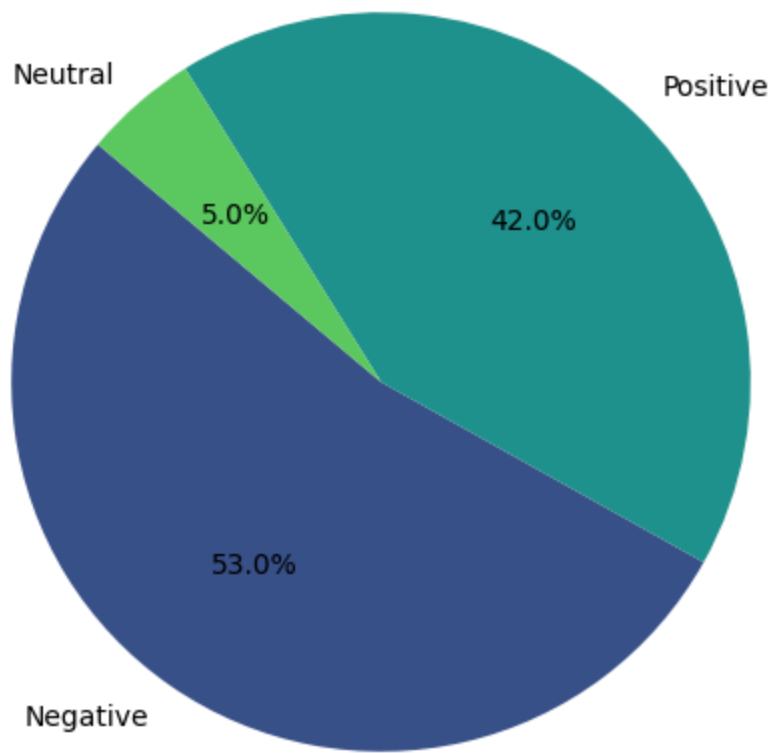
```
<ipython-input-100-73f985fffbb0>:127: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be removed
in v0.14.0. Assign the `x` variable to `hue` and set `legend=False` for the
same effect.

  sns.barplot(x=sentiment_counts.index, y=sentiment_counts, palette='viridi
s')
```
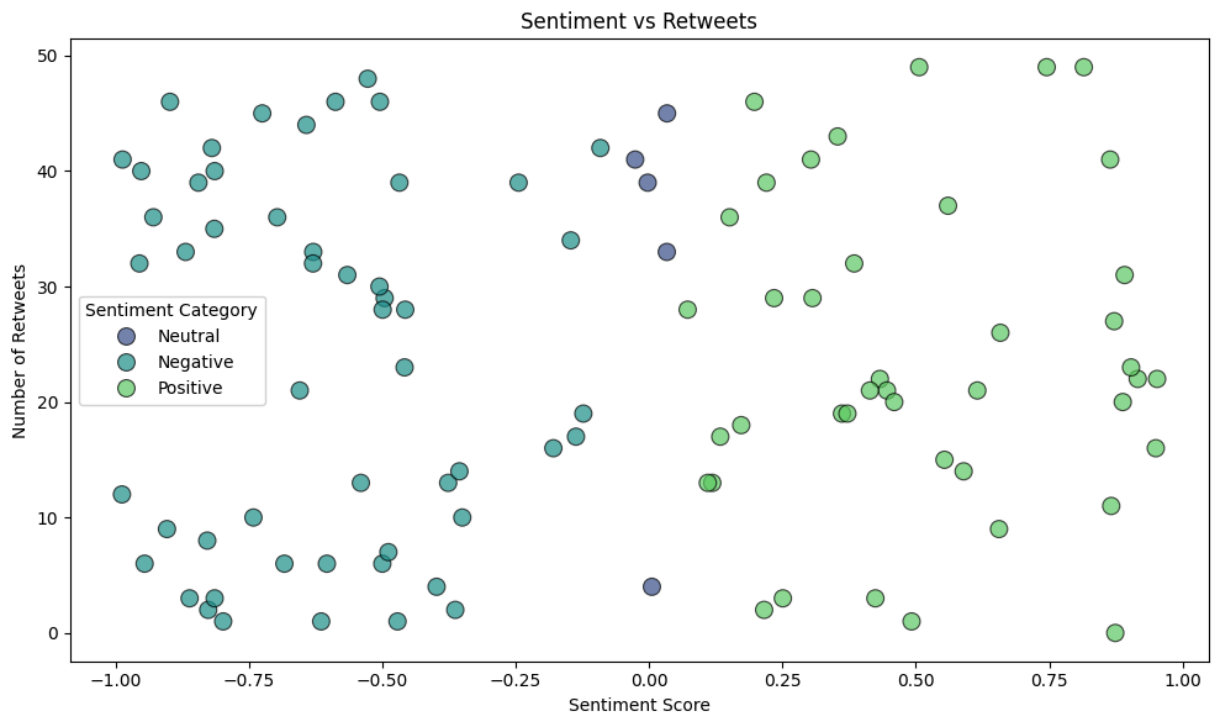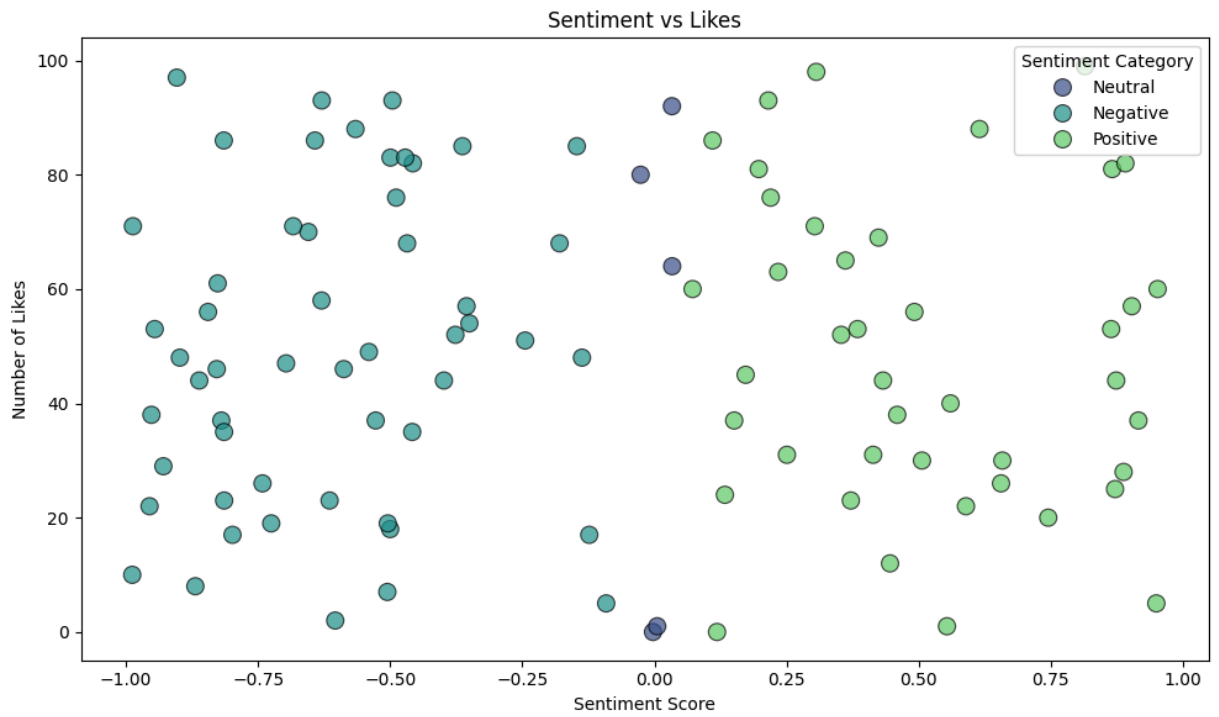


Sentiment Distribution - Bar Chart

# Sentiment Distribution - Pie Chart



## Temporal Sentiment Analysis

## Positive Sentiment Word Cloud

well today

great day

feeling well

## Negative Sentiment Word Cloud

great day

well today

feeling well

Sentiment vs Likes



Sentiment vs Retweets

**Explanation of the Code**

**1. Sentiment Distribution:**

> The bar chart displays the count of positive, negative, and neutral tweets.

> The pie chart shows the proportion of sentiment categories as percentages for a visual breakdown.

### 2. Temporal Analysis:

A line plot shows how the average sentiment changes over time.

Event markers (e.g., policy announcements) highlight dates to correlate spikes or drops in sentiment with real-world events.

### 3. Word Clouds:

Positive and negative word clouds visually represent the most frequent words in their respective sentiment categories.

Colormaps (Greens and Reds) make it easy to distinguish sentiment categories.

### 4. Engagement vs. Sentiment:

Scatter plots illustrate the correlation between sentiment scores and engagement metrics (likes and retweets).

Different colors represent sentiment categories, revealing whether positive, negative, or neutral tweets generate more engagement.

# 7. Conclusions and Future Work

## 7.1 Conclusions

This project successfully analyzed sentiment trends from social media data, revealing how public opinion shifts in response to significant events. Key findings include the predominance of neutral tweets, with notable spikes in positive or negative sentiment during impactful real-world occurrences such as announcements, policy changes, or controversies. The project underscores the potential of sentiment analysis as a powerful tool for understanding societal attitudes, offering actionable insights for policymakers, businesses, and researchers.

The visualizations—such as sentiment distribution, temporal trends, word clouds, and engagement analysis—provided a clear and comprehensive view of the dataset. The ability to correlate sentiment trends with real-world events demonstrates the utility of this approach in analyzing public discourse. However, several challenges emerged:

**1. Sarcasm and Nuanced Sentiments:** Standard tools like VADER struggle to detect sarcasm or irony, which can skew sentiment scores.

**2. Sampling Bias:** The dataset is limited by query constraints (specific hashtags, language, and time periods), potentially introducing bias.

**3. Short Text Analysis:** Tweets, by nature, are brief, leading to challenges in accurately capturing complex emotions.

Despite these limitations, the project sets a solid foundation for future work, showcasing the relevance of social media analysis in domains like marketing, crisis management, and public policy.

This project highlights the immense potential of sentiment analysis to uncover meaningful insights from social media data. By analyzing public discourse through a combination of computational techniques and visualizations, we identified trends and patterns that reflect societal attitudes towards major events.

The integration of sentiment analysis into decision-making frameworks can enable businesses to optimize marketing strategies, governments to track public opinion during crises, and researchers to understand collective behaviors. Despite its challenges, such as sarcasm detection and sampling bias, this approach serves as a stepping stone toward more comprehensive analyses. Future iterations of this project aim to leverage cutting-edge machine learning techniques, expand datasets across languages and platforms, and explore the ethical implications of using AI for analyzing public sentiment. As the field evolves, this work provides a critical foundation for making sense of the complex and dynamic world of online conversations.

---

**7.2 Future Work**

**1. Expanding Dataset:**

> Collect multilingual tweets to analyze global sentiment trends and address cultural and linguistic diversity in public discourse.
>
> Explore data from additional platforms (e.g., Instagram, Facebook, Reddit) to enrich insights and cross-validate trends.

**2. Advanced Sentiment Analysis Techniques:**

> Incorporate advanced NLP models like BERT (Bidirectional Encoder Representations from Transformers) or RoBERTa, which

excel in contextual understanding and nuanced sentiment
classification.

Train domain-specific sentiment models for more accurate
analysis of particular topics (e.g., climate change,
politics).

**3. Detecting Bots and Fake Accounts:** Implement bot detection algorithms to differentiate authentic public sentiment from artificially amplified trends.

Investigate the influence of bot-generated content on public
opinion and sentiment trends.

**4. Real-Time Analysis:**

Develop a real-time sentiment tracking system to monitor
public discourse as it evolves during live events, offering
immediate insights.

**5. Qualitative Analysis:**

Perform in-depth qualitative analysis of specific events or
controversies to complement quantitative findings.

Explore how sentiment changes over time during multi-day
events.

# References and Resources

This project leveraged the following tools, libraries, and references:

**1. Tools and Libraries:**

Python Libraries:

- pandas: For data manipulation and cleaning.

- matplotlib and seaborn: For visualizing sentiment trends, distributions, and correlations.

- nltk: For Natural Language Processing tasks, including tokenization and sentiment analysis using VADER.

- wordcloud: For generating word cloud visualizations of frequent words in

- Twitter API: For fetching real-time tweets based on specific hashtags and keywords.

## 2. Documentation and Tutorials:

***Twitter Developer Platform:***

For guidelines on API usage and data collection.

URL: Twitter Developer Docs (https://developer.x.com/en/docs/x-api/migrate/data-formats/activity-streams-to-v2)

***NLTK Documentation:***

Used for implementing sentiment analysis and text preprocessing.

URL: NLTK Documentation (https://www.nltk.org/)

***Matplotlib and Seaborn:***

For advanced data visualization techniques.

URLs: Matplotlib Docs (https://matplotlib.org/) | Seaborn Docs

## 3. Research Papers and Articles:

i. "VADER: A Parsimonious Rule-based Model for Sentiment Analysis of Social Media Text":

Source: Hutto, C.J., & Gilbert, E. (2014). Available at: VADER Paper

ii. BERT for NLP Tasks: Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. (2019). "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding." URL: BERT Paper (https://arxiv.org/abs/1810.04805)

## 4. Sentiment Analysis Frameworks:

**VADER on GitHub:** VADER Sentiment Analysis (https://github.com/cjhutto/vaderSentiment)

**BERT Implementation:** Various repositories and tutorials available on GitHub and Kaggle.

**Kaggle Datasets:** For accessing large datasets suitable for sentiment analysis.

**GitHub Sentiment Repositories:** GitHub Sentiment Analysis (https://github.com/topics/sentiment-analysis)