

▼ Datasets for further Analysis - Preprocessed

1 #Extracting the Preproceed Data for further Analysis

```

1 #Bitcoin Price
2 import pandas as pd
3
4 # URL to the raw CSV file
5 url = 'https://raw.githubusercontent.com/Amarpreet3/CIND-820-CAPSTONE/main/Sentimental%20Analysis/BitcoinPricePreprocessed.csv'
6
7 # Read the CSV file from the URL
8 crypto_usd = pd.read_csv(url)
9
10 # Display the first few rows of the data
11 print(crypto_usd.head())
12
13
```

	time	close	high	low	open	volumefrom	\
0	2023-02-19 13:00:00	24682.03	24715.82	24682.03	24707.39	903.97	
1	2023-02-19 14:00:00	24765.79	24792.85	24679.21	24682.03	1220.29	
2	2023-02-19 15:00:00	24928.21	25022.49	24751.96	24765.79	5074.50	
3	2023-02-19 16:00:00	24786.44	25175.28	24704.53	24928.21	7094.72	
4	2023-02-19 17:00:00	24364.95	24806.64	24346.17	24786.44	6896.84	

	volumeto	Date	Time	volume	marketcap	price_delta
0	2.233594e+07	2023-02-19	13:00:00	2.233504e+07	5.512964e+11	NaN
1	3.020300e+07	2023-02-19	14:00:00	3.020178e+07	7.480012e+11	83.76
2	1.263085e+08	2023-02-19	15:00:00	1.263034e+08	3.148644e+12	162.42
3	1.770671e+08	2023-02-19	16:00:00	1.770600e+08	4.388863e+12	-141.77
4	1.693379e+08	2023-02-19	17:00:00	1.693310e+08	4.125910e+12	-421.49

```

1 import pandas as pd
2
3 file_urls = [
4     'https://github.com/Amarpreet3/CIND-820-CAPSTONE/raw/main/Sentimental%20Analysis/BitcoinTweetsPreprocessed_1.csv',
5     'https://github.com/Amarpreet3/CIND-820-CAPSTONE/raw/main/Sentimental%20Analysis/BitcoinTweetsPreprocessed_2.csv',
6     'https://github.com/Amarpreet3/CIND-820-CAPSTONE/raw/main/Sentimental%20Analysis/BitcoinTweetsPreprocessed_3.csv',
7     'https://github.com/Amarpreet3/CIND-820-CAPSTONE/raw/main/Sentimental%20Analysis/BitcoinTweetsPreprocessed_4.csv',
8     'https://github.com/Amarpreet3/CIND-820-CAPSTONE/raw/main/Sentimental%20Analysis/BitcoinTweetsPreprocessed_5.csv',
9     'https://github.com/Amarpreet3/CIND-820-CAPSTONE/raw/main/Sentimental%20Analysis/BitcoinTweetsPreprocessed_6.csv'
10 ]
11
12 dfs = []
13
14 for url in file_urls:
15     # Read the CSV file
16     df = pd.read_csv(url)
17
18     # Append the DataFrame to the list
19     dfs.append(df)
20
21 # Combine all DataFrames into a single DataFrame
22 combined_df = pd.concat(dfs)
23
24 # Display the first few rows of the combined DataFrame
25 print(combined_df.head())
26
```

	user_name	user_location	\
0	Irk	Vancouver, WA	
1	Xiang Zhang	NaN	
2	Rhizoo	NaN	
3	Hari Marquez	Las Vegas, NV	
4	Bitcoin Candle Bot	Brazil	

	user_description	user_created	\
0	Irk started investing in 1...	2018-08-11 03:17:00	
1	Professional Software Engineer 🚀Crypto ...	2011-01-11 01:37:00	
2	researcher. local maxima dunningàkruger spec...	2019-04-03 18:09:00	
3	Donât trust, verify. #Bitcoin El Salvador ...	2014-01-17 23:04:00	
4	Robot that posts the closure of the bitcoin da...	2021-01-06 01:36:00	

	user_followers	user_friends	user_favourites	user_verified	\
--	----------------	--------------	-----------------	---------------	---

```

0      116.0        8.0     4580.0    False
1      42.0         22.0       5.0    False
2     778.0        627.0   32005.0    False
3     222.0        521.0  13052.0    False
4      40.0          4.0       1.0    False

           date                  text \
0  2023-02-25 23:59:00 bitcoin btc rest crypto ye bitcoin cryptocurr ...
1  2023-02-25 23:59:00 retriev invest fund current ongo tidexcoin kic...
2  2023-02-25 23:59:00 bull save monthli thread today good shit bitco...
3  2023-02-25 23:59:00      el salvador shape futur bitcoin membvk32cn
4  2023-02-25 23:59:00 cndl day 25022023 close open 2319406 high 232...

           hashtags            source \
0  ['Bitcoin', 'crypto', 'NeedsMoreCrash']  Twitter Web App
1  ['Tidexcoin', 'Kicurrency', 'LMY', 'GMK', 'SYR...  Twitter for iPhone
2  ['bitcoin']  Twitter Web App
3  ['Bitcoin']  Twitter Web App
4  ['Bitcoin', 'Candle', 'BearMarket']  Bitcoin Candle Bot

  is_retweet compound      score sentiment_level polarity subjectivity
0      0.0  -0.4019 -2.154092e+05      Negative  0.000000  0.000000
1      0.0   0.0000  0.000000e+00      Neutral  0.000000  0.400000
2      0.0   0.3612  9.005682e+06      Positive  0.250000  0.700000
3      0.0   0.0000  0.000000e+00      Neutral  0.000000  0.000000
4      0.0  -0.2732 -2.240240e+01      Negative  0.053333  0.446667

```

```
1 tweets = combined_df.copy()
```

```
1 tweets.head()
```

0	Irk	Vancouver, WA	Irk started investing in the stock market in 1...	2018-08-11 03:17:00	116.0	8.0	4580.0	False	2023-02-25 23:59:00	b re y cry		
1	Xiang Zhang	NaN	Professional Software Engineer ð¤»ð¤»Crypto ...	2011-01-11 01:37:00	42.0	22.0	5.0	False	2023-02-25 23:59:00	retri fun		
2	Rhizoo	NaN	researcher. local maxima dunningâ¬kruger spec...	2019-04-03 18:09:00	778.0	627.0	32005.0	False	2023-02-25 23:59:00	thre !		
3	Hari Marquez	Las Vegas, NV	Donâ¬t trust, verify. #Bitcoin El Salvador ...	2014-01-17 23:04:00	222.0	521.0	13052.0	False	2023-02-25 23:59:00	el sh mer		
4	Bitcoin Candle Bot	Brazil	Robot that posts the closure of the bitcoin da...	2021-01-06 01:36:00	40.0	4.0	1.0	False	2023-02-25 23:59:00	c 2 ck hi		



```
1 print(tweets.columns)
```

```

Index(['user_name', 'user_location', 'user_description', 'user_created',
       'user_followers', 'user_friends', 'user_favourites', 'user_verified',
       'date', 'text', 'hashtags', 'source', 'is_retweet', 'compound', 'score',
       'sentiment_level', 'polarity', 'subjectivity'],
      dtype='object')

```

```

1 import pandas as pd
2
3
4 # Check the shape of the dataset
5 print("Shape of the dataset:", tweets.shape)
6
7 # Check the size of the dataset

```

```
8 print("Size of the dataset (number of elements):", tweets.size)
```

```
9
```

```
Shape of the dataset: (167652, 18)
```

```
Size of the dataset (number of elements): 3017736
```

```
1 import pandas as pd
```

```
2 import os
```

```
3
```

```
4
```

```
5 # Check the shape of the data
```

```
6 print("Shape of the data:", tweets.shape)
```

```
7
```

```
Shape of the data: (167652, 18)
```

```
1 label_counts = tweets['sentiment_level'].value_counts()
```

```
2 print(label_counts)
```

```
Neutral 93169
```

```
Positive 35921
```

```
Extreme Positive 17343
```

```
Negative 15903
```

```
Extreme Negative 5316
```

```
Name: sentiment_level, dtype: int64
```

```
1 #Crypto - Bitcoin
```

```
2 crypto_usd.head(2)
```

```
3
```

	time	close	high	low	open	volumefrom	volumeto	Date	Time	volume	marketcap	price_delta
0	2023-02-19 13:00:00	24682.03	24715.82	24682.03	24707.39	903.97	22335943.28	2023-02-19	13:00:00	22335039.31	5.512964e+11	NaN
1	2023-02-19 14:00:00	24705.70	24702.95	24670.21	24692.02	1220.20	20202001.55	2023-02-19	14:00:00	20201791.26	7.180012e+11	0.27%

```
1 #Tweets-Bitcoin
```

```
2
```

```
3 tweets.head(2)
```

	user_name	user_location	user_description	user_created	user_followers	user_friends	user_favourites	user_verified	date
0	Irk	Vancouver, WA	Irk started investing in the stock market in 1...	2018-08-11 03:17:00	116.0	8.0	4580.0	False	2023-02-25 23:59:00
1	Xiang Zhang	NaN	Professional Software Engineer at Crypto	2011-01-11 01:37:00	42.0	22.0	5.0	False	2023-02-25 23:59:00



▼ Exploratory Data Analysis (EDA)

▼ Data Description

```
1 print(type(df))
```

```
<class 'pandas.core.frame.DataFrame'>
```

```
1 df.head(2)
```

	user_name	user_location	user_description	user_created	user_followers	user_friends	user_favourites	user_verified	date
0	Bitcoin Treasures δ□□□	Bitcoin City	I am a bot tracking the USD value of various b...	2020-08-23 16:20:00	2603.0	8.0	48.0	False	2023-03-04 20:00:00 world
1	Bitcoin Price	Order Book	\$BTC price updated every 4 hours\nPrices in ...	2022-03-25 08:31:00	16.0	0.0	2.0	False	2023-03-04 20:00:00 ;

2 rows × 21 columns

1

```

1 import pandas as pd
2
3 # Iterate through each attribute
4 for column in df.columns:
5     unique_count = df[column].nunique()
6     print(f"Unique count for attribute '{column}': {unique_count}")
7

```

```

Unique count for attribute 'user_name': 36053
Unique count for attribute 'user_location': 9432
Unique count for attribute 'user_description': 33094
Unique count for attribute 'user_created': 36531
Unique count for attribute 'user_followers': 11793
Unique count for attribute 'user_friends': 5190
Unique count for attribute 'user_favourites': 17427
Unique count for attribute 'user_verified': 2
Unique count for attribute 'date': 10215
Unique count for attribute 'text': 162425
Unique count for attribute 'hashtags': 40623
Unique count for attribute 'source': 719
Unique count for attribute 'is_retweet': 1
Unique count for attribute 'compound': 1716
Unique count for attribute 'score': 56006
Unique count for attribute 'sentiment_level': 5
Unique count for attribute 'polarity': 2584
Unique count for attribute 'subjectivity': 2594

```

1 df.columns

```

Index(['user_name', 'user_location', 'user_description', 'user_created',
       'user_followers', 'user_friends', 'user_favourites', 'user_verified',
       'date', 'text', 'hashtags', 'source', 'is_retweet', 'compound', 'score',
       'sentiment_level', 'polarity', 'subjectivity'],
      dtype='object')

```

1 df['sentiment_level'].value_counts()

Neutral	93169
Positive	35921
Extreme Positive	17343
Negative	15903
Extreme Negative	5316
Name: sentiment_level, dtype:	int64

```

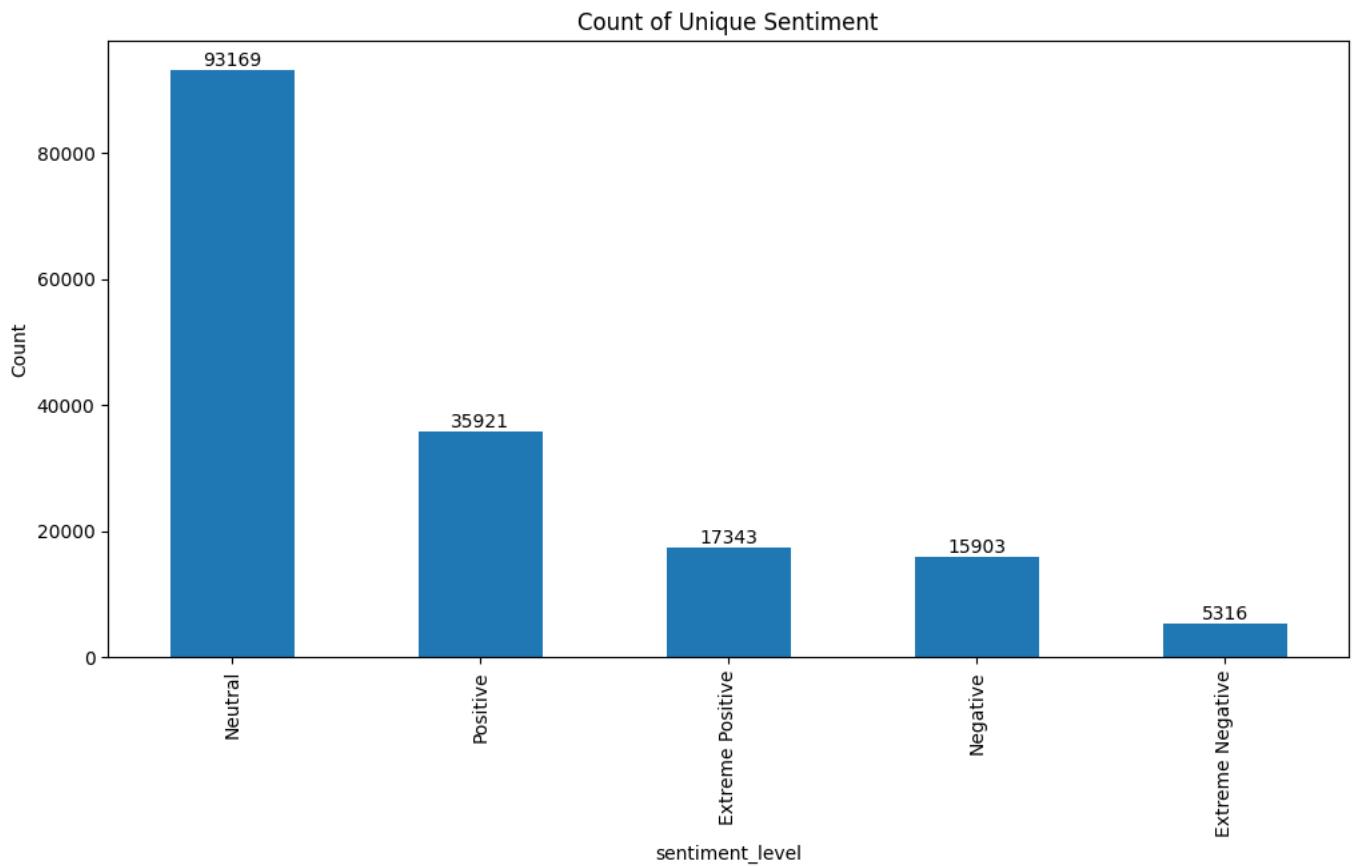
1 import matplotlib.pyplot as plt
2
3 counts = df['sentiment_level'].value_counts()
4
5 # Plotting the bar chart
6 plt.figure(figsize=(12, 6))
7 ax = counts.plot(kind='bar')
8 plt.xlabel('sentiment_level')
9 plt.ylabel('Count')
10 plt.title('Count of Unique Sentiment')
11 plt.xticks(rotation=90)
12

```

```

13 # Labeling the height on top of each bar
14 for p in ax.patches:
15     ax.annotate(str(p.get_height()), (p.get_x() + p.get_width() / 2, p.get_height()), ha='center', va='bottom')
16
17 plt.show()
18

```



```

1
2 print(df.dtypes)

user_name      object
user_location   object
user_description object
user_created    object
user_followers  float64
user_friends    float64
user_favourites float64
user_verified   bool
date           object
text            object
hashtags       object
source          object
is_retweet      float64
compound        float64
score           float64
sentiment_level object
polarity        float64
subjectivity    float64
dtype: object

```

```

1
2 # Get the shape of the merged dataset
3 rows, columns = df.shape
4 print("Number of rows:", rows)
5 print("Number of columns:", columns)

```

Number of rows: 167652
 Number of columns: 18

```

1 # Displaying the summary information
2 summary = df.info()

```

```

3 print(summary)
4

<class 'pandas.core.frame.DataFrame'>
Int64Index: 167652 entries, 0 to 27941
Data columns (total 18 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   user_name        167642 non-null   object  
 1   user_location     83349 non-null   object  
 2   user_description  157323 non-null   object  
 3   user_created      167652 non-null   object  
 4   user_followers    167652 non-null   float64 
 5   user_friends      167652 non-null   float64 
 6   user_favourites   167652 non-null   float64 
 7   user_verified     167652 non-null   bool    
 8   date              167652 non-null   object  
 9   text               167652 non-null   object  
 10  hashtags          167142 non-null   object  
 11  source             167142 non-null   object  
 12  is_retweet        167142 non-null   float64 
 13  compound           167652 non-null   float64 
 14  score              167652 non-null   float64 
 15  sentiment_level   167652 non-null   object  
 16  polarity            167652 non-null   float64 
 17  subjectivity       167652 non-null   float64 
dtypes: bool(1), float64(8), object(9)
memory usage: 23.2+ MB
None

```

```

1 # Let's see meta information about numeric data, we can also see if there any extreme values
2 df.describe()

```

	user_followers	user_friends	user_favourites	is_retweet	compound	score	polarity	subjectivity
count	1.676520e+05	167652.000000	1.676520e+05	167142.0	167652.000000	1.676520e+05	167652.000000	167652.000000
mean	1.069845e+04	774.246069	6.267079e+03	0.0	0.105530	3.260548e+07	0.069989	0.246614
std	1.317029e+05	2691.620939	2.192673e+04	0.0	0.339531	1.440558e+09	0.223690	0.271300
min	0.000000e+00	0.000000	0.000000e+00	0.0	-0.991300	-9.459832e+10	-1.000000	0.000000
25%	1.190000e+02	9.000000	1.100000e+01	0.0	0.000000	0.000000e+00	0.000000	0.000000
50%	5.560000e+02	123.000000	2.810000e+02	0.0	0.000000	0.000000e+00	0.000000	0.223016
75%	1.956000e+03	606.000000	3.300000e+03	0.0	0.361200	7.648344e+03	0.100000	0.418182
max	1.878937e+07	254276.000000	1.083014e+06	0.0	0.984000	1.350093e+11	1.000000	1.000000

```
1 df.describe().transpose()
```

	count	mean	std	min	25%	50%	75%	max
user_followers	167652.0	1.069845e+04	1.317029e+05	0.000000e+00	119.0	556.000000	1956.000000	1.878937e+07
user_friends	167652.0	7.742461e+02	2.691621e+03	0.000000e+00	9.0	123.000000	606.000000	2.542760e+05
user_favourites	167652.0	6.267079e+03	2.192673e+04	0.000000e+00	11.0	281.000000	3300.000000	1.083014e+06
is_retweet	167142.0	0.000000e+00	0.000000e+00	0.000000e+00	0.0	0.000000	0.000000	0.000000e+00
compound	167652.0	1.055299e-01	3.395308e-01	-9.913000e-01	0.0	0.000000	0.361200	9.840000e-01
score	167652.0	3.260548e+07	1.440558e+09	-9.459832e+10	0.0	0.000000	7648.344025	1.350093e+11
polarity	167652.0	6.998889e-02	2.236899e-01	-1.000000e+00	0.0	0.000000	0.100000	1.000000e+00
subjectivity	167652.0	2.466140e-01	2.712997e-01	0.000000e+00	0.0	0.223016	0.418182	1.000000e+00

▼ Analysing Null/NAN values

```
1 df.isnull().sum()
```

user_name	10
user_location	84303
user_description	10329
user_created	0
user_followers	0

```

user_friends      0
user_favourites   0
user_verified     0
date              0
text              0
hashtags         510
source            510
is_retweet        510
compound          0
score             0
sentiment_level   0
polarity          0
subjectivity      0
dtype: int64

```

```
1 #Null values for attributes which are not significant, so keepin the dataset unaltered
```

```
1 df.isnull().count()
```

```

user_name         27942
user_location     27942
user_description  27942
user_created      27942
user_followers    27942
user_friends      27942
user_favourites   27942
user_verified     27942
date              27942
text              27942
hashtags         27942
source            27942
is_retweet        27942
compound          27942
score             27942
sentiment_level   27942
polarity          27942
subjectivity      27942
dtype: int64

```

▼ Statistics Summary

```

1 # Display summary statistics of numeric columns
2 numeric_columns = df.select_dtypes(include=[int, float])
3 summary_statistics = numeric_columns.describe()
4 print("Summary Statistics:")
5 print(summary_statistics)

```

```

Summary Statistics:
   user_followers  user_friends  user_favourites  is_retweet \
count  1.676520e+05  167652.000000  1.676520e+05  167142.0
mean   1.069845e+04   774.246069   6.267079e+03   0.0
std    1.317029e+05   2691.620939   2.192673e+04   0.0
min    0.000000e+00   0.000000   0.000000e+00   0.0
25%   1.190000e+02    9.000000   1.100000e+01   0.0
50%   5.560000e+02   123.000000   2.810000e+02   0.0
75%   1.956000e+03   606.000000   3.300000e+03   0.0
max   1.878937e+07  254276.000000   1.083014e+06   0.0

   compound       score      polarity  subjectivity
count  167652.000000  1.676520e+05  167652.000000  167652.000000
mean   0.105530  3.260548e+07   0.069989   0.246614
std    0.339531  1.440558e+09   0.223690   0.271300
min   -0.991300 -9.459832e+10  -1.000000   0.000000
25%   0.000000  0.000000e+00   0.000000   0.000000
50%   0.000000  0.000000e+00   0.000000   0.223016
75%   0.361200  7.648344e+03   0.100000   0.418182
max   0.984000  1.350093e+11   1.000000   1.000000

```

▼ Categorical to One-Hot (numeric) Encoding

```

1 #Let's create a list for our categorical columns
2 #cat_cols=["Name", "Symbol", "Date"]

```

```
1 # Create a copy of the data frame in memory with a different name
2 #df_onehot=df.copy()
```

```

3 #convert only categorical variables/features to dummy/one-hot features
4 #df_onehot = pd.get_dummies(df, columns=cat_cols, prefix = cat_cols)
5 #print the dataset
6 #df_onehot

1 #Will do Unigram and Bigram for 'text'

```

▼ Pandas Profiling

```

1 #!pip install pandas-profiling
2

1 #from pandas_profiling import ProfileReport
2 #profile = ProfileReport(df)

1 #profile.to_notebook_iframe()

```

▼ Univariate Analysis

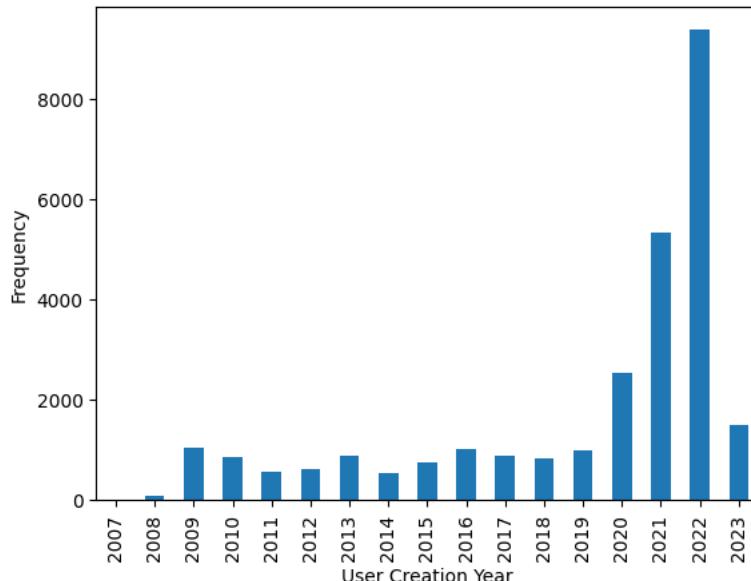
```

1 import matplotlib.pyplot as plt
2 # Univariate analysis for user_created
3 df['user_created'] = pd.to_datetime(df['user_created'])
4 df['user_created_year'] = df['user_created'].dt.year
5 user_created_counts = df['user_created_year'].value_counts().sort_index()
6 user_created_counts.plot(kind='bar')
7 plt.xlabel('User Creation Year')
8 plt.ylabel('Frequency')
9 plt.title('Distribution of User Creation Years')
10 plt.show()
11
12 # Univariate analysis for user_followers
13 user_followers_stats = df['user_followers'].describe()
14 print(user_followers_stats)
15
16 plt.hist(df['user_followers'], bins=20)
17 plt.xlabel('User Followers')
18 plt.ylabel('Frequency')
19 plt.title('Distribution of User Followers')
20 plt.show()
21
22 # Univariate analysis for user_friends
23 user_friends_stats = df['user_friends'].describe()
24 print(user_friends_stats)
25
26 plt.hist(df['user_friends'], bins=20)
27 plt.xlabel('User Friends')
28 plt.ylabel('Frequency')
29 plt.title('Distribution of User Friends')
30 plt.show()
31
32 # Univariate analysis for user_favourites
33 user_favourites_stats = df['user_favourites'].describe()
34 print(user_favourites_stats)
35
36 plt.hist(df['user_favourites'], bins=20)
37 plt.xlabel('User Favorites')
38 plt.ylabel('Frequency')
39 plt.title('Distribution of User Favorites')
40 plt.show()
41
42 # Univariate analysis for user_verified
43 user_verified_counts = df['user_verified'].value_counts()
44 user_verified_counts.plot(kind='bar')
45 plt.xlabel('User Verified')
46 plt.ylabel('Frequency')
47 plt.title('Distribution of User Verification')
48 plt.show()
49
50 # Univariate analysis for date

```

```
51 df['date'] = pd.to_datetime(df['date'])
52 df['date_year'] = df['date'].dt.year
53 date_counts = df['date_year'].value_counts().sort_index()
54 date_counts.plot(kind='bar')
55 plt.xlabel('Year')
56 plt.ylabel('Frequency')
57 plt.title('Distribution of Tweets over Time')
58 plt.show()
59
60 # Univariate analysis for text
61 df['text_length'] = df['text'].apply(lambda x: len(x))
62 text_length_stats = df['text_length'].describe()
63 print(text_length_stats)
64
65 plt.hist(df['text_length'], bins=20)
66 plt.xlabel('Text Length')
67 plt.ylabel('Frequency')
68 plt.title('Distribution of Text Lengths')
69 plt.show()
70
71 # Univariate analysis for hashtags - Will Perform word cloud analysis
```

Distribution of User Creation Years

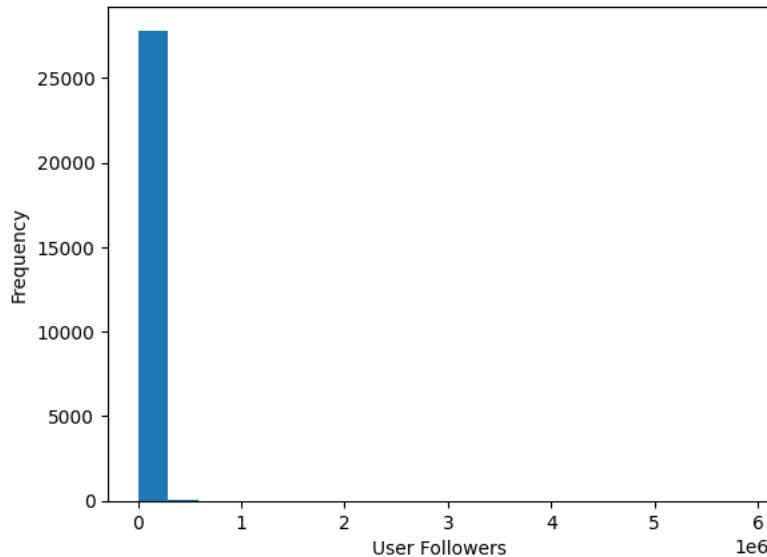


```

count    2.794200e+04
mean     8.843726e+03
std      9.777185e+04
min      0.000000e+00
25%     1.130000e+02
50%     5.430000e+02
75%     1.852000e+03
max     5.825939e+06
Name: user_followers, dtype: float64

```

Distribution of User Followers

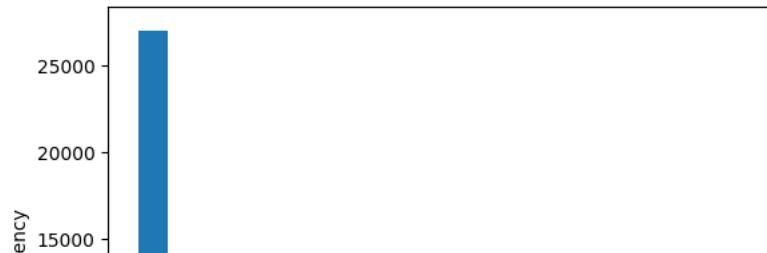


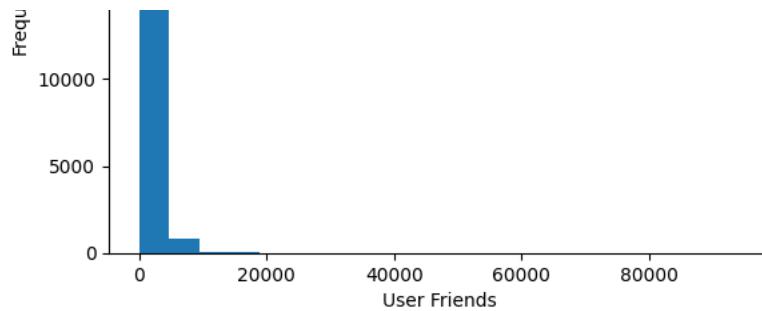
```

count    27942.000000
mean     722.446353
std      2094.308428
min      0.000000
25%     9.000000
50%    121.000000
75%    606.000000
max    94353.000000
Name: user_friends, dtype: float64

```

Distribution of User Friends



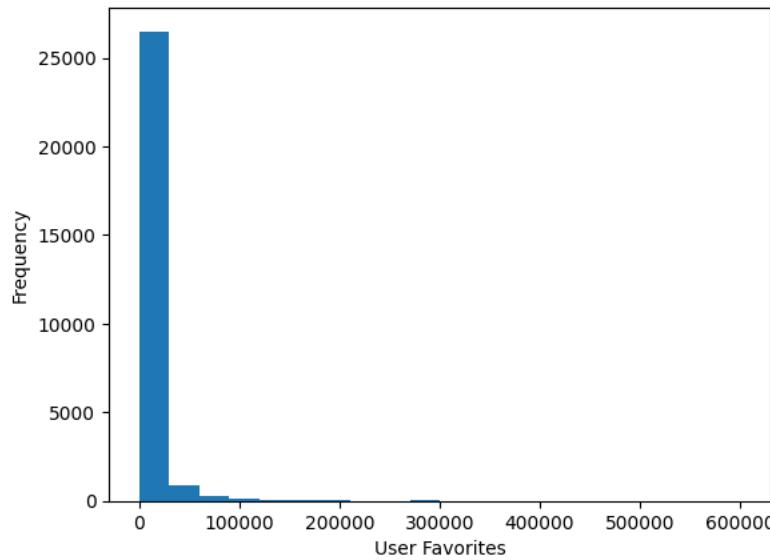


```

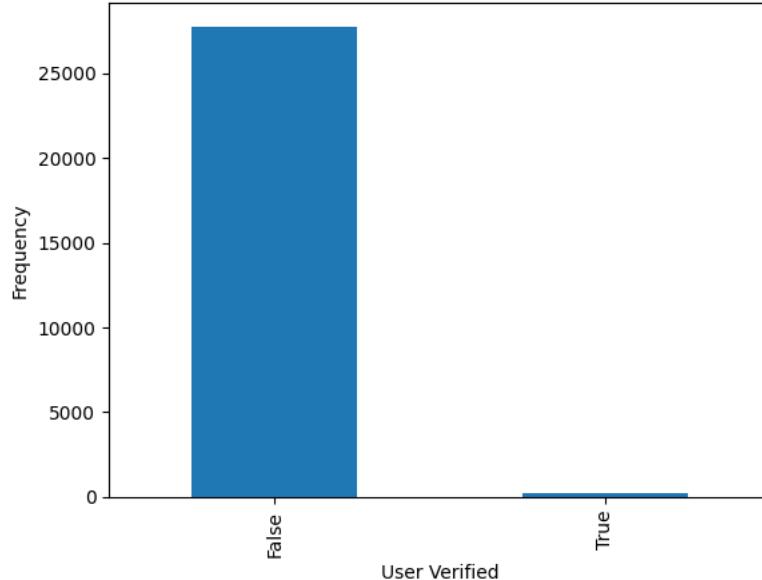
count      27942.000000
mean       6427.776358
std        21926.397398
min        0.000000
25%       11.000000
50%       281.000000
75%      3561.750000
max      600534.000000
Name: user_favourites, dtype: float64

```

Distribution of User Favorites



Distribution of User Verification



Distribution of Tweets over Time



```

1 # Univariate analysis for source
2 #source_counts = df['source'].value_counts()

```

```
3 #source_counts.plot(kind='bar')
4 #plt.xlabel('Tweet Source')
5 #plt.ylabel('Frequency')
6 #plt.title('Distribution of Tweet Sources')
7 #plt.show()

1
2 # Univariate analysis for is_retweet
3 is_retweet_counts = df['is_retweet'].value_counts()
4 is_retweet_counts.plot(kind='bar')
5 plt.xlabel('Is Retweet')
6 plt.ylabel('Frequency')
7 plt.title('Distribution of Retweets')
8 plt.show()
9

10 # Univariate analysis for compound
11 compound_stats = df['compound'].describe()
12 print(compound_stats)
13

14 plt.hist(df['compound'], bins=20)
15 plt.xlabel('Compound Sentiment Score')
16 plt.ylabel('Frequency')
17 plt.title('Distribution of Compound Sentiment Scores')
18 plt.show()
19

20 # Univariate analysis for score
21 score_stats = df['score'].describe()
22 print(score_stats)
23

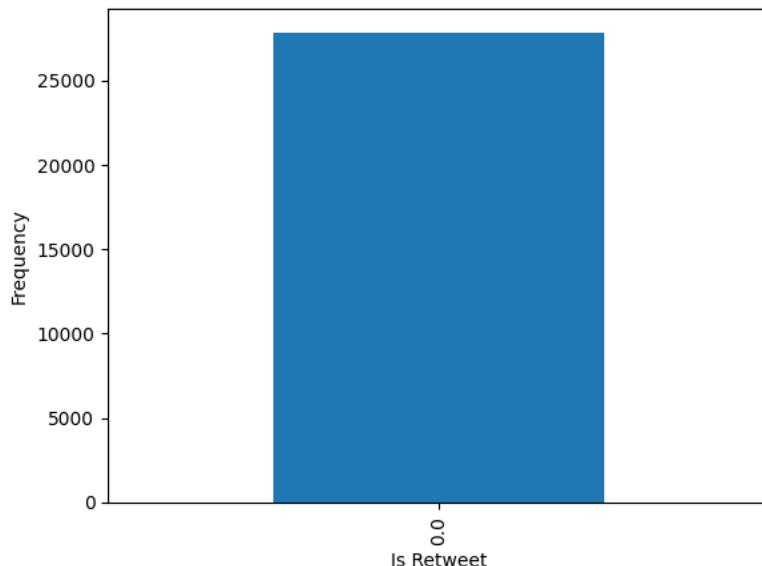
24 plt.hist(df['score'], bins=20)
25 plt.xlabel('Sentiment Score')
26 plt.ylabel('Frequency')
27 plt.title('Distribution of Sentiment Scores')
28 plt.show()
29

30 # Univariate analysis for sentiment_level
31 sentiment_level_counts = df['sentiment_level'].value_counts()
32 sentiment_level_counts.plot(kind='bar')
33 plt.xlabel('Sentiment Level')
34 plt.ylabel('Frequency')
35 plt.title('Distribution of Sentiment Levels')
36 plt.show()
37

38 # Univariate analysis for polarity
39 polarity_stats = df['polarity'].describe()
40 print(polarity_stats)
41

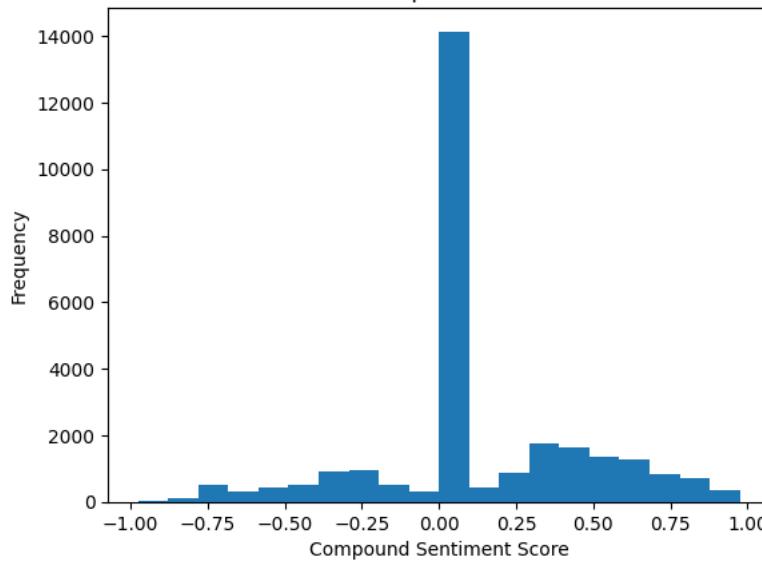
42 plt.hist(df['polarity'], bins=20)
43 plt.xlabel('Polarity')
44 plt.ylabel('Frequency')
45 plt.title('Distribution of Polarity')
46 plt.show()
47
```

Distribution of Retweets



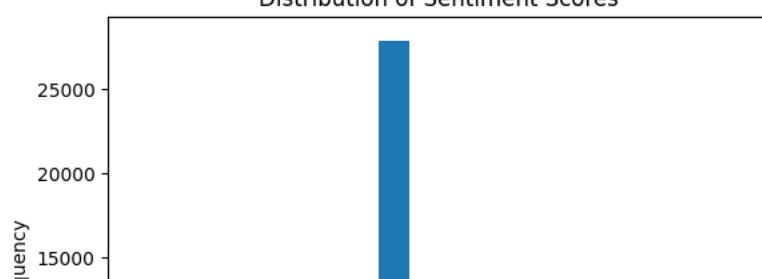
```
count    27942.000000
mean     0.102366
std      0.343487
min     -0.976100
25%     0.000000
50%     0.000000
75%     0.340000
max     0.975300
Name: compound, dtype: float64
```

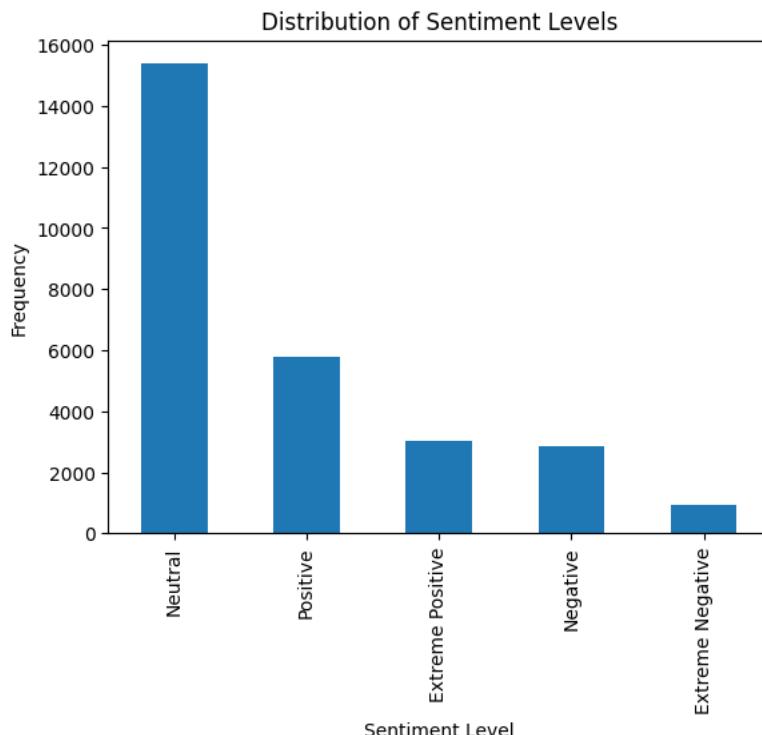
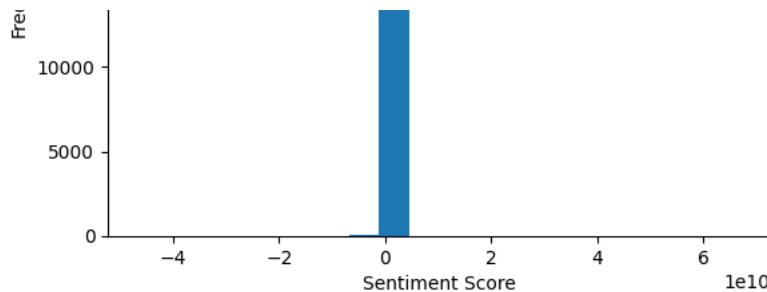
Distribution of Compound Sentiment Scores



```
count    2.794200e+04
mean    2.813255e+07
std     1.346742e+09
min    -4.657105e+10
25%     0.000000e+00
50%     0.000000e+00
75%     8.164702e+03
max     6.683636e+10
Name: score, dtype: float64
```

Distribution of Sentiment Scores





```

count    27942.000000
mean      0.069336
std       0.221637
min     -1.000000
25%      0.000000
50%      0.000000
75%      0.095238
max      1.000000
Name: polarity, dtype: float64

```

Distribution of Polarity



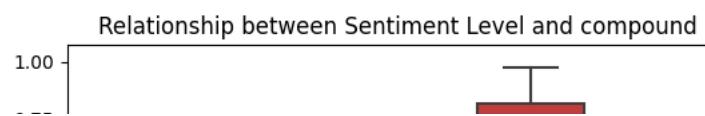
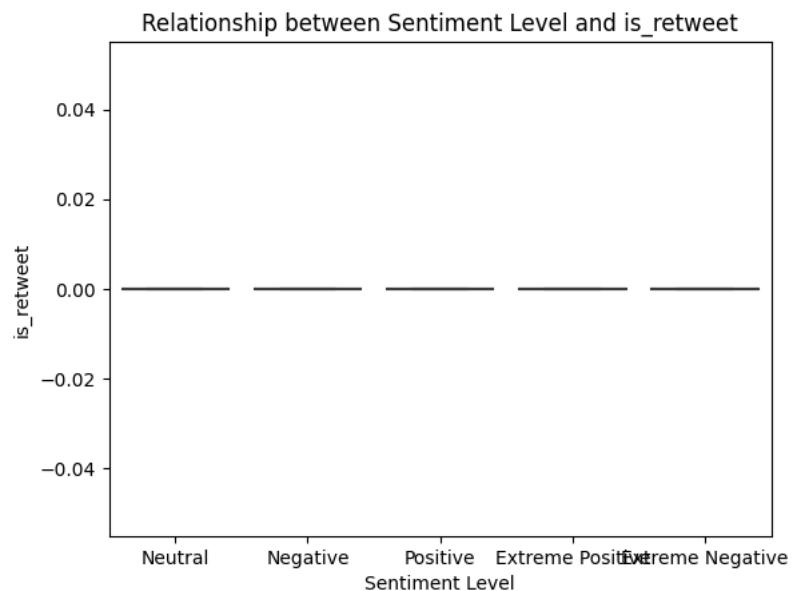
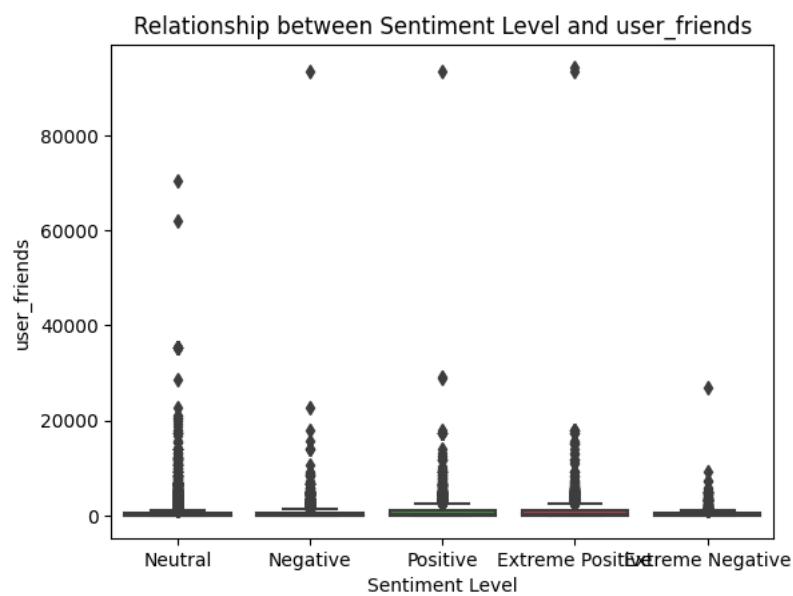
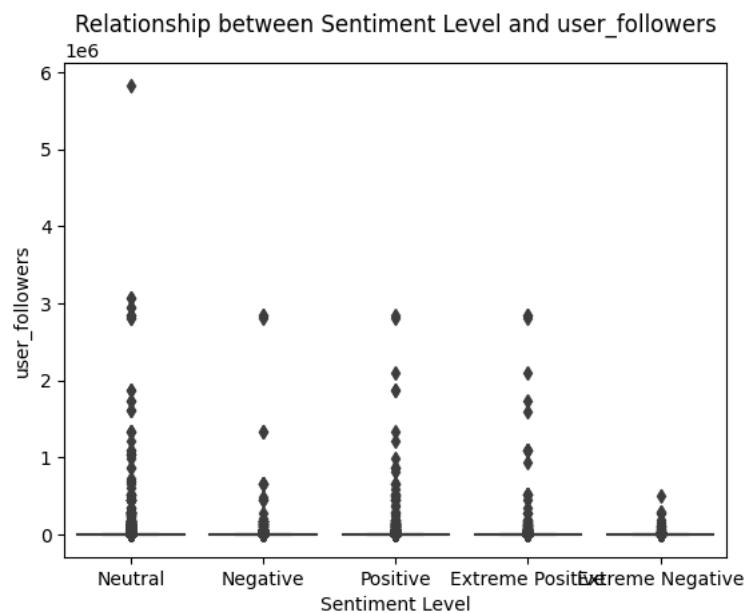
▼ Bivariate Analysis

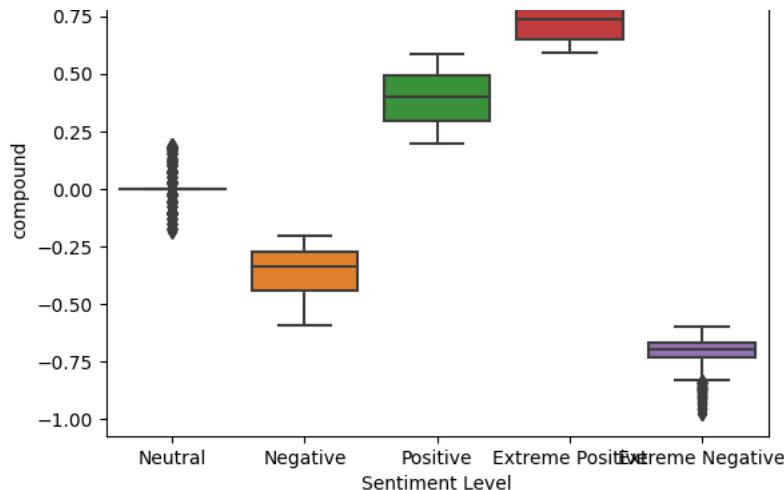
```

1 #Bivariate Analysis of sentiment_level with other attributes
2
3 # Select pairs of attributes for bivariate analysis
4 attribute_pairs = [
5     ('sentiment_level', 'user_followers'),
6     ('sentiment_level', 'user_friends'),
7     ('sentiment_level', 'is_retweet'),
8     ('sentiment_level', 'compound'),
9     ('sentiment_level', 'score'),
10    ('sentiment_level', 'polarity'),
11    ('sentiment_level', 'subjectivity'),
12
13 ]
14
15 # Perform bivariate analysis for each attribute pair
16 for attr_x, attr_y in attribute_pairs:
17     df['sentiment_level'] = df['sentiment_level'].astype(str) # Convert sentiment_level to string type

```

```
18     sns.boxplot(x=attr_x, y=attr_y, data=df)
19     plt.xlabel('Sentiment Level')
20     plt.ylabel(attr_y)
21     plt.title(f'Relationship between Sentiment Level and {attr_y}')
22     plt.show()
23
```





1e10 Relationship between Sentiment Level and score

▼ Data visualization

```

1 !pip install seaborn
2 !pip install matplotlib
3

Looking in indexes: https://pypi.org/simple, https://us-python.pkg.dev/colab-wheels/public/simple/
Requirement already satisfied: seaborn in /usr/local/lib/python3.10/dist-packages (0.12.2)
Requirement already satisfied: numpy!=1.24.0,>=1.17 in /usr/local/lib/python3.10/dist-packages (from seaborn) (1.22.4)
Requirement already satisfied: pandas>=0.25 in /usr/local/lib/python3.10/dist-packages (from seaborn) (1.5.3)
Requirement already satisfied: matplotlib!=3.6.1,>=3.1 in /usr/local/lib/python3.10/dist-packages (from seaborn) (3.7.1)
Requirement already satisfied: contourpy>=1.0.1 in /usr/local/lib/python3.10/dist-packages (from matplotlib!=3.6.1,>=3.1->seaborn) (1.1)
Requirement already satisfied: cycler>=0.10 in /usr/local/lib/python3.10/dist-packages (from matplotlib!=3.6.1,>=3.1->seaborn) (0.11.0)
Requirement already satisfied: fonttools>=4.22.0 in /usr/local/lib/python3.10/dist-packages (from matplotlib!=3.6.1,>=3.1->seaborn) (4.)
Requirement already satisfied: kiwisolver>=1.0.1 in /usr/local/lib/python3.10/dist-packages (from matplotlib!=3.6.1,>=3.1->seaborn) (1.)
Requirement already satisfied: packaging>=20.0 in /usr/local/lib/python3.10/dist-packages (from matplotlib!=3.6.1,>=3.1->seaborn) (23.1)
Requirement already satisfied: pillow>=6.2.0 in /usr/local/lib/python3.10/dist-packages (from matplotlib!=3.6.1,>=3.1->seaborn) (8.4.0)
Requirement already satisfied: pyparsing>=2.3.1 in /usr/local/lib/python3.10/dist-packages (from matplotlib!=3.6.1,>=3.1->seaborn) (3.1)
Requirement already satisfied: python-dateutil>=2.7 in /usr/local/lib/python3.10/dist-packages (from matplotlib!=3.6.1,>=3.1->seaborn)
Requirement already satisfied: pytz>=2020.1 in /usr/local/lib/python3.10/dist-packages (from pandas>=0.25->seaborn) (2022.7.1)
Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.10/dist-packages (from python-dateutil>=2.7->matplotlib!=3.6.1,>=3.1-
Looking in indexes: https://pypi.org/simple, https://us-python.pkg.dev/colab-wheels/public/simple/
Requirement already satisfied: matplotlib in /usr/local/lib/python3.10/dist-packages (3.7.1)
Requirement already satisfied: contourpy>=1.0.1 in /usr/local/lib/python3.10/dist-packages (from matplotlib) (1.1.0)
Requirement already satisfied: cycler>=0.10 in /usr/local/lib/python3.10/dist-packages (from matplotlib) (0.11.0)
Requirement already satisfied: fonttools>=4.22.0 in /usr/local/lib/python3.10/dist-packages (from matplotlib) (4.40.0)
Requirement already satisfied: kiwisolver>=1.0.1 in /usr/local/lib/python3.10/dist-packages (from matplotlib) (1.4.4)
Requirement already satisfied: numpy>=1.20 in /usr/local/lib/python3.10/dist-packages (from matplotlib) (1.22.4)
Requirement already satisfied: packaging>=20.0 in /usr/local/lib/python3.10/dist-packages (from matplotlib) (23.1)
Requirement already satisfied: pillow>=6.2.0 in /usr/local/lib/python3.10/dist-packages (from matplotlib) (8.4.0)
Requirement already satisfied: pyparsing>=2.3.1 in /usr/local/lib/python3.10/dist-packages (from matplotlib) (3.1.0)
Requirement already satisfied: python-dateutil>=2.7 in /usr/local/lib/python3.10/dist-packages (from matplotlib) (2.8.2)
Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.10/dist-packages (from python-dateutil>=2.7->matplotlib) (1.16.0)

```

```
1 #Some Visualizations of attributes for analysis
```

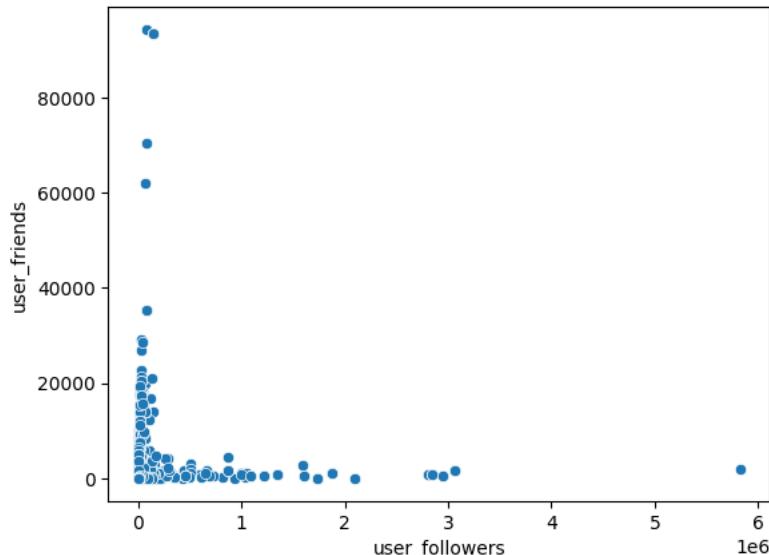
```

1 import seaborn as sns
2 import matplotlib.pyplot as plt
3
4 # Scatter plot
5 sns.scatterplot(data=df, x='user_followers', y='user_friends')
6 plt.title('Scatter Plot: User Followers vs User Friends')
7 plt.show()
8
9 # Box plot
10 sns.boxplot(data=df, x='user_verified', y='user_followers')
11 plt.title('Box Plot: User Verified vs User Followers')
12 plt.show()
13
14 # Bar chart
15 sns.countplot(data=df, x='sentiment_level')
16 plt.title('Bar Chart: Sentiment Levels')

```

```
17 plt.show()
18
19 # Histogram
20 sns.histplot(data=df, x='user_favourites', bins=20)
21 plt.title('Histogram: User Favorites')
22 plt.show()
23
24 # Heatmap
25 heatmap_data = df[['user_followers', 'user_friends', 'user_favourites', 'compound', 'polarity', 'subjectivity']].corr()
26 sns.heatmap(heatmap_data, annot=True, cmap='coolwarm')
27 plt.title('Heatmap: Correlation Matrix')
28 plt.show()
29
```

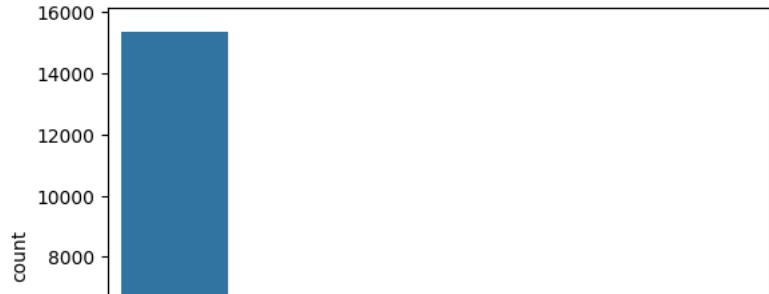
Scatter Plot: User Followers vs User Friends



Box Plot: User Verified vs User Followers



Bar Chart: Sentiment Levels

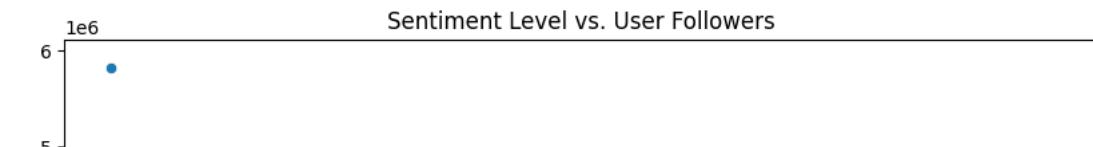
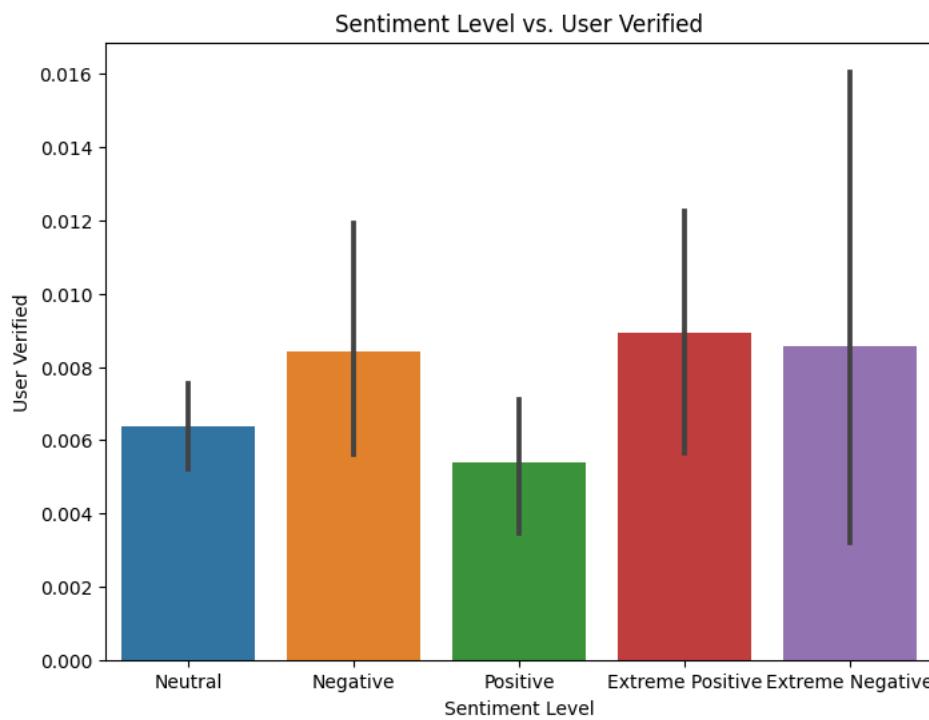
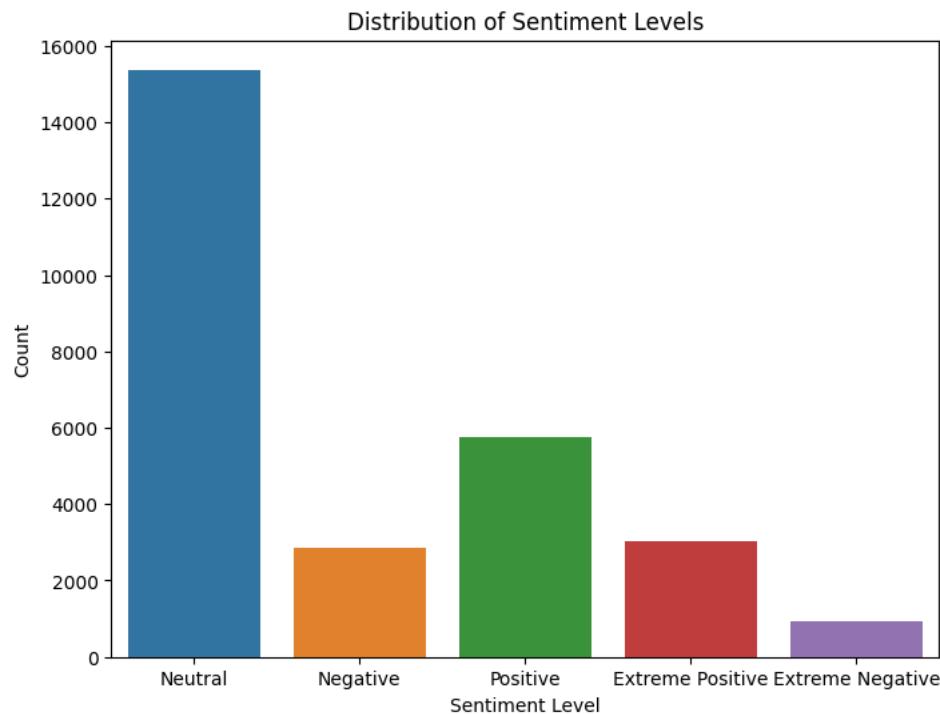


```

1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4
5 # Convert 'user_verified' column to numeric representation
6 df['user_verified_numeric'] = df['user_verified'].astype(int)
7
8 # EDA for sentiment_level
9
10 # Countplot of sentiment_level
11 plt.figure(figsize=(8, 6))
12 sns.countplot(x='sentiment_level', data=df)
13 plt.title('Distribution of Sentiment Levels')
14 plt.xlabel('Sentiment Level')
15 plt.ylabel('Count')
16 plt.show()

```

```
10 plt.show()
11
12 # Relationship with other attributes
13
14 # Bar plot of sentiment_level vs. user_verified
15 plt.figure(figsize=(8, 6))
16 sns.barplot(x='sentiment_level', y='user_verified_numeric', data=df)
17 plt.title('Sentiment Level vs. User Verified')
18 plt.xlabel('Sentiment Level')
19 plt.ylabel('User Verified')
20 plt.show()
21
22 # Bar plot of sentiment_level vs. user_location
23 plt.figure(figsize=(12, 6))
24 sns.countplot(x='sentiment_level', hue='user_location', data=df)
25 plt.title('Sentiment Level vs. User Location')
26 plt.xlabel('Sentiment Level')
27 plt.ylabel('Count')
28 plt.legend(title='User Location')
29 plt.show()
30
31 # Scatter plot of sentiment_level vs. user_followers
32 plt.figure(figsize=(10, 6))
33 sns.scatterplot(x='sentiment_level', y='user_followers', data=df)
34 plt.title('Sentiment Level vs. User Followers')
35 plt.xlabel('Sentiment Level')
36 plt.ylabel('User Followers')
37 plt.show()
38
39 # Summary statistics of sentiment_level
40 sentiment_stats = df.groupby('sentiment_level').describe()['user_followers']
41 print("\nSummary Statistics of User Followers by Sentiment Level:")
42 print(sentiment_stats)
43
44 # Remove the temporary column
45 df.drop('user_verified_numeric', axis=1, inplace=True)
46
47
```



▼ Feature engineering

```
4 +
```

▼ Scaling and Encoding

```
5 # Creating interaction variables
```

```
5 df['interaction_var'] = df['user_followers'] * df['user_friends']
```

```

6
7 # Scaling or normalizing variables
8 scaler = MinMaxScaler()
9 df['user_followers_scaled'] = scaler.fit_transform(df[['user_followers']])
10 df['user_friends_scaled'] = scaler.fit_transform(df[['user_friends']])
11
12 # Binning variables
13 df['user_favourites_bins'] = pd.cut(df['user_favourites'], bins=[0, 100, 1000, float('inf')], labels=['low', 'medium', 'high'])
14
15 # Encoding categorical variables
16 label_encoder = LabelEncoder()
17 df['sentiment_level_encoded'] = label_encoder.fit_transform(df['sentiment_level'])
18
19 # Display the updated DataFrame
20 print(df.head())
21

          user_name user_location \
0  Bitcoin Treasures 🇺🇸 Bitcoin City
1      Bitcoin Price     Order Book
2        OPQRSTEE       Miami
3      BNB Price Tracker      NaN
4        Frases_k3_      NaN

          user_description      user_created \
0  I am a bot tracking the USD value of various b... 2020-08-23 16:20:00
1  $BTC price updated every 4 hours\n\nPrices in ... 2022-03-25 08:31:00
2                           NaN 2019-01-22 07:36:00
3           #BNB #BNBTracker #BinanceCoin #Binance 2022-07-11 05:24:00
4                           NaN 2011-02-25 15:56:00

    user_followers  user_friends  user_favourites  user_verified \
0        2603.0        8.0        48.0      False
1         16.0        0.0        2.0      False
2         47.0       169.0        65.0      False
3        483.0        5.0        0.0      False
4        267.0       754.0       727.0      False

          date                  text \
0  2023-03-04 20:00:00  cynthia lummi bitcoin treasur worth36279 usd 5...
1  2023-03-04 20:00:00  btc price sat 04 mar 2023 200042 gmt 222579 🇺...
2  2023-03-04 20:00:00  check product 🇺🇸 hulk logo men premium tshir...
3  2023-03-04 20:00:00  binancecoin price updat bnb bnb 28882 usd bitc...
4  2023-03-04 20:00:00  know peopl say bitcoin

    ... compound    score sentiment_level  polarity subjectivity \
0 ...  0.0000  0.0000      Neutral     0.00      0.00
1 ... -0.2732 -13.9332     Negative    0.08      0.42
2 ...  0.0000  0.0000      Neutral   -0.40      0.40
3 ...  0.0000  0.0000      Neutral     0.00      0.25
4 ...  0.0000  0.0000      Neutral     0.00      0.00

  interaction_var  user_followers_scaled  user_friends_scaled \
0        20824.0            0.000447        0.000085
1          0.0            0.000003        0.000000
2        7943.0            0.000008        0.001791
3        2415.0            0.000083        0.000053
4       201318.0            0.000046        0.007991

  user_favourites_bins  sentiment_level_encoded
0             low                      3
1             low                      2
2             low                      3
3             NaN                      3
4        medium                      3

[5 rows x 23 columns]

```

▼ Correlation analysis

```
1 #Correlation Matrix for all attributes
```

```

1 import seaborn as sns
2
3 # Generate correlation matrix
4 correlation_matrix = df.corr()
5
6 # Create heatmap of correlation matrix

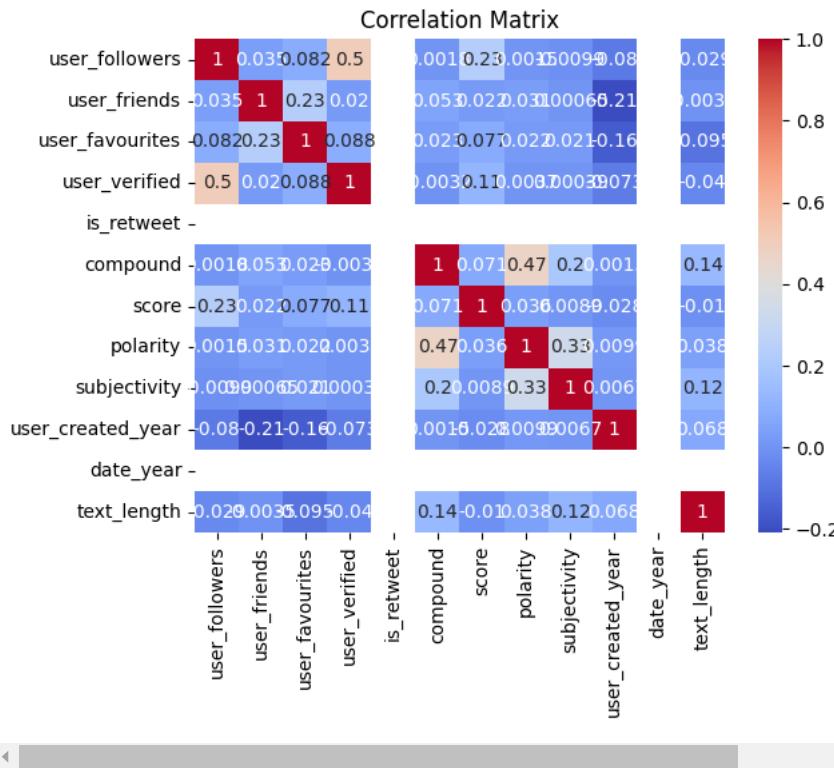
```

```

7 sns.heatmap(correlation_matrix, annot=True, cmap="coolwarm")
8
9 # Display the plot
10 plt.title("Correlation Matrix")
11 plt.show()
12

<ipython-input-22-dae9ae9f4ad1>:4: FutureWarning: The default value of numeric_only in DataFrame.corr is deprecated. In a future version
correlation_matrix = df.corr()

```



```
1 #Correlation Matrix for Selected features
```

```

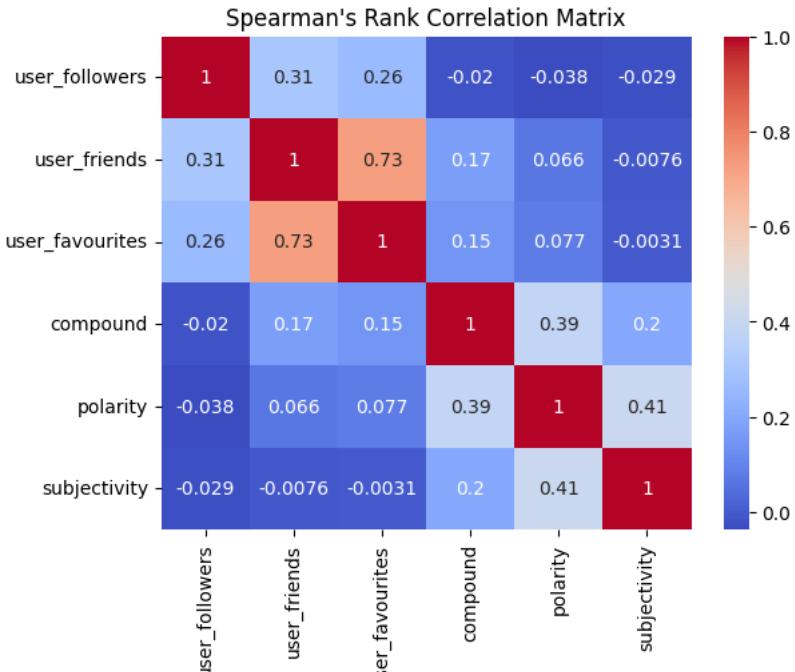
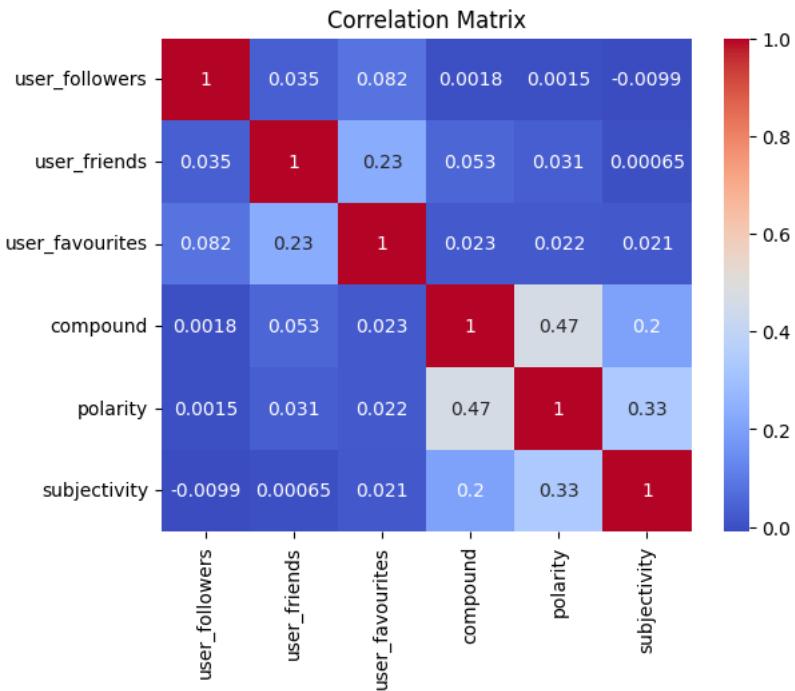
1 import pandas as pd
2 import seaborn as sns
3 import matplotlib.pyplot as plt
4
5 # Select the relevant columns for correlation analysis
6 selected_columns = ['user_followers', 'user_friends', 'user_favourites', 'compound', 'polarity', 'subjectivity']
7 correlation_data = df[selected_columns]
8
9 # Calculate correlation matrix
10 correlation_matrix = correlation_data.corr()
11
12 # Visualize correlation matrix as a heatmap
13 sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm')
14 plt.title('Correlation Matrix')
15 plt.show()
16
17 # Calculate Spearman's rank correlation matrix
18 spearman_corr_matrix = correlation_data.corr(method='spearman')
19
20 # Visualize Spearman's rank correlation matrix as a heatmap
21 sns.heatmap(spearman_corr_matrix, annot=True, cmap='coolwarm')
22 plt.title('Spearman's Rank Correlation Matrix')
23 plt.show()
24
25 # Identify highly correlated variables
26 threshold = 0.7 # Set the correlation threshold as per your requirement
27 highly_correlated_pairs = set()
28 for i in range(len(correlation_matrix.columns)):
29     for j in range(i):
30         if abs(correlation_matrix.iloc[i, j]) >= threshold:
31             col_i = correlation_matrix.columns[i]
32             col_j = correlation_matrix.columns[j]
33             highly_correlated_pairs.add((col_i, col_j))
34

```

```

35 print("Highly Correlated Variable Pairs:")
36 for pair in highly_correlated_pairs:
37     print(pair)
38

```



Highly Correlated Variable Pairs:

▼ Outlier and anomaly detection

1 #Outliers in categorical or non-numeric columns might not be meaningful or applicable in the same way as with numerical data

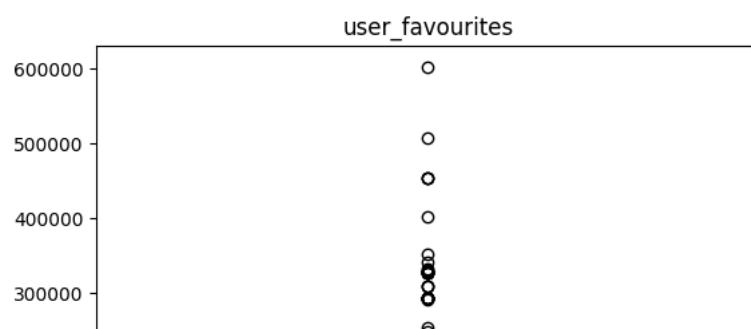
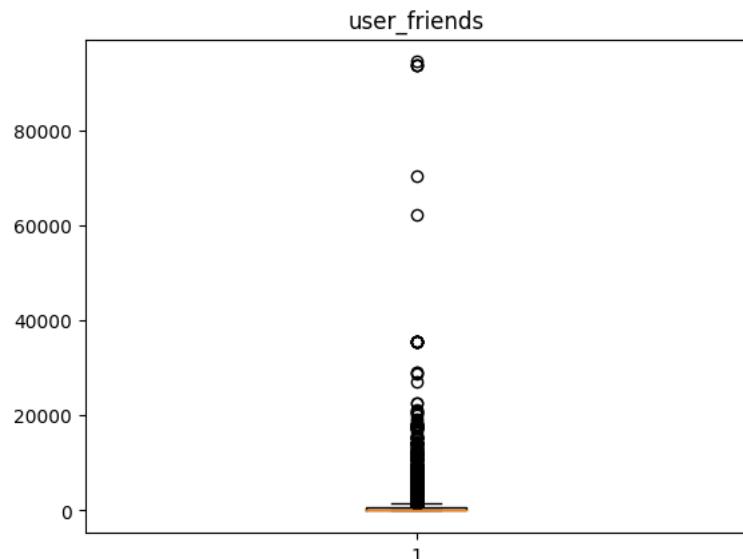
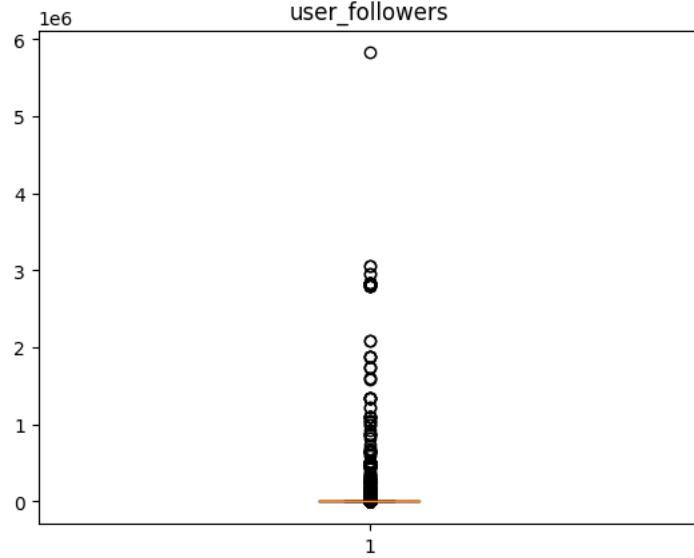
```

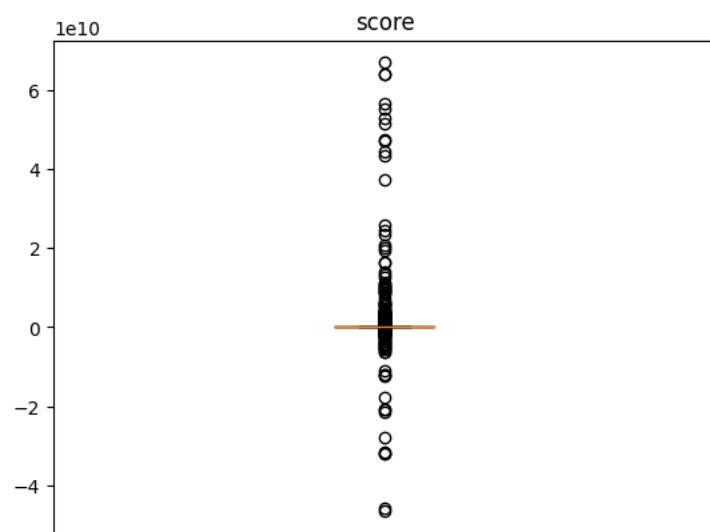
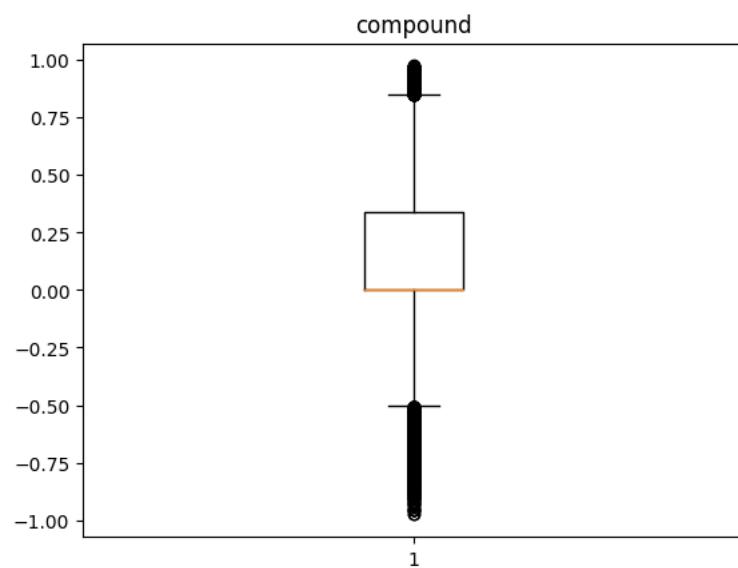
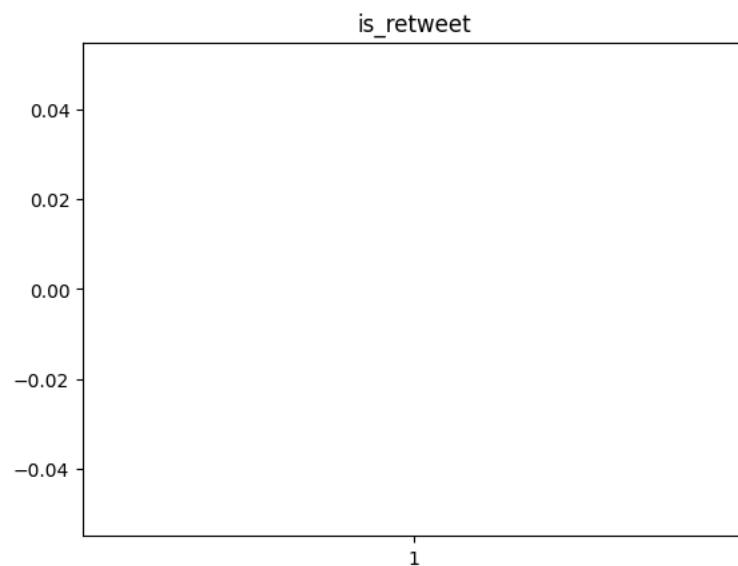
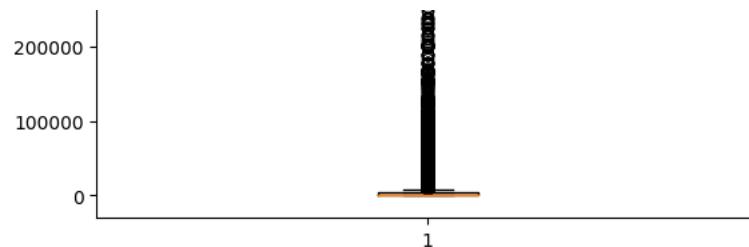
1 import pandas as pd
2 import matplotlib.pyplot as plt
3
4 # Select numerical columns for outlier analysis
5 numerical_columns = df.select_dtypes(include='number').columns
6
7 # Compute descriptive statistics for numerical columns

```

```
8 df_describe = df[numerical_columns].describe()
9
10 # Identify and store outliers
11 outliers = pd.DataFrame(columns=numerical_columns)
12 for column in numerical_columns:
13     q1 = df_describe.loc['25%', column]
14     q3 = df_describe.loc['75%', column]
15     iqr = q3 - q1
16     threshold = 1.5 * iqr
17     column_outliers = df[(df[column] < q1 - threshold) | (df[column] > q3 + threshold)]
18     outliers = pd.concat([outliers, column_outliers])
19
20 # Visualize outliers
21 for column in numerical_columns:
22     plt.figure()
23     plt.boxplot(df[column])
24     plt.title(column)
25
26 plt.show()
27
```

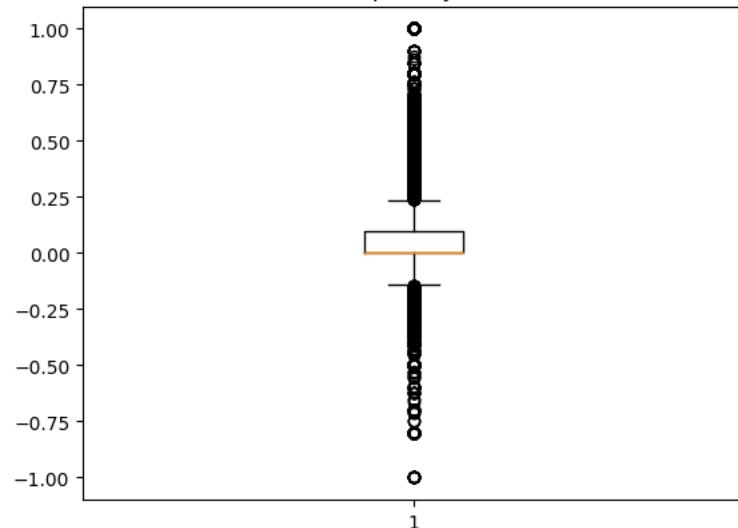
```
<ipython-input-20-245c3cb061a7>:18: FutureWarning: In a future version, object-dtype columns with all-bool values will not be included
outliers = pd.concat([outliers, column_outliers])
<ipython-input-20-245c3cb061a7>:18: FutureWarning: In a future version, object-dtype columns with all-bool values will not be included
outliers = pd.concat([outliers, column_outliers])
<ipython-input-20-245c3cb061a7>:18: FutureWarning: In a future version, object-dtype columns with all-bool values will not be included
outliers = pd.concat([outliers, column_outliers])
<ipython-input-20-245c3cb061a7>:18: FutureWarning: In a future version, object-dtype columns with all-bool values will not be included
outliers = pd.concat([outliers, column_outliers])
<ipython-input-20-245c3cb061a7>:18: FutureWarning: In a future version, object-dtype columns with all-bool values will not be included
outliers = pd.concat([outliers, column_outliers])
<ipython-input-20-245c3cb061a7>:18: FutureWarning: In a future version, object-dtype columns with all-bool values will not be included
outliers = pd.concat([outliers, column_outliers])
<ipython-input-20-245c3cb061a7>:18: FutureWarning: In a future version, object-dtype columns with all-bool values will not be included
outliers = pd.concat([outliers, column_outliers])
<ipython-input-20-245c3cb061a7>:18: FutureWarning: In a future version, object-dtype columns with all-bool values will not be included
outliers = pd.concat([outliers, column_outliers])
<ipython-input-20-245c3cb061a7>:18: FutureWarning: In a future version, object-dtype columns with all-bool values will not be included
outliers = pd.concat([outliers, column_outliers])
<ipython-input-20-245c3cb061a7>:18: FutureWarning: In a future version, object-dtype columns with all-bool values will not be included
outliers = pd.concat([outliers, column_outliers])
```



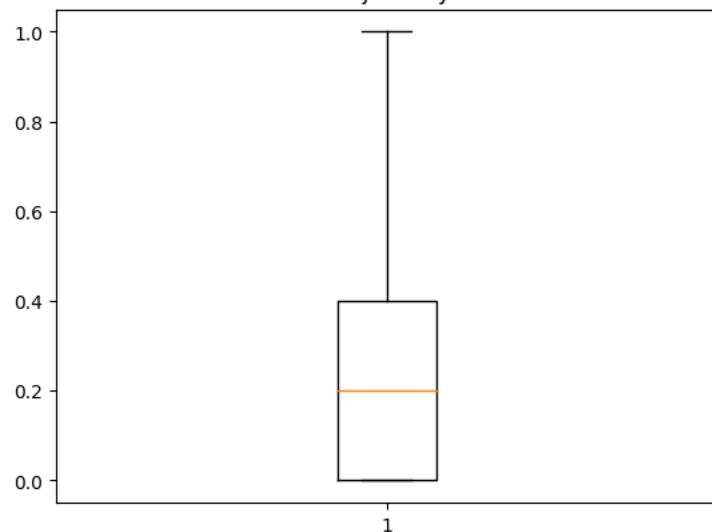


1

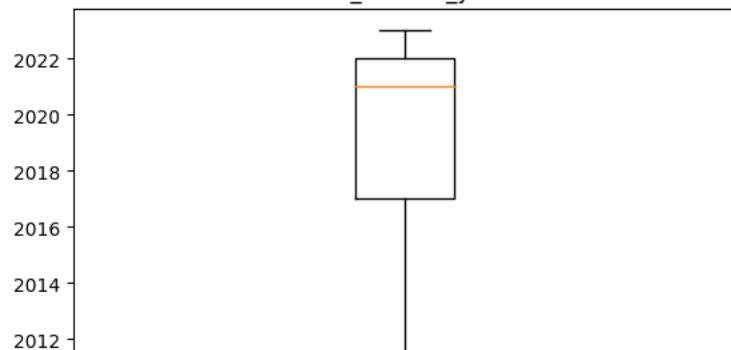
polarity



subjectivity



user_created_year



1 #Outlier for User Followers

---- | - |

```

1 import pandas as pd
2 import numpy as np
3
4 # Detect outliers using z-score
5 def detect_outliers_zscore(data, threshold=3):
6     z_scores = (data - data.mean()) / data.std()
7     return np.abs(z_scores) > threshold
8
9 # Detect outliers using IQR (Interquartile Range)
10 def detect_outliers_iqr(data, k=1.5):
11     Q1 = data.quantile(0.25)

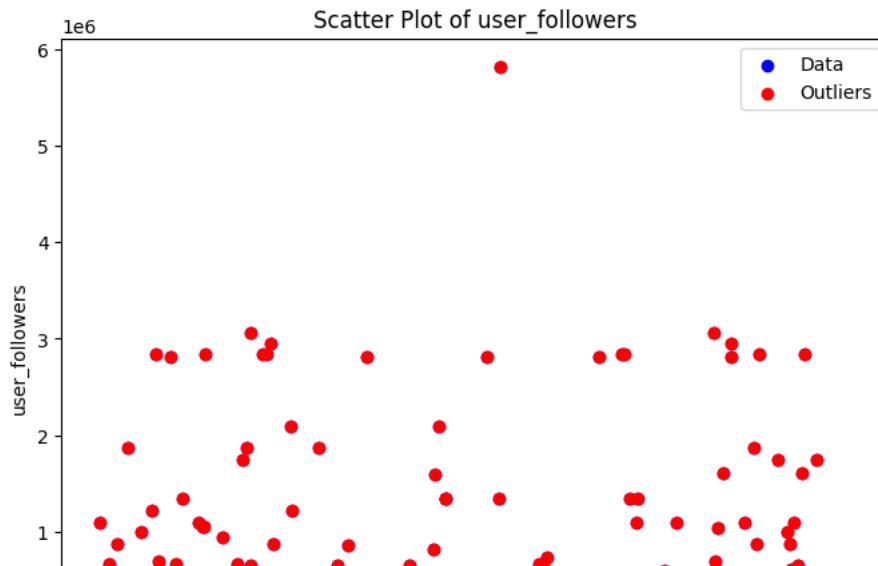
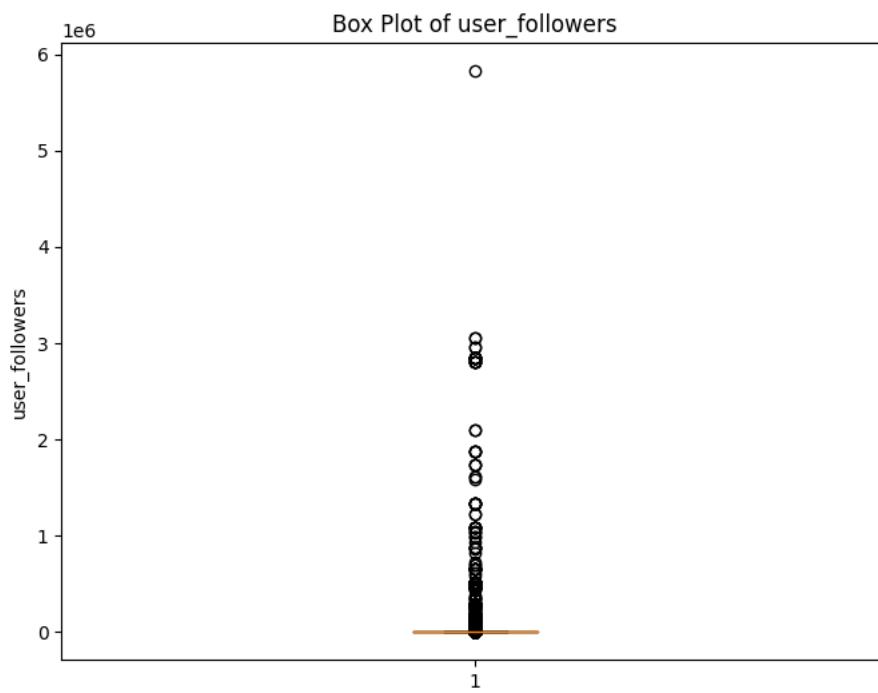
```

```

12     Q3 = data.quantile(0.75)
13     IQR = Q3 - Q1
14     lower_bound = Q1 - k * IQR
15     upper_bound = Q3 + k * IQR
16     return (data < lower_bound) | (data > upper_bound)
17
18 # Apply outlier detection on a specific column
19 outliers_zscore = detect_outliers_zscore(df['user_followers'])
20 outliers_iqr = detect_outliers_iqr(df['user_followers'])
21
22 # Display the outliers
23 outliers_zscore_indices = df[outliers_zscore].index
24 outliers_iqr_indices = df[outliers_iqr].index
25
26 print("Outliers detected using Z-score:")
27 print(outliers_zscore_indices)
28
29 print("Outliers detected using IQR:")
30 print(outliers_iqr_indices)
31

    Outliers detected using Z-score:
    Int64Index([ 22,    364,    666,   1010,   1079,   1379,   1523,   1542,   1671,
                 1935,
                 ...
                 25448, 25536, 25670, 25723, 25814, 25962, 26053, 26126, 26190,
                 26641],
                dtype='int64', length=116)
    Outliers detected using IQR:
    Int64Index([ 17,    22,    35,    40,    41,    42,    46,    51,    66,
                 77,
                 ...
                 27881, 27886, 27902, 27903, 27913, 27917, 27925, 27930, 27934,
                 27937],
                dtype='int64', length=3730)
    |           |           |
1 import pandas as pd
2 import matplotlib.pyplot as plt
3
4 # Select the column for outlier detection
5 column_name = 'user_followers'
6
7 # Detect outliers using z-score
8 def detect_outliers_zscore(data, threshold=3):
9     z_scores = (data - data.mean()) / data.std()
10    return z_scores.abs() > threshold
11
12 # Apply outlier detection on the specified column
13 outliers_zscore = detect_outliers_zscore(df[column_name])
14
15 # Get the outliers indices
16 outliers_indices = df[outliers_zscore].index
17
18 # Create a box plot
19 plt.figure(figsize=(8, 6))
20 plt.boxplot(df[column_name])
21 plt.title('Box Plot of {}'.format(column_name))
22 plt.ylabel(column_name)
23 plt.show()
24
25 # Create a scatter plot highlighting the outliers
26 plt.figure(figsize=(8, 6))
27 plt.scatter(df.index, df[column_name], color='blue', label='Data')
28 plt.scatter(outliers_indices, df.loc[outliers_indices, column_name], color='red', label='Outliers')
29 plt.title('Scatter Plot of {}'.format(column_name))
30 plt.xlabel('Index')
31 plt.ylabel(column_name)
32 plt.legend()
33 plt.show()
34

```

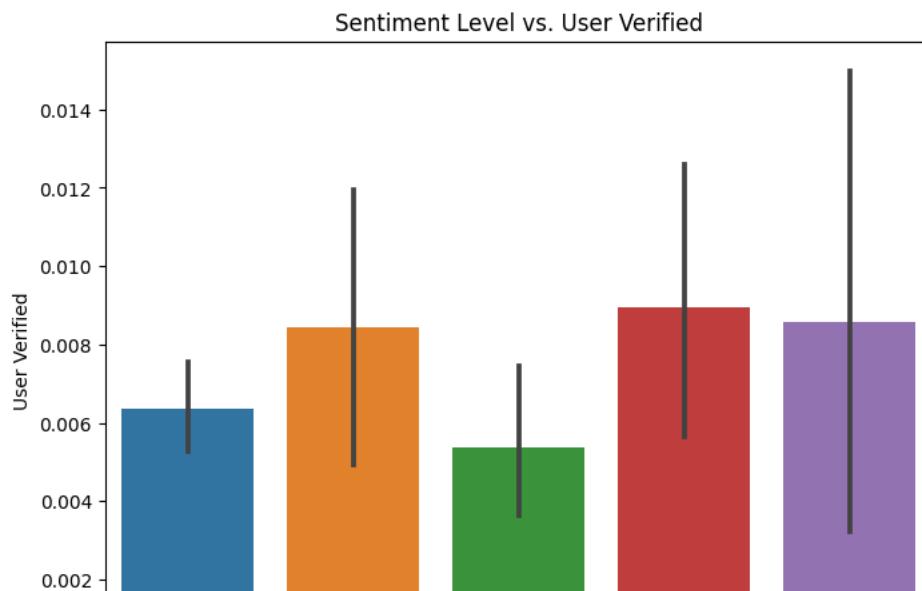
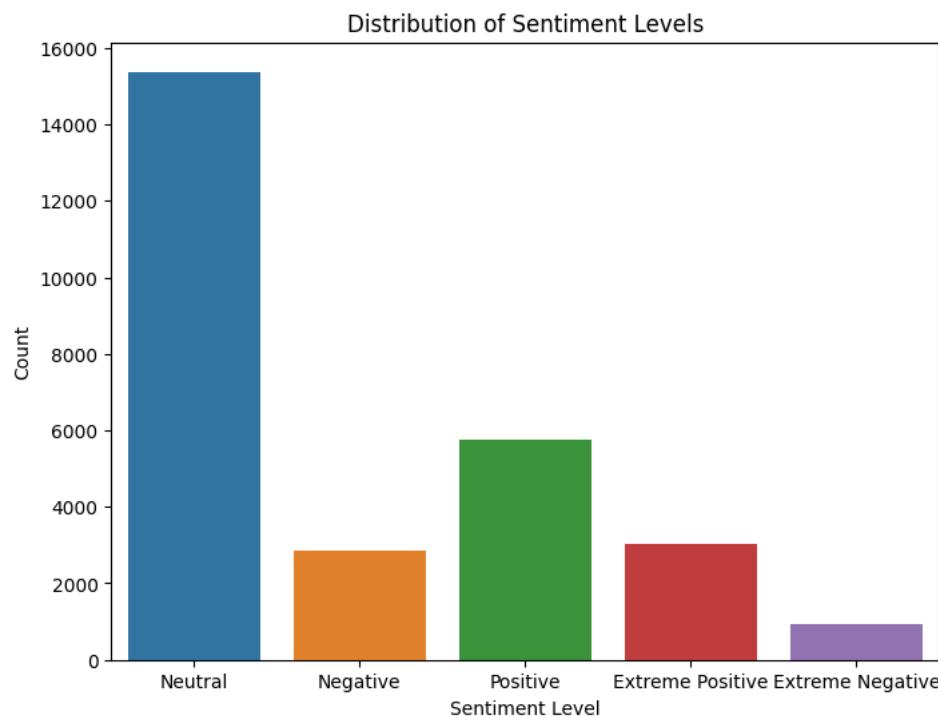


```

1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4
5 # Convert 'user_verified' column to numeric representation
6 df['user_verified_numeric'] = df['user_verified'].astype(int)
7
8 # EDA for sentiment_level
9
10 # Countplot of sentiment_level
11 plt.figure(figsize=(8, 6))
12 sns.countplot(x='sentiment_level', data=df)
13 plt.title('Distribution of Sentiment Levels')
14 plt.xlabel('Sentiment Level')
15 plt.ylabel('Count')
16 plt.show()
17
18 # Relationship with other attributes
19
20 # Bar plot of sentiment_level vs. user_verified
21 plt.figure(figsize=(8, 6))
22 sns.barplot(x='sentiment_level', y='user_verified_numeric', data=df)
23 plt.title('Sentiment Level vs. User Verified')
24 plt.xlabel('Sentiment Level')
25 plt.ylabel('User Verified')
26 plt.show()

```

```
27
28 # Bar plot of sentiment_level vs. user_location
29 #plt.figure(figsize=(12, 6))
30 #sns.countplot(x='sentiment_level', hue='user_location', data=df)
31 #plt.title('Sentiment Level vs. User Location')
32 #plt.xlabel('Sentiment Level')
33 #plt.ylabel('Count')
34 #plt.legend(title='User Location')
35 #plt.show()
36
37 # Scatter plot of sentiment_level vs. user_followers
38 plt.figure(figsize=(10, 6))
39 sns.scatterplot(x='sentiment_level', y='user_followers', data=df)
40 plt.title('Sentiment Level vs. User Followers')
41 plt.xlabel('Sentiment Level')
42 plt.ylabel('User Followers')
43 plt.show()
44
45 # Summary statistics of sentiment_level
46 sentiment_stats = df.groupby('sentiment_level').describe()['user_followers']
47 print("\nSummary Statistics of User Followers by Sentiment Level:")
48 print(sentiment_stats)
49
50 # Remove the temporary column
51 df.drop('user_verified_numeric', axis=1, inplace=True)
52
```



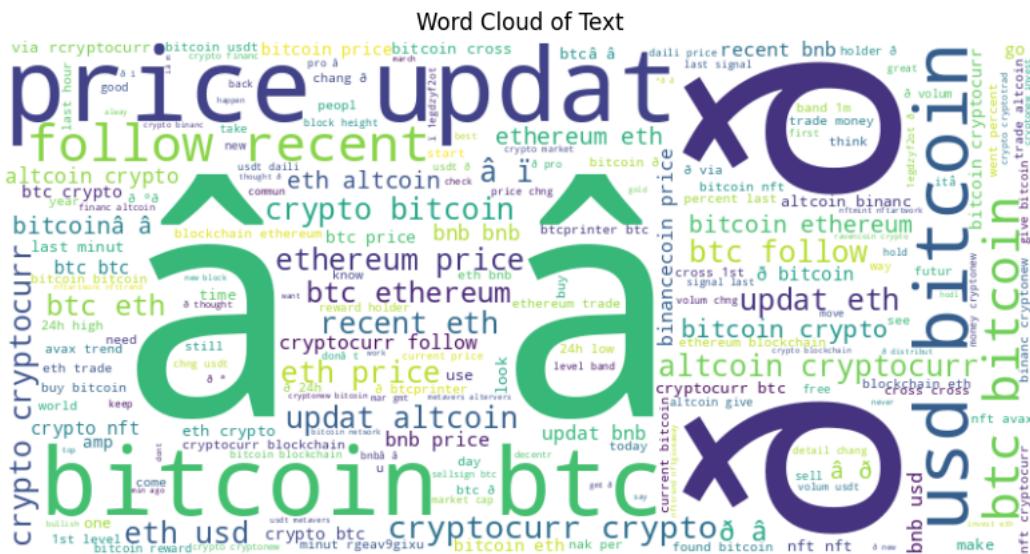
```
1 #Word Cloud
```

```
1 !pip install wordcloud
2
```

```
Looking in indexes: https://pypi.org/simple, https://us-python.pkg.dev/colab-wheels/public/simple/
Requirement already satisfied: wordcloud in /usr/local/lib/python3.10/dist-packages (1.8.2.2)
Requirement already satisfied: numpy>=1.6.1 in /usr/local/lib/python3.10/dist-packages (from wordcloud) (1.22.4)
Requirement already satisfied: pillow in /usr/local/lib/python3.10/dist-packages (from wordcloud) (8.4.0)
Requirement already satisfied: matplotlib in /usr/local/lib/python3.10/dist-packages (from wordcloud) (3.7.1)
Requirement already satisfied: contourpy>=1.0.1 in /usr/local/lib/python3.10/dist-packages (from matplotlib->wordcloud) (1.1.0)
Requirement already satisfied: cycler>=0.10 in /usr/local/lib/python3.10/dist-packages (from matplotlib->wordcloud) (0.11.0)
Requirement already satisfied: fonttools>=4.22.0 in /usr/local/lib/python3.10/dist-packages (from matplotlib->wordcloud) (4.40.0)
Requirement already satisfied: kiwisolver>=1.0.1 in /usr/local/lib/python3.10/dist-packages (from matplotlib->wordcloud) (1.4.4)
Requirement already satisfied: packaging>=20.0 in /usr/local/lib/python3.10/dist-packages (from matplotlib->wordcloud) (23.1)
Requirement already satisfied: pyparsing>=2.3.1 in /usr/local/lib/python3.10/dist-packages (from matplotlib->wordcloud) (3.1.0)
Requirement already satisfied: python-dateutil>=2.7 in /usr/local/lib/python3.10/dist-packages (from matplotlib->wordcloud) (2.8.2)
Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.10/dist-packages (from python-dateutil>=2.7->matplotlib->wordcloud) (
```

Word Cloud

```
1 from wordcloud import WordCloud
2
3 # Concatenate all the text data
4 text_data = ' '.join(df['text'].astype(str))
5
6 # Create a word cloud
7 wordcloud = WordCloud(width=800, height=400, background_color='white').generate(text_data)
8
9 # Plot the word cloud
10 plt.figure(figsize=(10, 6))
11 plt.imshow(wordcloud, interpolation='bilinear')
12 plt.axis('off')
13 plt.title('Word Cloud of Text')
14 plt.show()
15
```



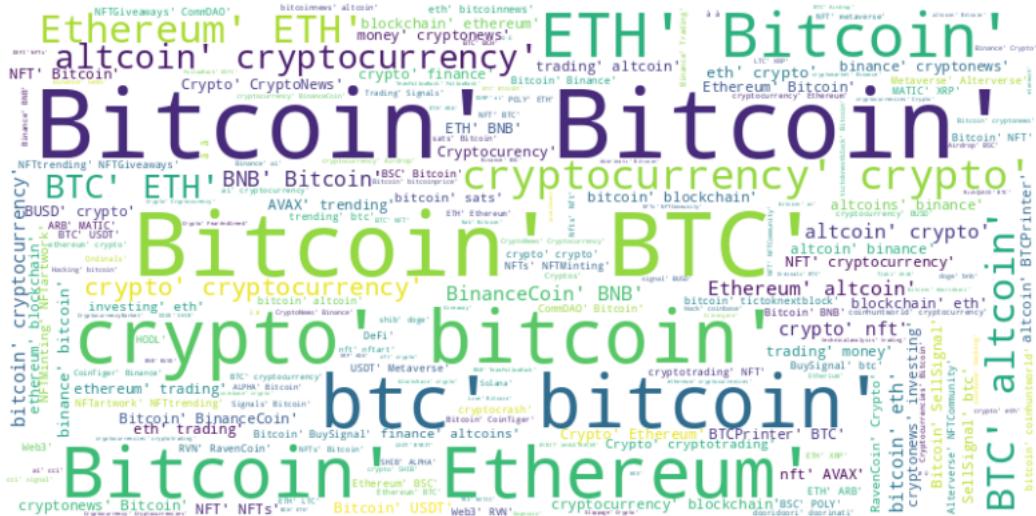
```
1 from wordcloud import WordCloud
2
3 # Concatenate all the user location data
4 user_location_data = ' '.join(df['user_location'].astype(str))
5
6 # Create a word cloud
7 wordcloud = WordCloud(width=800, height=400, background_color='white').generate(user_location_data)
8
9 # Plot the word cloud
10 plt.figure(figsize=(10, 6))
11 plt.imshow(wordcloud, interpolation='bilinear')
12 plt.axis('off')
13 plt.title('Word Cloud of User Locations')
14 plt.show()
15
```

Word Cloud of User Locations



```
1 from wordcloud import WordCloud
2
3 # Concatenate all the hashtags data
4 hashtags_data = ' '.join(df['hashtags'].astype(str))
5
6 # Create a word cloud
7 wordcloud = WordCloud(width=800, height=400, background_color='white').generate(hashtags_data)
8
9 # Plot the word cloud
10 plt.figure(figsize=(10, 6))
11 plt.imshow(wordcloud, interpolation='bilinear')
12 plt.axis('off')
13 plt.title('Word Cloud of Hashtags')
14 plt.show()
15
```

Word Cloud of Hashtags



```
1 from wordcloud import WordCloud
2
3 # Create a dictionary with sentiment levels and their counts
4 label_counts = {
5     'Neutral': 93169,
6     'Positive': 35921,
7     'Extreme Positive': 17343,
8     'Negative': 15903,
9     'Extreme Negative': 5316
10 }
11
12 # Generate the word cloud
13 wordcloud = WordCloud(width=800, height=400, background_color='white').generate_from_frequencies(label_counts)
14
15 # Plot the word cloud
16 plt.figure(figsize=(10, 6))
17 plt.imshow(wordcloud, interpolation='bilinear')
18 plt.axis('off')
19 plt.title('Word Cloud of Sentiment Levels')
20 plt.show()
21
```



▼ Sentiment Analysis Vs Price [Close]

EXTREME NEGATIVE

1 #Grouping the data by hours allows you to analyze the sentiment trends at a more granular level and grouping by hours can help in aggregat

EXTREME POSITIVE

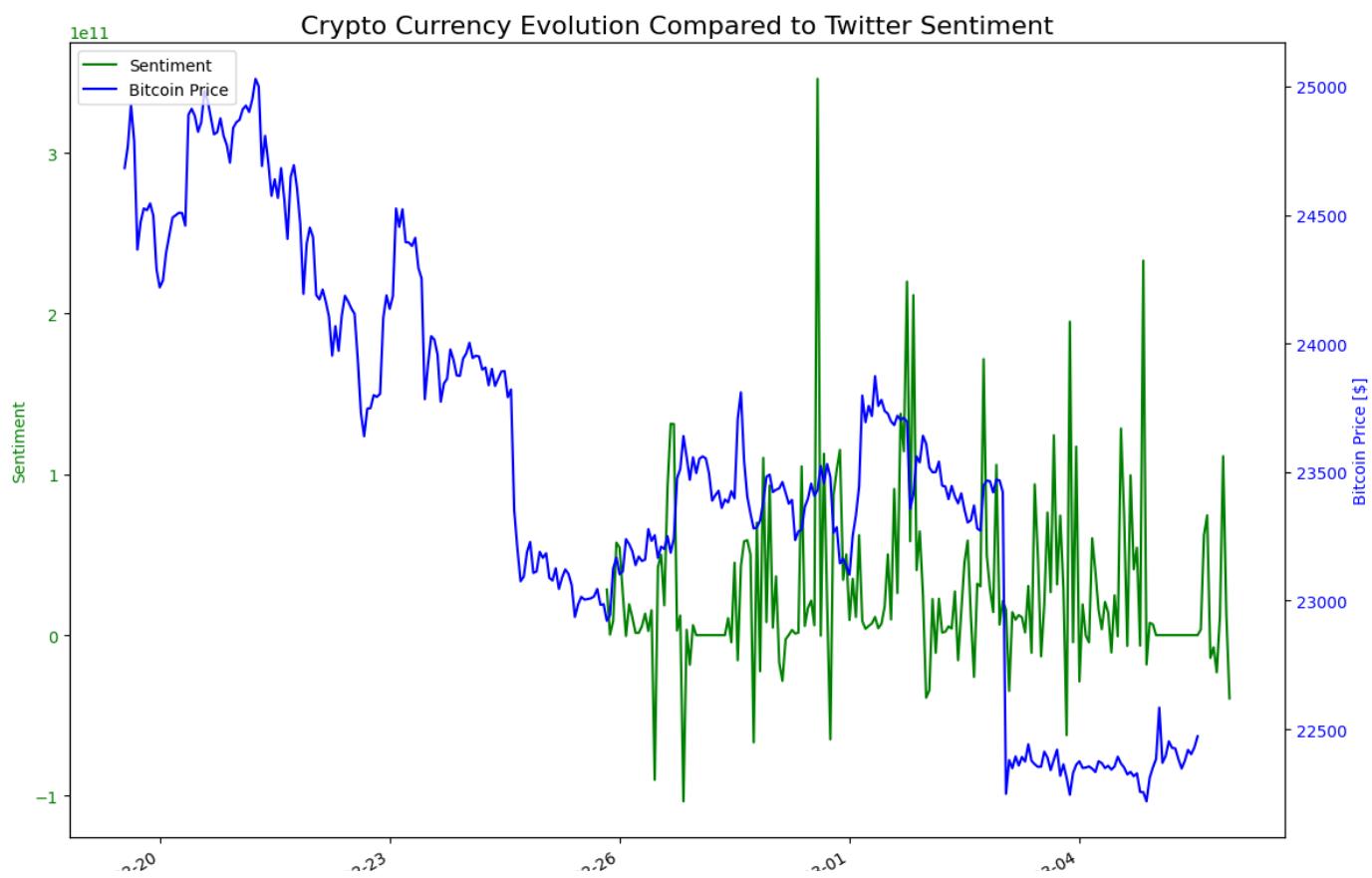
```

1 import pandas as pd
2
3 # Assuming you have a DataFrame named 'tweets'
4 tweets['date'] = pd.to_datetime(tweets['date'])
5 tweets.set_index('date', inplace=True)
6
7 # Group by hour and sum the 'score' column
8 tweets_grouped = tweets['score'].resample('1H').sum()

1 #'crypto_usd' DataFrame
2 crypto_usd['time'] = pd.to_datetime(crypto_usd['time'])
3 crypto_usd.set_index('time', inplace=True)
4
5 # Group by hour and calculate the mean of the 'close' column
6 crypto_usd_grouped = crypto_usd['close'].resample('1H').mean()

1 import matplotlib.pyplot as plt
2
3 fig, ax1 = plt.subplots(figsize=(12, 8))
4 ax2 = ax1.twinx()
5
6 # Plot sentiment analysis data
7 ax1.plot_date(tweets_grouped.index, tweets_grouped, 'g-', label='Sentiment')
8 ax1.set_ylabel("Sentiment", color='g')
9 ax1.tick_params(axis='y', labelcolor='g')
10
11 # Plot cryptocurrency price data
12 ax2.plot_date(crypto_usd_grouped.index, crypto_usd_grouped, 'b-', label='Bitcoin Price')
13 ax2.set_ylabel("Bitcoin Price [$]", color='b')
14 ax2.tick_params(axis='y', labelcolor='b')
15
16 # Set title and x-axis label
17 plt.title("Crypto Currency Evolution Compared to Twitter Sentiment", fontsize=16)
18 plt.xlabel("Date")
19
20 # Add legend
21 lines = ax1.get_lines() + ax2.get_lines()
22 labels = [line.get_label() for line in lines]
23 plt.legend(lines, labels, loc='upper left')
24
25 # Format the x-axis tick labels
26 ax1.xaxis.set_major_locator(plt.MaxNLocator(6)) # Set maximum number of x-axis tick labels
27 fig.autofmt_xdate() # Automatically format x-axis tick labels as dates
28
29 plt.tight_layout()
30 plt.show()
31

```



```

1 # Define the cross-correlation function
2
3 def crosscorr(datax, datay, lag=0, method="pearson"):
4     """ Lag-N cross correlation.
5     Parameters
6     -----
7     lag : int, default 0
8     datax, datay : pandas.Series objects of equal length
9
10    Returns
11    -----
12    crosscorr : float
13    """
14    return datax.corr(datay.shift(lag), method=method)

1 #Crop the time series to match the time frames

1 beginning = max(tweets_grouped.index.min(), crypto_usd_grouped.index.min())
2 end = min(tweets_grouped.index.max(), crypto_usd_grouped.index.max())
3 tweets_grouped = tweets_grouped[begginning:end]
4 crypto_usd_grouped = crypto_usd_grouped[begginning:end]

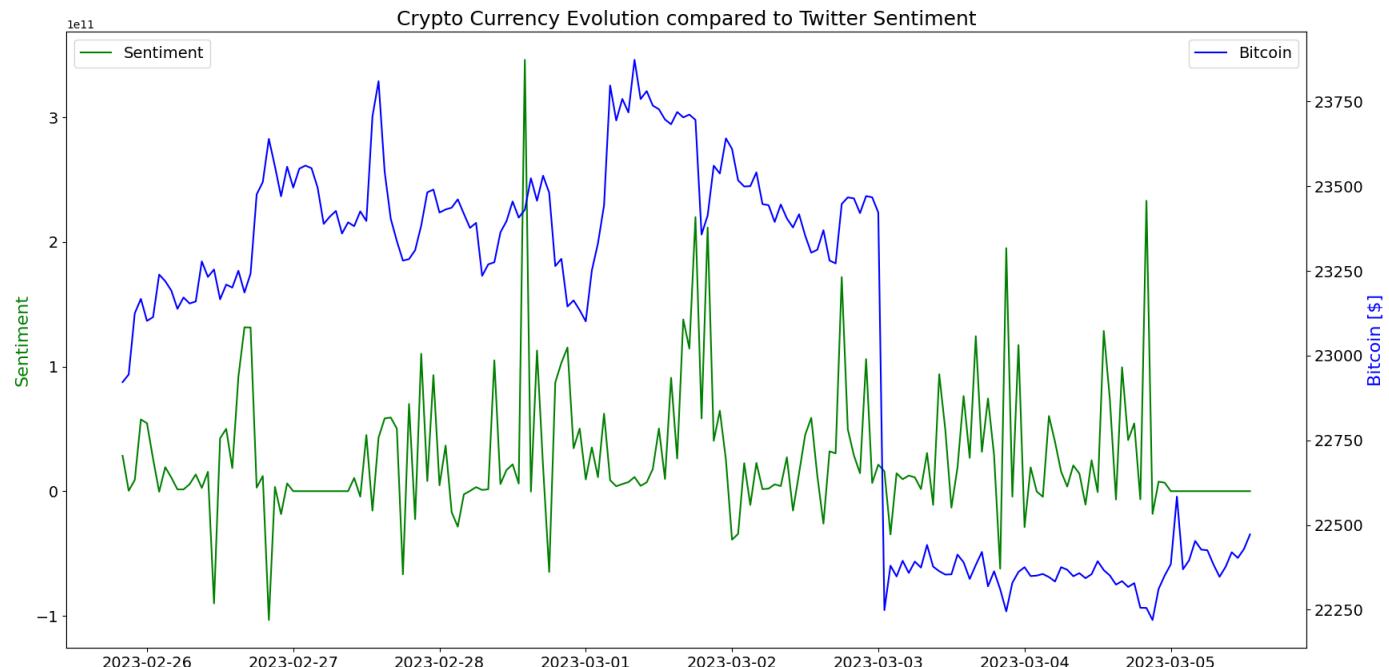
1 import matplotlib.pyplot as plt
2 import matplotlib.dates as mdates
3
4 # Create figure and axes
5 fig, ax1 = plt.subplots(figsize=(20, 10))
6 ax2 = ax1.twinx()
7
8 # Set title and axis labels
9 ax1.set_title("Crypto Currency Evolution compared to Twitter Sentiment", fontsize=18)
10 ax1.set_ylabel("Sentiment", color='g', fontsize=16)
11 ax2.set_ylabel("Bitcoin [$]", color='b', fontsize=16)
12
13 # Set tick label size
14 ax1.tick_params(labelsize=14)
15 ax2.tick_params(labelsize=14)

```

```

16
17 # Plot sentiment data
18 ax1.plot_date(tweets_grouped.index, tweets_grouped, 'g-', label="Sentiment")
19
20 # Plot cryptocurrency data
21 ax2.plot_date(crypto_usd_grouped.index, crypto_usd_grouped, 'b-', label="Bitcoin")
22
23 # Format x-axis as dates
24 ax1.xaxis.set_major_locator(mdates.AutoDateLocator())
25 ax1.xaxis.set_major_formatter(mdates.DateFormatter('%Y-%m-%d'))
26
27 # Rotate x-axis tick labels for better visibility
28 plt.xticks(rotation=45)
29
30 # Display legend
31 ax1.legend(loc='upper left', fontsize=14)
32 ax2.legend(loc='upper right', fontsize=14)
33
34 # Show the plot
35 plt.show()
36

```



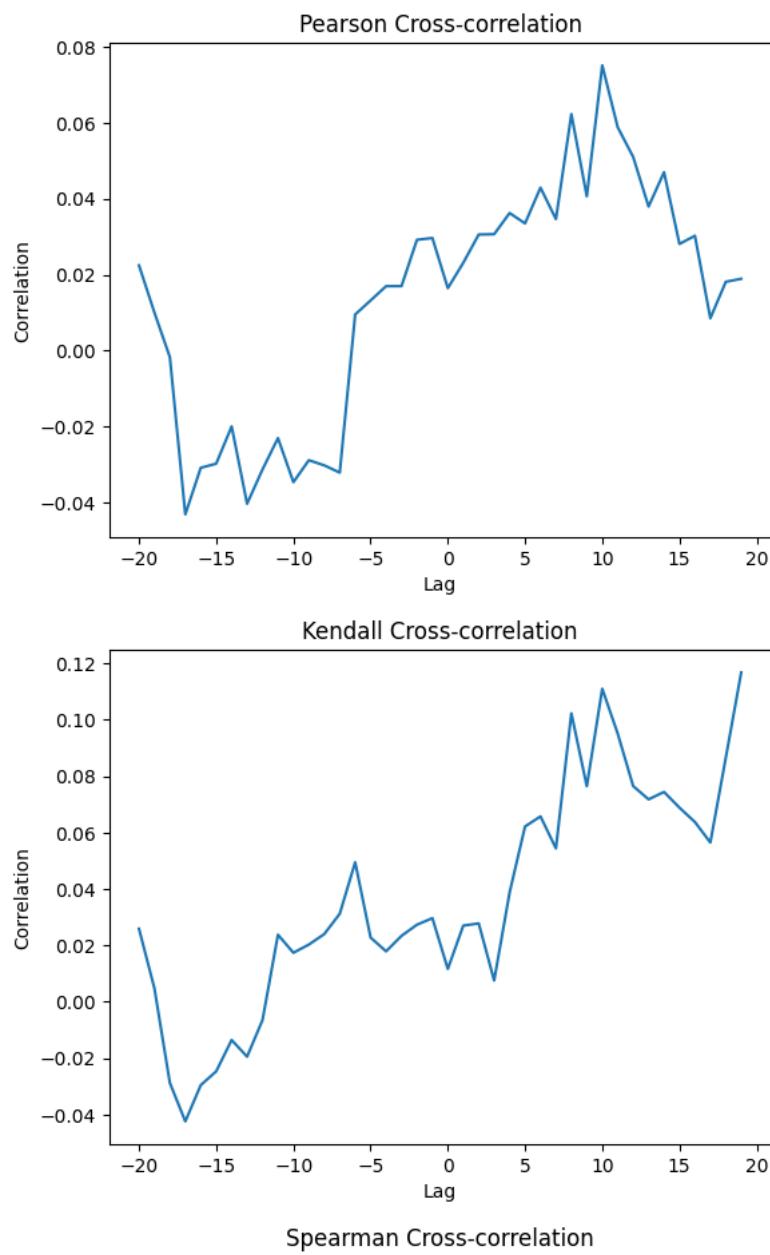
```

1 #The code imports the matplotlib.pyplot library for plotting and matplotlib.dates library for date formatting

1 #The plot visually suggests a potential relationship between sentiment extracted from Twitter data and the price evolution of Bitcoin, in

1 import matplotlib.pyplot as plt
2
3 methods = ["pearson", "kendall", "spearman"]
4
5 for method in methods:
6     xcov = [crosscorr(tweets_grouped, crypto_usd_grouped, lag=i, method=method) for i in range(-20, 20)]
7     plt.plot(range(-20, 20), xcov)
8     plt.title(f"{method.capitalize()} Cross-correlation")
9     plt.xlabel("Lag")
10    plt.ylabel("Correlation")
11    plt.show()
12

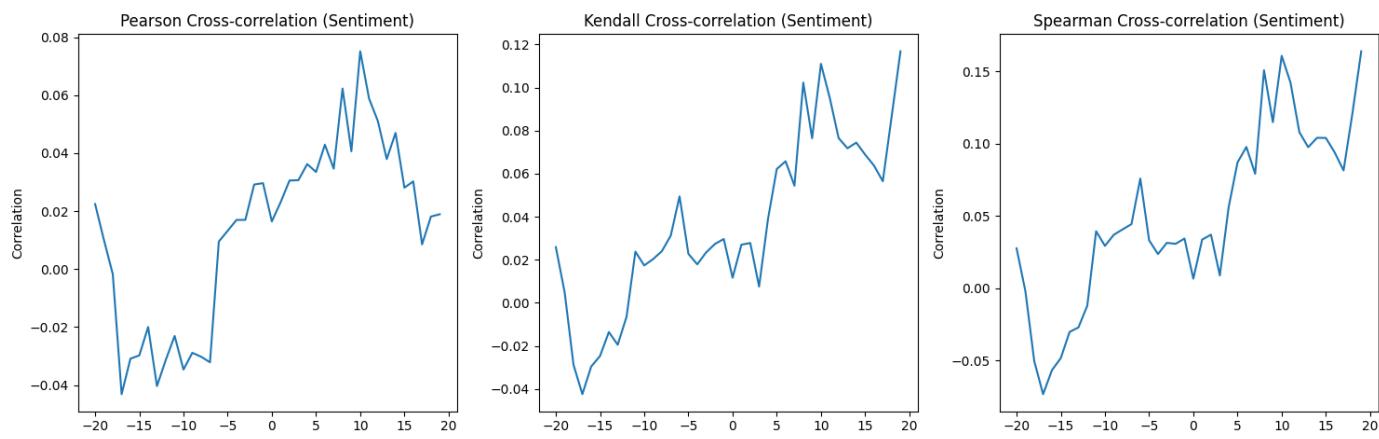
```



```

1 #To analyze the results and interpret the cross-correlation plots
2
3 methods = ["pearson", "kendall", "spearman"]
4
5 fig, axs = plt.subplots(1, 3, figsize=(15, 5))
6
7 # Sentiment Analysis
8 for i, method in enumerate(methods):
9     xcov_sentiment = [crosscorr(tweets_grouped, crypto_usd_grouped, lag=j, method=method) for j in range(-20, 20)]
10    axs[i].plot(range(-20, 20), xcov_sentiment)
11    axs[i].set_title(f"{method.capitalize()} Cross-correlation (Sentiment)")
12    axs[i].set_xlabel("Lag")
13    axs[i].set_ylabel("Correlation")
14
15 plt.tight_layout()
16 plt.show()
17

```



```

1 from scipy.stats import pearsonr, kendalltau, spearmanr
2
3 methods = ["pearson", "kendall", "spearman"]
4
5 for method in methods:
6     if method == "pearson":
7         correlation_coefficient, p_value = pearsonr(tweets_grouped, crypto_usd_grouped)
8     elif method == "kendall":
9         correlation_coefficient, p_value = kendalltau(tweets_grouped, crypto_usd_grouped)
10    elif method == "spearman":
11        correlation_coefficient, p_value = spearmanr(tweets_grouped, crypto_usd_grouped)
12
13    print(f"{method.capitalize()} correlation coefficient:", correlation_coefficient)
14    print(f"{method.capitalize()} p-value:", p_value)
15
16    if p_value < 0.05:
17        print(f"The {method} correlation is statistically significant at this confidence level.")
18    else:
19        print(f"The {method} correlation is not statistically significant at this confidence level.")
20
21

```

Pearson correlation coefficient: 0.01643626585537878
 Pearson p-value: 0.8237964874502892
 The pearson correlation is not statistically significant at this confidence level.

Kendall correlation coefficient: 0.01166030969779228
 Kendall p-value: 0.8144610941602477
 The kendall correlation is not statistically significant at this confidence level.

Spearman correlation coefficient: 0.006553754143596371
 Spearman p-value: 0.9292569945760621
 The spearman correlation is not statistically significant at this confidence level.

```
1 #The Pearson, Kendall, and Spearman correlation coefficients between the sentiment and Bitcoin prices are all very close to zero, indicating no strong linear relationship.
```

```

1 # Identify interesting correlations
2 interesting_correlations = []
3
4 # Loop through different lag values
5 for lag in range(-20, 20):
6     # Calculate cross-correlation
7     correlation = crosscorr(tweets_grouped, crypto_usd_grouped, lag=lag, method="pearson")
8
9     # Check if correlation is interesting (adjust criteria as needed)
10    if abs(correlation) > 0.5:
11        interesting_correlations.append((lag, correlation))
12
13 # Print the interesting correlations
14 for lag, correlation in interesting_correlations:
15     print(f"Correlation at lag {lag}: {correlation}")
16
17 # Conduct additional analysis for the interesting correlations
18 for lag, correlation in interesting_correlations:
19     # Perform specific analysis for each interesting correlation
20     # Example: Analyze the relationship during specific time periods or incorporate additional variables

```

```
21     print(f"Additional analysis for correlation at lag {lag}")
22

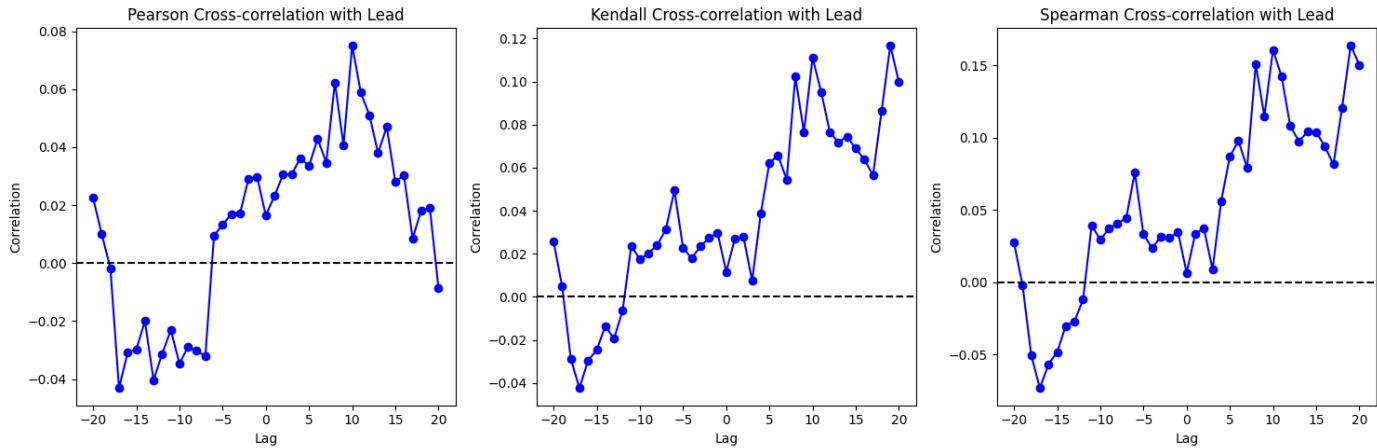
1 #The cross-correlation with different lags for the lead (future) relationship between sentiment and cryptocurrency prices

1 import matplotlib.pyplot as plt
2
3 methods = ["pearson", "kendall", "spearman"]
4 lags = range(-20, 21) # Include both positive and negative lags
5
6 for method in methods:
7     for lag in lags:
8         xcov = crosscorr(tweets_grouped, crypto_usd_grouped, lag=lag, method=method)
9         plt.plot(lag, xcov, marker='o', color='b')
10
11     plt.title(f"{method.capitalize()} Cross-correlation with Lead")
12     plt.xlabel("Lag")
13     plt.ylabel("Correlation")
14     plt.axhline(0, color='black', linestyle='--') # Add a horizontal line at y=0 for reference
15     plt.show()
16
```

```

1 import matplotlib.pyplot as plt
2
3 methods = ["pearson", "kendall", "spearman"]
4 lags = range(-20, 21) # Include both positive and negative lags
5
6 fig, axs = plt.subplots(1, 3, figsize=(15, 5)) # Create a 1x3 grid of subplots
7
8 for i, method in enumerate(methods):
9     correlations = [crosscorr(tweets_grouped, crypto_usd_grouped, lag=lag, method=method) for lag in lags]
10    axs[i].plot(lags, correlations, marker='o', color='b')
11    axs[i].set_title(f"{method.capitalize()} Cross-correlation with Lead")
12    axs[i].set_xlabel("Lag")
13    axs[i].set_ylabel("Correlation")
14    axs[i].axhline(0, color='black', linestyle='--') # Add a horizontal line at y=0 for reference
15
16 plt.tight_layout() # Adjust spacing between subplots
17 plt.show()
18

```

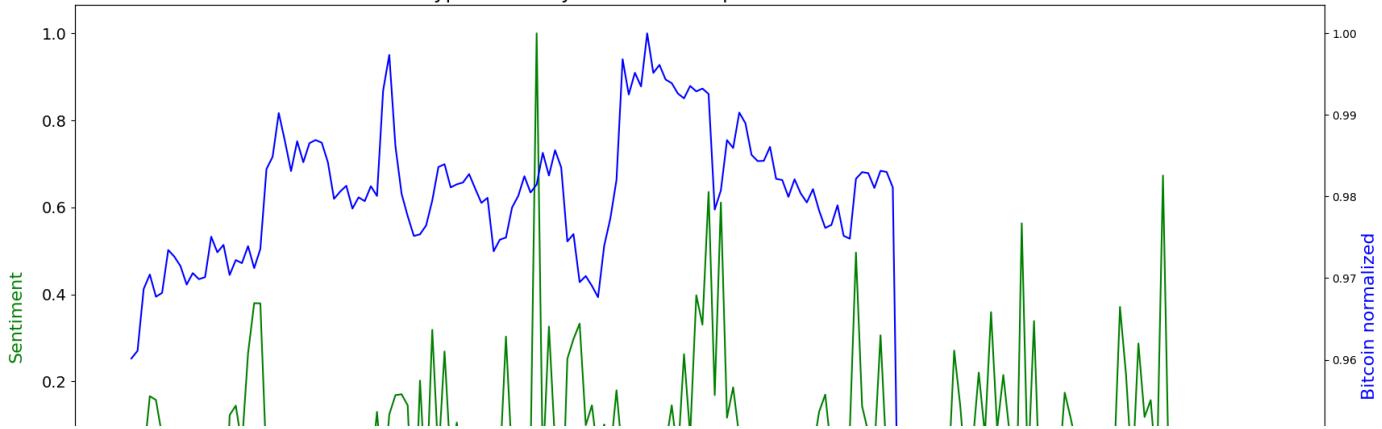


```

1 #The term "lead" suggests that the correlation is examined for the sentiment leading the Bitcoin prices at various time shifts. The plot h
| | |
1 #Normalize
| | |
1 import matplotlib.pyplot as plt
2
3 # Normalize time series data
4 tweets_grouped_normalized = tweets_grouped / max(tweets_grouped.max(), abs(tweets_grouped.min()))
5 crypto_usd_grouped_normalized = crypto_usd_grouped / max(crypto_usd_grouped.max(), abs(crypto_usd_grouped.min()))
6
7 # Create the plot
8 fig, ax1 = plt.subplots(figsize=(20, 10))
9 ax1.set_title("Normalized Crypto currency evolution compared to normalized Twitter sentiment", fontsize=18)
10 ax1.tick_params(labelsize=14)
11
12 ax2 = ax1.twinx()
13 ax1.plot_date(tweets_grouped_normalized.index, tweets_grouped_normalized, 'g-')
14 ax2.plot_date(crypto_usd_grouped_normalized.index, crypto_usd_grouped_normalized, 'b-')
15
16 ax1.set_ylabel("Sentiment", color='g', fontsize=16)
17 ax2.set_ylabel("Bitcoin normalized", color='b', fontsize=16)
18
19 plt.show()
20

```

Normalized Crypto currency evolution compared to normalized Twitter sentiment

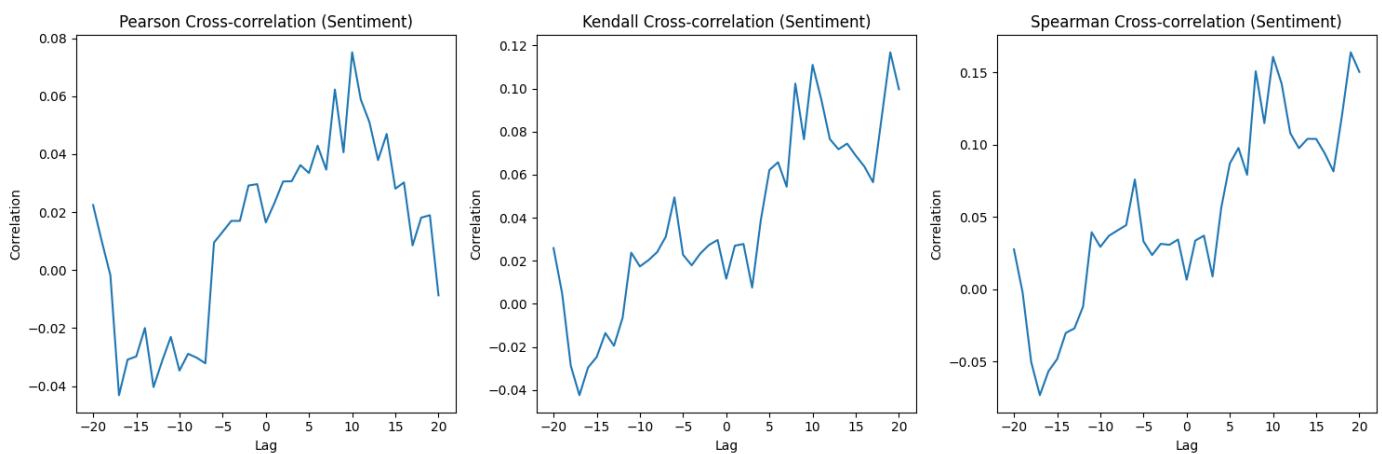


1 #The plot allows for visual analysis of the relationship between the two variables, after normalizing the data to a common scale.

```

1 import matplotlib.pyplot as plt
2
3 methods = ["pearson", "kendall", "spearman"]
4 lags = range(-20, 21) # Include both positive and negative lags
5
6 fig, axs = plt.subplots(1, 3, figsize=(15, 5))
7
8 # Normalize time series data
9 tweets_grouped_norm = (tweets_grouped - tweets_grouped.mean()) / tweets_grouped.std()
10 crypto_usd_grouped_norm = (crypto_usd_grouped - crypto_usd_grouped.mean()) / crypto_usd_grouped.std()
11
12 # Sentiment Analysis
13 for i, method in enumerate(methods):
14     xcov_sentiment = [crosscorr(tweets_grouped_norm, crypto_usd_grouped_norm, lag=j, method=method) for j in lags]
15     axs[i].plot(lags, xcov_sentiment)
16     axs[i].set_title(f"{method.capitalize()} Cross-correlation (Sentiment)")
17     axs[i].set_xlabel("Lag")
18     axs[i].set_ylabel("Correlation")
19
20 plt.tight_layout()
21 plt.show()
22

```



```

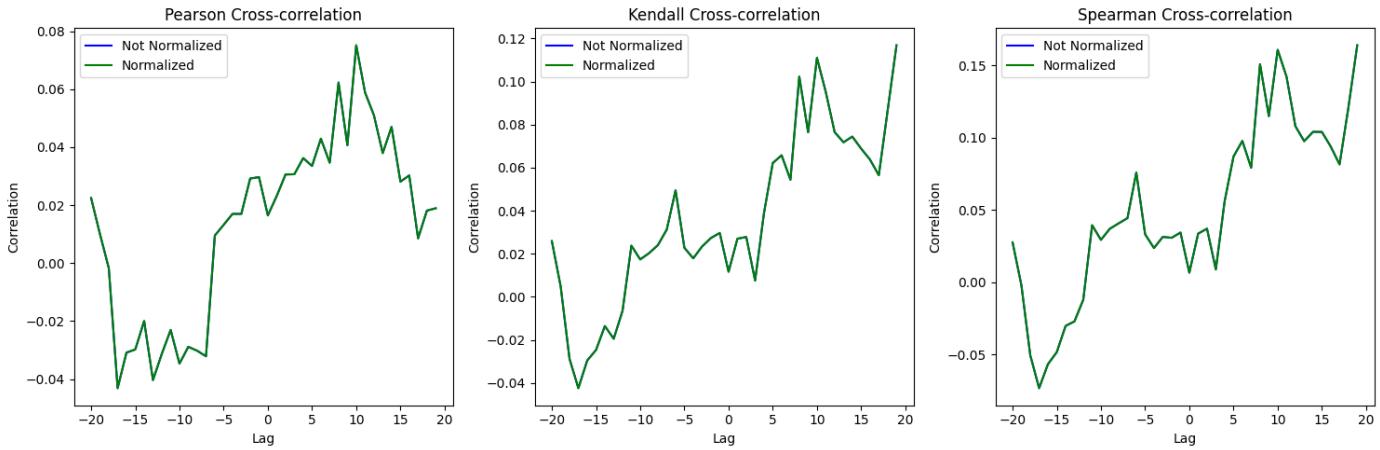
1 import matplotlib.pyplot as plt
2
3 methods = ["pearson", "kendall", "spearman"]
4
5 fig, axs = plt.subplots(1, 3, figsize=(15, 5))
6
7 for i, method in enumerate(methods):
8     # Calculate cross-correlation without normalization
9     xcov = [crosscorr(tweets_grouped, crypto_usd_grouped, lag=j, method=method) for j in range(-20, 20)]
10    axs[i].plot(range(-20, 20), xcov, color='blue', label='Not Normalized')
11
12 # Calculate cross-correlation with normalization

```

```

13     xcov_norm = [crosscorr(tweets_grouped_norm, crypto_usd_grouped_norm, lag=j, method=method) for j in range(-20, 20)]
14     axs[i].plot(range(-20, 20), xcov_norm, color='green', label='Normalized')
15
16     axs[i].set_title(f"{method.capitalize()} Cross-correlation")
17     axs[i].set_xlabel("Lag")
18     axs[i].set_ylabel("Correlation")
19     axs[i].legend()
20
21 plt.tight_layout()
22 plt.show()
23

```



▼ Sentiment Analysis Vs Price Change Delta

```

1 print(crypto_usd.columns)

Index(['close', 'high', 'low', 'open', 'volumefrom', 'volumeto', 'Date',
       'Time', 'volume', 'marketcap', 'price_delta'],
      dtype='object')

1 # Group by hour and calculate the mean of the 'close' column
2 crypto_usd_grouped2 = crypto_usd['price_delta'].resample('1H').mean()

1 beggining = max(tweets_grouped.index.min(), crypto_usd_grouped2.index.min())
2 end = min(tweets_grouped.index.max(), crypto_usd_grouped2.index.max())
3 tweets_grouped = tweets_grouped[beggining:end]
4 crypto_usd_grouped2 = crypto_usd_grouped2[beggining:end]

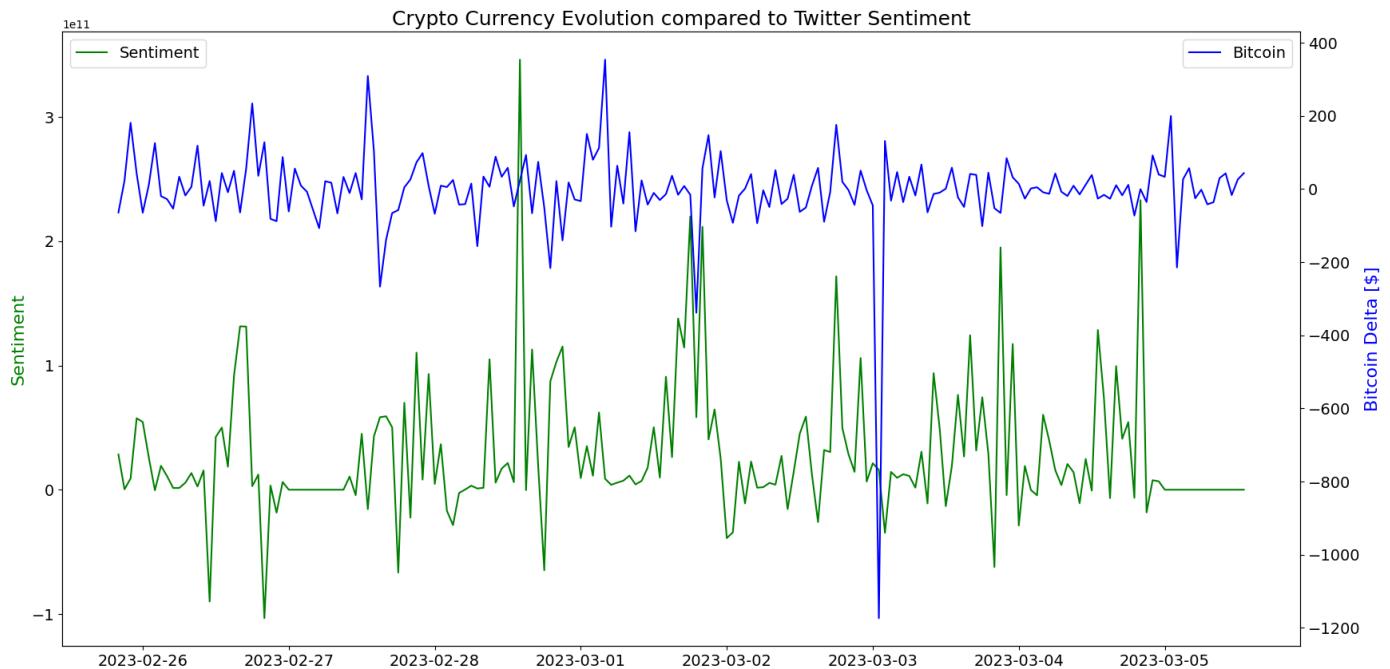
1 import matplotlib.pyplot as plt
2 import matplotlib.dates as mdates
3
4 # Create figure and axes
5 fig, ax1 = plt.subplots(figsize=(20, 10))
6 ax2 = ax1.twinx()
7
8 # Set title and axis labels
9 ax1.set_title("Crypto Currency Evolution compared to Twitter Sentiment", fontsize=18)
10 ax1.set_ylabel("Sentiment", color='g', fontsize=16)
11 ax2.set_ylabel("Bitcoin Delta [$]", color='b', fontsize=16)
12
13 # Set tick label size
14 ax1.tick_params(labelsize=14)
15 ax2.tick_params(labelsize=14)
16
17 # Plot sentiment data
18 ax1.plot_date(tweets_grouped.index, tweets_grouped, 'g-', label="Sentiment")
19
20 # Plot cryptocurrency data
21 ax2.plot_date(crypto_usd_grouped2.index, crypto_usd_grouped2, 'b-', label="Bitcoin")
22
23 # Format x-axis as dates
24 ax1.xaxis.set_major_locator(mdates.AutoDateLocator())

```

```

25 ax1.xaxis.set_major_formatter(mdates.DateFormatter('%Y-%m-%d'))
26
27 # Rotate x-axis tick labels for better visibility
28 plt.xticks(rotation=45)
29
30 # Display legend
31 ax1.legend(loc='upper left', fontsize=14)
32 ax2.legend(loc='upper right', fontsize=14)
33
34 # Show the plot
35 plt.show()
36

```



```

1 # Group by hour and sum the 'score' column
2 tweets_grouped = tweets['score'].resample('1H').sum()
3 # Group by hour and calculate the mean of the 'close' column
4 crypto_usd_grouped = crypto_usd['close'].resample('1H').mean()
5 # Group by hour and calculate the mean of the 'close' column
6 crypto_usd_grouped2 = crypto_usd['price_delta'].resample('1H').mean()

1 beggining = max(tweets_grouped.index.min(), crypto_usd_grouped.index.min())
2 end = min(tweets_grouped.index.max(), crypto_usd_grouped.index.max())
3 tweets_grouped = tweets_grouped[beggining:end]
4 crypto_usd_grouped = crypto_usd_grouped[beggining:end]
5 crypto_usd_grouped2 = crypto_usd_grouped2[beggining:end]

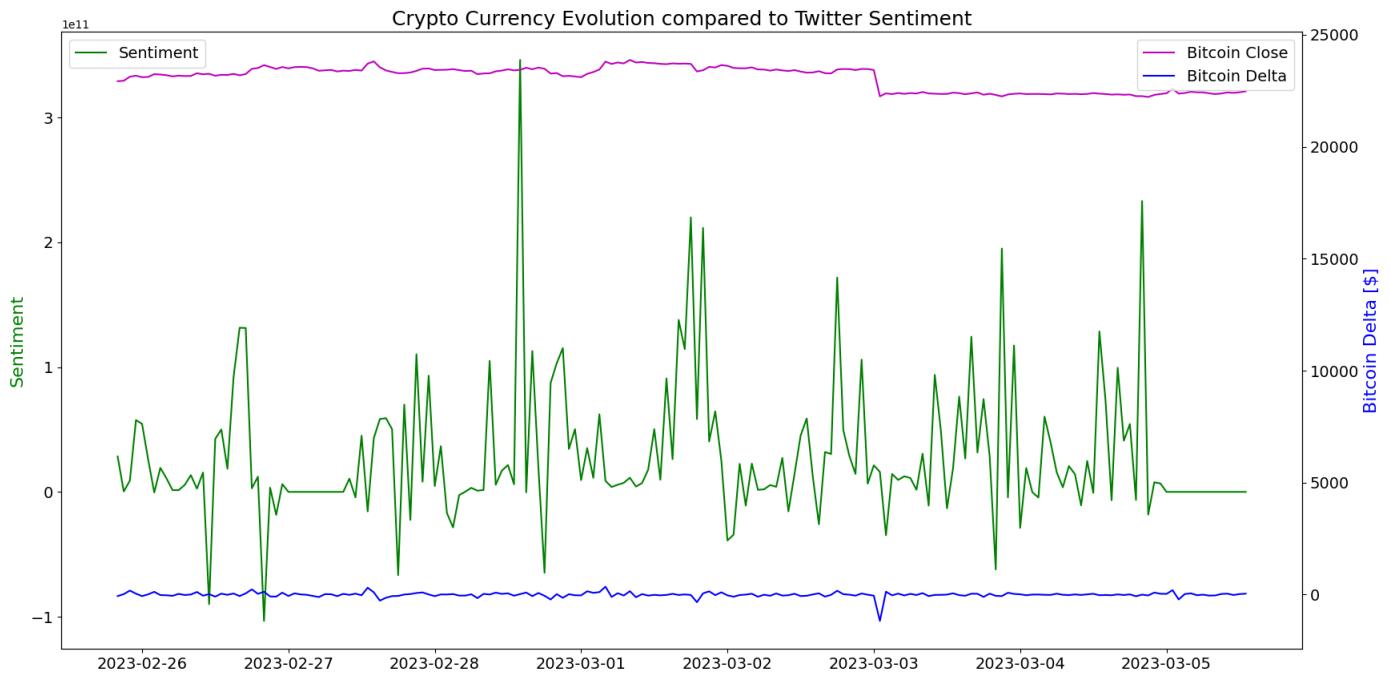
1 import matplotlib.pyplot as plt
2 import matplotlib.dates as mdates
3
4 # Create figure and axes
5 fig, ax1 = plt.subplots(figsize=(20, 10))
6 ax2 = ax1.twinx()
7
8 # Set title and axis labels
9 ax1.set_title("Crypto Currency Evolution compared to Twitter Sentiment", fontsize=18)
10 ax1.set_ylabel("Sentiment", color='g', fontsize=16)
11 ax2.set_ylabel("Bitcoin Delta [$]", color='b', fontsize=16)
12
13 # Set tick label size
14 ax1.tick_params(labelsize=14)
15 ax2.tick_params(labelsize=14)
16
17 # Plot sentiment data
18 ax1.plot_date(tweets_grouped.index, tweets_grouped, 'g-', label="Sentiment")
19
20 # Plot cryptocurrency data (crypto_usd_grouped)

```

```

21 ax2.plot_date(crypto_usd_grouped.index, crypto_usd_grouped, 'm-', label="Bitcoin Close")
22
23 # Plot additional cryptocurrency data (crypto_usd_grouped2)
24 ax2.plot_date(crypto_usd_grouped2.index, crypto_usd_grouped2, 'b-', label="Bitcoin Delta")
25
26 # Format x-axis as dates
27 ax1.xaxis.set_major_locator(mdates.AutoDateLocator())
28 ax1.xaxis.set_major_formatter(mdates.DateFormatter('%Y-%m-%d'))
29
30 # Rotate x-axis tick labels for better visibility
31 plt.xticks(rotation=45)
32
33 # Display legend
34 ax1.legend(loc='upper left', fontsize=14)
35 ax2.legend(loc='upper right', fontsize=14)
36
37 # Show the plot
38 plt.show()
39

```



```

1 #the sentiment line is plotted based on the direction of the sentiment slope. If the sign of the sentiment slope at a
2
3 #This way, the areas where the sentiment line has the same direction as the delta line will be highlighted in red for better observation.

```

```
1 #consider the sentiment slope 1 hour before the delta slope to observe if the sentiment affects the price change
```

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3 import matplotlib.dates as mdates
4
5 # Create figure and axes
6 fig, ax1 = plt.subplots(figsize=(20, 10))
7 ax2 = ax1.twinx()
8
9 # Set title and axis labels
10 ax1.set_title("Crypto Currency Evolution compared to Twitter Sentiment", fontsize=18)
11 ax1.set_ylabel("Sentiment", color='g', fontsize=16)
12 ax2.set_ylabel("Bitcoin Delta [$]", color='b', fontsize=16)
13
14 # Set tick label size
15 ax1.tick_params(labelsize=14)
16 ax2.tick_params(labelsize=14)
17
18 # Calculate sentiment slope and delta slope
19 sentiment_slope = np.gradient(tweets_grouped)
20 delta_slope = np.gradient(crypto_usd_grouped2)

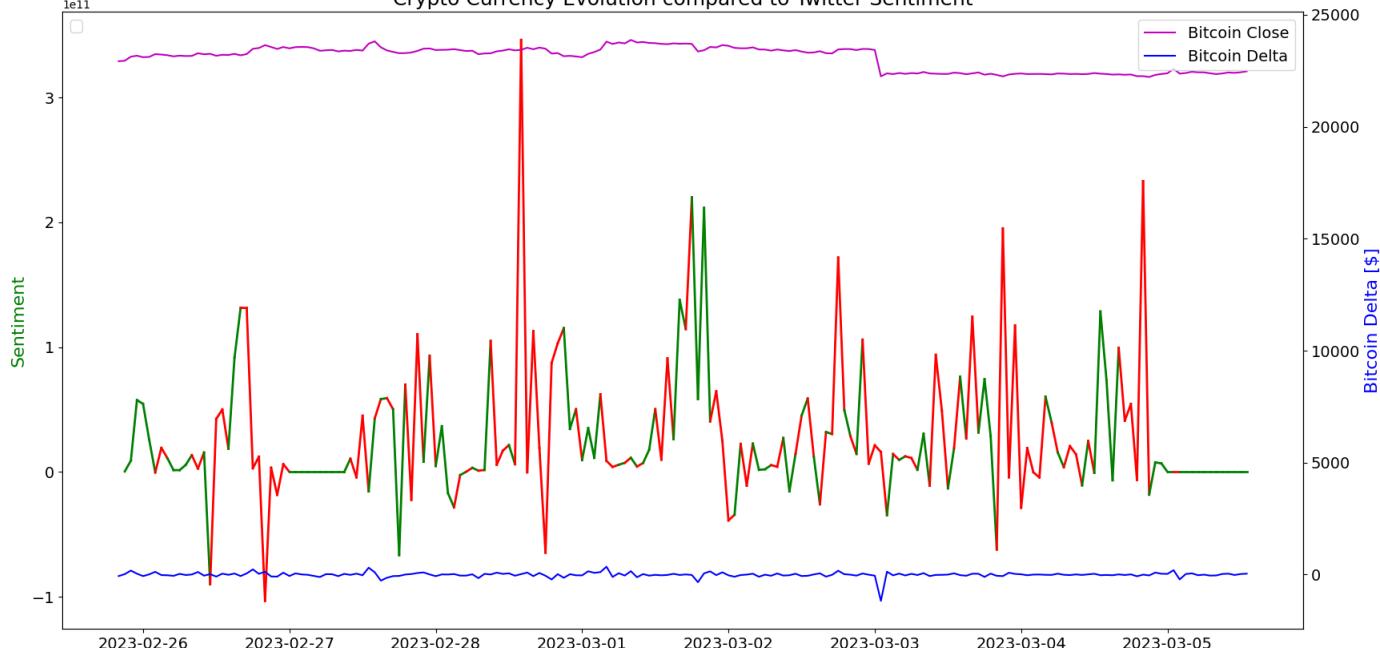
```

```

21
22 # Plot sentiment data
23 for i in range(1, len(tweets_grouped)-1):
24     if np.sign(sentiment_slope[i-1]) == np.sign(delta_slope[i]):
25         ax1.plot([tweets_grouped.index[i], tweets_grouped.index[i+1]], [tweets_grouped[i], tweets_grouped[i+1]], 'r-', linewidth=2)
26     else:
27         ax1.plot([tweets_grouped.index[i], tweets_grouped.index[i+1]], [tweets_grouped[i], tweets_grouped[i+1]], 'g-', linewidth=2)
28
29 # Plot cryptocurrency data (crypto_usd_grouped)
30 ax2.plot_date(crypto_usd_grouped.index, crypto_usd_grouped, 'm-', label="Bitcoin Close")
31
32 # Plot additional cryptocurrency data (crypto_usd_grouped2)
33 ax2.plot_date(crypto_usd_grouped2.index, crypto_usd_grouped2, 'b-', label="Bitcoin Delta")
34
35 # Format x-axis as dates
36 ax1.xaxis.set_major_locator(mdates.AutoDateLocator())
37 ax1.xaxis.set_major_formatter(mdates.DateFormatter('%Y-%m-%d'))
38
39 # Rotate x-axis tick labels for better visibility
40 plt.xticks(rotation=45)
41
42 # Display legend
43 ax1.legend(loc='upper left', fontsize=14)
44 ax2.legend(loc='upper right', fontsize=14)
45
46 # Show the plot
47 plt.show()
48

```

WARNING:matplotlib.legend:No artists with labels found to put in legend. Note that artists whose label start with an underscore are ignored.



```
1 #The plot visualizes the comparison between sentiment and Bitcoin delta, indicating the segments where their slopes have the same sign using
```

```
1 #Delta and Sentiment - consider the sentiment slope 1 hour before the delta slope to observe if the sentiment affects the price change
```

```
1 # If the sign of the sentiment slope at index i-2 is the same as the sign of the delta slope at index i, a red line ('r-') is plotted between
```

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3 import matplotlib.dates as mdates
4
5 # Create figure and axes
6 fig, ax1 = plt.subplots(figsize=(20, 10))
7 ax2 = ax1.twinx()
8

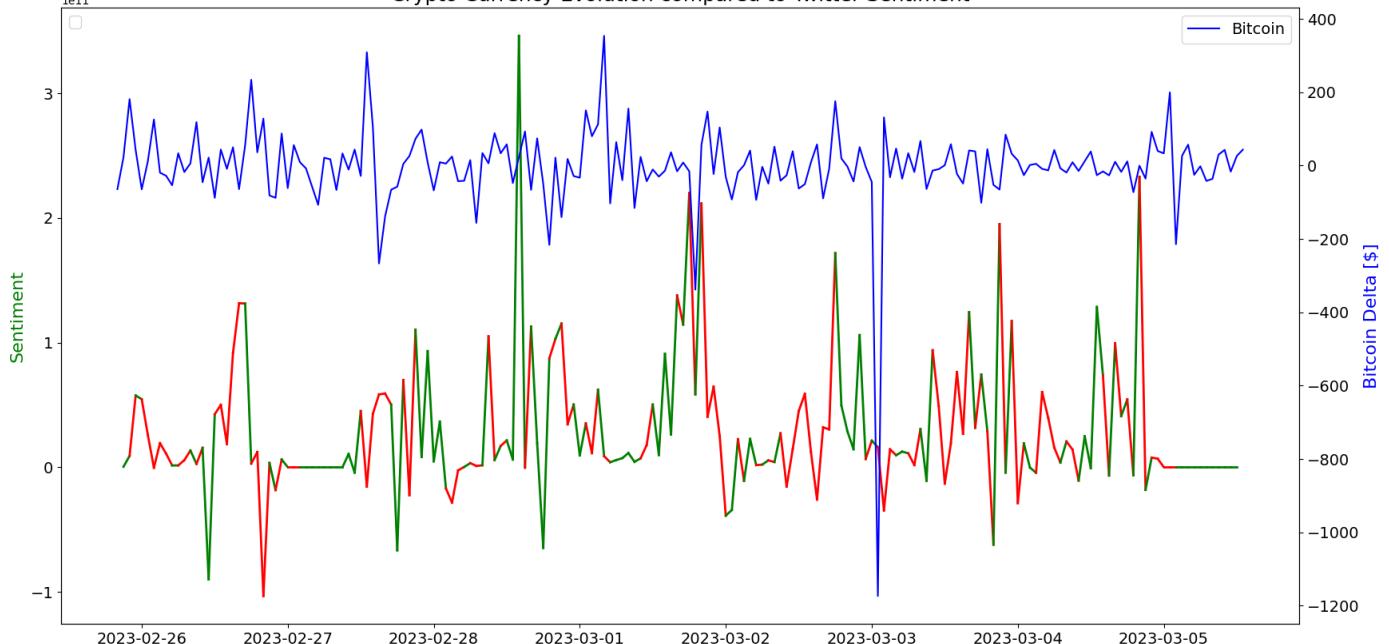
```

```

9 # Set title and axis labels
10 ax1.set_title("Crypto Currency Evolution compared to Twitter Sentiment", fontsize=18)
11 ax1.set_ylabel("Sentiment", color='g', fontsize=16)
12 ax2.set_ylabel("Bitcoin Delta [$]", color='b', fontsize=16)
13
14 # Set tick label size
15 ax1.tick_params(labelsize=14)
16 ax2.tick_params(labelsize=14)
17
18 # Calculate sentiment slope and delta slope
19 sentiment_slope = np.gradient(tweets_grouped)
20 delta_slope = np.gradient(crypto_usd_grouped2)
21
22 # Plot sentiment data
23 for i in range(2, len(tweets_grouped)-1):
24     if np.sign(sentiment_slope[i-2]) == np.sign(delta_slope[i]):
25         ax1.plot_date([tweets_grouped.index[i-1], tweets_grouped.index[i]], [tweets_grouped[i-1], tweets_grouped[i]], 'r-', linewidth=2)
26     else:
27         ax1.plot_date([tweets_grouped.index[i-1], tweets_grouped.index[i]], [tweets_grouped[i-1], tweets_grouped[i]], 'g-', linewidth=2)
28
29 # Plot cryptocurrency data
30 ax2.plot_date(crypto_usd_grouped2.index, crypto_usd_grouped2, 'b-', label="Bitcoin")
31
32 # Format x-axis as dates
33 ax1.xaxis.set_major_locator(mdates.AutoDateLocator())
34 ax1.xaxis.set_major_formatter(mdates.DateFormatter('%Y-%m-%d'))
35
36 # Rotate x-axis tick labels for better visibility
37 plt.xticks(rotation=45)
38
39 # Display legend
40 ax1.legend(loc='upper left', fontsize=14)
41 ax2.legend(loc='upper right', fontsize=14)
42
43 # Show the plot
44 plt.show()
45

```

WARNING:matplotlib.legend:No artists with labels found to put in legend. Note that artists whose label start with an underscore are ignored.



```
1 #To consider both the slope direction and magnitude proportionality
```

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3 import matplotlib.dates as mdates
4

```

```
5 # Create figure and axes
6 fig, ax1 = plt.subplots(figsize=(20, 10))
7 ax2 = ax1.twinx()
8
9 # Set title and axis labels
10 ax1.set_title("Crypto Currency Evolution compared to Twitter Sentiment", fontsize=18)
11 ax1.set_ylabel("Sentiment", color='g', fontsize=16)
12 ax2.set_ylabel("Bitcoin Delta [$]", color='b', fontsize=16)
13
14 # Set tick label size
15 ax1.tick_params(labelsize=14)
16 ax2.tick_params(labelsize=14)
17
18 # Calculate sentiment slope and delta slope
19 sentiment_slope = np.gradient(tweets_grouped)
20 delta_slope = np.gradient(crypto_usd_grouped2)
21
22 # Calculate proportionality factor
23 proportionality = 0.2 # Adjust this value as desired
24
25 # Plot sentiment data
26 for i in range(2, len(tweets_grouped)-1):
27     if np.sign(sentiment_slope[i-2]) == np.sign(delta_slope[i]):
28         if np.abs(sentiment_slope[i-2]) >= proportionality * np.abs(delta_slope[i]):
29             ax1.plot_date([tweets_grouped.index[i-1], tweets_grouped.index[i]], [tweets_grouped[i-1], tweets_grouped[i]], 'r-', linewidth=
30         else:
31             ax1.plot_date([tweets_grouped.index[i-1], tweets_grouped.index[i]], [tweets_grouped[i-1], tweets_grouped[i]], 'g-', linewidth=
32     else:
33         ax1.plot_date([tweets_grouped.index[i-1], tweets_grouped.index[i]], [tweets_grouped[i-1], tweets_grouped[i]], 'g-', linewidth=2)
34
35 # Plot cryptocurrency data
36 ax2.plot_date(crypto_usd_grouped2.index, crypto_usd_grouped2, 'b-', label="Bitcoin")
37
38 # Format x-axis as dates
39 ax1.xaxis.set_major_locator(mdates.AutoDateLocator())
40 ax1.xaxis.set_major_formatter(mdates.DateFormatter('%Y-%m-%d'))
41
42 # Rotate x-axis tick labels for better visibility
43 plt.xticks(rotation=45)
44
45 # Display legend
46 ax1.legend(loc='upper left', fontsize=14)
47 ax2.legend(loc='upper right', fontsize=14)
48
49 # Show the plot
50 plt.show()
51
```