

▼ Dataset- Tweets

1 #Dataset Link: <https://www.kaggle.com/datasets/kaushiksuresh147/bitcoin-tweets>
 2 #Divide the large csv file into smaller chunks for easier storage due to its big size.

```
1 import pandas as pd
2 import os
3
4 file_urls = [
5     'https://github.com/Amarpreet3/CIND-820-CAPSTONE/raw/main/Sentimental%20Analysis/bitcoinTweets_2023-02-25.csv',
6     'https://github.com/Amarpreet3/CIND-820-CAPSTONE/raw/main/Sentimental%20Analysis/bitcoinTweets_2023-02-26.csv',
7     'https://github.com/Amarpreet3/CIND-820-CAPSTONE/raw/main/Sentimental%20Analysis/bitcoinTweets_2023-02-27.csv',
8     'https://github.com/Amarpreet3/CIND-820-CAPSTONE/raw/main/Sentimental%20Analysis/bitcoinTweets_2023-02-28.csv',
9     'https://github.com/Amarpreet3/CIND-820-CAPSTONE/raw/main/Sentimental%20Analysis/bitcoinTweets_2023-03-01.csv',
10    'https://github.com/Amarpreet3/CIND-820-CAPSTONE/raw/main/Sentimental%20Analysis/bitcoinTweets_2023-03-02.csv',
11    'https://github.com/Amarpreet3/CIND-820-CAPSTONE/raw/main/Sentimental%20Analysis/bitcoinTweets_2023-03-03.csv',
12    'https://github.com/Amarpreet3/CIND-820-CAPSTONE/raw/main/Sentimental%20Analysis/bitcoinTweets_2023-03-04.csv',
13    'https://github.com/Amarpreet3/CIND-820-CAPSTONE/raw/main/Sentimental%20Analysis/bitcoinTweets_2023-03-05.csv'
14 ]
15
16 for url in file_urls:
17     # Extract the file name
18     file_name = url.split('/')[-1]
19
20     # Read the CSV file
21     df = pd.read_csv(url)
22
23     # Save the DataFrame as CSV with a different file name
24     file_name_without_extension = file_name.split('.')[0]
25     new_file_name = f"{file_name_without_extension}.csv"
26     df.to_csv(new_file_name, index=False)
27
28     # Check the size of the saved file
29     file_size_mb = os.path.getsize(new_file_name) / (1024 * 1024) # in MB
30     print(f"File {new_file_name} size: {file_size_mb:.2f} MB")
31
32
33 File bitcoinTweets_2023-02-25.csv size: 0.97 MB
34 File bitcoinTweets_2023-02-26.csv size: 9.32 MB
35 File bitcoinTweets_2023-02-27.csv size: 6.40 MB
36 File bitcoinTweets_2023-02-28.csv size: 10.86 MB
37 File bitcoinTweets_2023-03-01.csv size: 13.57 MB
38 File bitcoinTweets_2023-03-02.csv size: 9.28 MB
39 File bitcoinTweets_2023-03-03.csv size: 9.85 MB
40 File bitcoinTweets_2023-03-04.csv size: 9.72 MB
41 File bitcoinTweets_2023-03-05.csv size: 3.69 MB
```

1 #Concatenate the Tweets into one file for Analysis

```
1 import pandas as pd
2
3 file_urls = [
4     'https://github.com/Amarpreet3/CIND-820-CAPSTONE/raw/main/Sentimental%20Analysis/bitcoinTweets_2023-02-25.csv',
5     'https://github.com/Amarpreet3/CIND-820-CAPSTONE/raw/main/Sentimental%20Analysis/bitcoinTweets_2023-02-26.csv',
6     'https://github.com/Amarpreet3/CIND-820-CAPSTONE/raw/main/Sentimental%20Analysis/bitcoinTweets_2023-02-27.csv',
7     'https://github.com/Amarpreet3/CIND-820-CAPSTONE/raw/main/Sentimental%20Analysis/bitcoinTweets_2023-02-28.csv',
8     'https://github.com/Amarpreet3/CIND-820-CAPSTONE/raw/main/Sentimental%20Analysis/bitcoinTweets_2023-03-01.csv',
9     'https://github.com/Amarpreet3/CIND-820-CAPSTONE/raw/main/Sentimental%20Analysis/bitcoinTweets_2023-03-02.csv',
10    'https://github.com/Amarpreet3/CIND-820-CAPSTONE/raw/main/Sentimental%20Analysis/bitcoinTweets_2023-03-03.csv',
11    'https://github.com/Amarpreet3/CIND-820-CAPSTONE/raw/main/Sentimental%20Analysis/bitcoinTweets_2023-03-04.csv',
12    'https://github.com/Amarpreet3/CIND-820-CAPSTONE/raw/main/Sentimental%20Analysis/bitcoinTweets_2023-03-05.csv'
13 ]
14
15 dfs = []
16
17 for url in file_urls:
18     # Read the CSV file
19     df = pd.read_csv(url)
20
21     # Append the DataFrame to the list
22     dfs.append(df)
23
```

```
24 # Combine all DataFrames into a single DataFrame
25 combined_df = pd.concat(dfs)
26
27 # Save the combined DataFrame as a CSV file
28 combined_df.to_csv('bitcoinTweets.csv', index=False)
29

1 bitcoinTweets= pd.read_csv("/content/bitcoinTweets.csv")
2 bitcoinTweets.head()

1 df = bitcoinTweets.copy()

1 df.head()

1 df.info()

1 print(df.columns)

1 import pandas as pd
2 import numpy as np
3 import datetime as dt
4
5 # Convert 'date' column to datetime64[ns]
6 df['date'] = pd.to_datetime(df['date'], format='%Y-%m-%d %H:%M:%S')
7
8 # Display the updated dataset
9 df.info()
10
11
12

1 # Check the size of the dataset
2 num_rows, num_cols = df.shape
3 print("Number of rows:", num_rows)
4 print("Number of columns:", num_cols)

1 import pandas as pd
2 import sys
3 # Calculate the size of the dataset in megabytes
4 dataset_size_mb = sys.getsizeof(df) / (1024 * 1024)
5 print("Dataset size:", round(dataset_size_mb, 2), "MB")

1 column_names = df.columns
2 print(column_names)

1 # Change the data type of date
2 date_dtype = df['date'].dtypes
3 print(date_dtype)

1 # Convert the 'date' column to a date format with errors set to 'coerce'
2 df['date'] = pd.to_datetime(df['date'], errors='coerce')
3
4 # Save the modified dataset
5 df.to_csv('/content/bitcoinTweets_updated.csv', index=False)

1 df = pd.read_csv('/content/bitcoinTweets_updated.csv')
2 df['date'] = pd.to_datetime(df['date'], errors='coerce')
3 date_dtype = df['date'].dtypes
4 print(date_dtype)

1 column_names = df.columns
2 print(column_names)

1 df.head()
```

▼ Bitcoin price

```

1 #Cryptocompare Website
2 #Dataset Link # https://min-api.cryptocompare.com/data/v2/histohour

1 import requests
2 import pandas as pd
3
4 # API endpoint for CryptoCompare
5 url = "https://min-api.cryptocompare.com/data/v2/histohour"
6
7 # Parameters for API request
8 coin_symbol = "BTC"
9 currency = "USD"
10 timestamp_start = 1677110400 # February 22, 2023, 00:00:00 in Unix timestamp
11 timestamp_end = 1678022399 # March 10, 2023, 23:59:59 in Unix timestamp
12
13 # API request payload
14 payload = {
15     "fsym": coin_symbol,
16     "tsym": currency,
17     "toTs": timestamp_end,
18     "limit": 336, # Number of hours between the two timestamps
19     "aggregate": 1,
20     "e": "CCCAGG" # CryptoCompare aggregate index
21 }
22
23 # Send API request
24 response = requests.get(url, params=payload)
25
26 # Check if the request was successful
27 if response.status_code == 200:
28     # Process the response data
29     data = response.json()
30     # The 'Data' field contains the hourly Bitcoin prices
31     bitcoin_prices = data['Data']['Data']
32
33     # Create a DataFrame from the Bitcoin price data
34     df = pd.DataFrame(bitcoin_prices, columns=["time", "close", "high", "low", "open", "volumefrom", "volumeto"])
35     df["time"] = pd.to_datetime(df["time"], unit="s") # Convert Unix timestamp to datetime
36
37     # Save DataFrame to CSV file
38     df.to_csv("bitcoin_prices_2023-02-22_to_2023-03-10.csv", index=False)
39     print("Bitcoin data from February 22 to March 10, 2023, saved to 'bitcoin_prices_2023-02-22_to_2023-03-10.csv' file.")
40 else:
41     print("Failed to fetch Bitcoin data.")
42

```

Bitcoin data from February 22 to March 10, 2023, saved to 'bitcoin_prices_2023-02-22_to_2023-03-10.csv' file.

```
1 df_bitcoin = pd.read_csv("/content/bitcoin_prices_2023-02-22_to_2023-03-10.csv")
```

```

1 df_price = df_bitcoin.copy()
2 df_price['time'] = pd.to_datetime(df_price['time'])
3
4 df_price['Date'] = df_price['time'].dt.date
5 df_price['Time'] = df_price['time'].dt.time
6
7 df_price.head(2)
8

```

	time	close	high	low	open	volumefrom	volumeto	Date	Time
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

```

1
2 df_price['Date'] = pd.to_datetime(df_price['Date'], format='%Y-%m-%d').dt.date
3

```

```
4 df_price.info()
5
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 337 entries, 0 to 336
Data columns (total 9 columns):
#   Column      Non-Null Count  Dtype
---  ---
0    time        337 non-null    datetime64[ns]
1    close        337 non-null    float64
2    high         337 non-null    float64
3    low          337 non-null    float64
4    open         337 non-null    float64
5    volumefrom   337 non-null    float64
6    volumeto     337 non-null    float64
7    Date         337 non-null    object
8    Time         337 non-null    object
dtypes: datetime64[ns](1), float64(6), object(2)
memory usage: 23.8+ KB
```

```
1 crypto_usd = df_price.copy()
2 # Drop duplicates Currency
3 print('bitcoin shape before dropping duplicates', crypto_usd.shape)
4 duplicates_removed = crypto_usd.shape[0]
5 crypto_usd = crypto_usd.drop_duplicates(subset=['time'])
6 print('bitcoin shape after dropping duplicates', crypto_usd.shape)
7 duplicates_removed -= crypto_usd.shape[0]
8 print('duplicates removed', duplicates_removed)
9 crypto_usd.head(2)
```

```
bitcoin shape before dropping duplicates (337, 9)
bitcoin shape after dropping duplicates (337, 9)
duplicates removed 0
```

	time	close	high	low	open	volumefrom	volumeto	Date	Time
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

```
1 # Create the new 'volume' column
2 crypto_usd['volume'] = crypto_usd['volumeto'] - crypto_usd['volumefrom']
3
4 # Display the updated DataFrame
5 crypto_usd.head(2)
```

	time	close	high	low	open	volumefrom	volumeto	Date	Time
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

```
1 # Sort the dataframe by the 'time' column
2 crypto_usd = crypto_usd.sort_values('time')
3
4 # Calculate market cap
5 crypto_usd['marketcap'] = crypto_usd['close'] * crypto_usd['volumeto']
6
7 # Calculate price difference delta
8 crypto_usd['price_delta'] = crypto_usd['close'].diff()
9
10 # Display the updated dataframe
11 print(crypto_usd.head())
12
```

	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 #Saved the Data with all additional attributes so it is easier to extract for analysis
```

```
1 import pandas as pd
2
3 # Save the DataFrame as CSV
4 crypto_usd.to_csv('BitcoinPrice.csv', index=False)
5
6 # Print the file name
7 print("Data saved as BitcoinPricePreprocessed.csv")
8
```

```
Data saved as BitcoinPricePreprocessed.csv
```

▼ Tweets Cleaning

```
1 import pandas as pd
2
3 # Read the CSV file
4 tweets_raw_file = pd.read_csv('/content/bitcoinTweets_updated.csv')
5
6 # Convert 'date' column to datetime
7 tweets_raw_file['date'] = pd.to_datetime(tweets_raw_file['date'], errors='coerce')
8
9 # Check the data type of 'date' column
10 date_dtype = tweets_raw_file['date'].dtypes
11
12 # Print the data type
13 print(date_dtype)
14
```

```
datetime64[ns]
```

```
1 tweets_raw_file.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 169761 entries, 0 to 169760
Data columns (total 13 columns):
#   Column                Non-Null Count  Dtype
---  -
0   user_name              169751 non-null object
1   user_location          84090 non-null object
2   user_description       158666 non-null object
3   user_created           169761 non-null object
4   user_followers         169761 non-null float64
5   user_friends           169761 non-null float64
6   user_favourites        169761 non-null float64
7   user_verified          169761 non-null bool
8   date                   169761 non-null datetime64[ns]
9   text                   169761 non-null object
10  hashtags                168954 non-null object
11  source                  168954 non-null object
12  is_retweet              168954 non-null float64
dtypes: bool(1), datetime64[ns](1), float64(4), object(7)
memory usage: 15.7+ MB
```

```
1 import pandas as pd
2 import re
3 import string
4 from tqdm import tqdm
5 from nltk.corpus import stopwords
6 from nltk.stem import PorterStemmer
7 import nltk
8 nltk.download('stopwords')
9
10
11 tweets_raw_file = "/content/bitcoinTweets_updated.csv"
12 tweets_clean_file = "bitcoinTweets_cleaned.csv"
13
14 # Load the tweets DataFrame
15 tweets_df = pd.read_csv(tweets_raw_file)
16
17 # Preprocess the 'text' column
18 for i, text in tqdm(enumerate(tweets_df['text']), total=len(tweets_df['text'])):
```

```

19 # Remove hashtags
20 text = text.replace("#", "")
21
22 # Remove URLs
23 text = re.sub('https?:\/\/(?:[-\w.]|(?:%[da-fA-F]{2}))+', '', text)
24
25 # Remove mentions
26 text = re.sub('@\w+', '', text)
27
28 # Convert to lowercase
29 text = text.lower()
30
31 # Remove punctuation
32 text = text.translate(str.maketrans('', '', string.punctuation))
33
34 # Tokenize the text
35 tokens = text.split()
36
37 # Remove stopwords
38 stop_words = set(stopwords.words('english'))
39 tokens = [word for word in tokens if word not in stop_words]
40
41 # Apply stemming
42 stemmer = PorterStemmer()
43 tokens = [stemmer.stem(word) for word in tokens]
44
45 # Join the tokens back into a single string
46 preprocessed_text = ' '.join(tokens)
47
48 # Update the 'text' column with the preprocessed text
49 tweets_df.loc[i, 'text'] = preprocessed_text
50
51 # Save the cleaned data to a new CSV file
52 tweets_df.to_csv(tweets_clean_file, header=True, encoding='utf-8', index=False)
53
54 print("Cleaned tweets data saved to 'bitcoinTweets_cleaned.csv'.")
55

```

```

[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data] Unzipping corpora/stopwords.zip.
100%|██████████| 169761/169761 [11:21<00:00, 249.02it/s]
Cleaned tweets data saved to 'bitcoinTweets_cleaned.csv'.

```

```

1 import pandas as pd
2
3 tweets_clean_file = "bitcoinTweets_cleaned.csv"
4
5 # Load the cleaned tweets DataFrame
6 tweets_df = pd.read_csv(tweets_clean_file)
7 print('tweet shape before dropping duplicates', tweets_df.shape)
8 # Remove duplicate rows based on all columns
9 tweets_df.drop_duplicates(inplace=True)
10
11 # Save the deduplicated data to a new CSV file
12 tweets_df.to_csv(tweets_clean_file, header=True, encoding='utf-8', index=False)
13
14 print("Duplicates removed from the cleaned tweets data.")
15 print('tweet shape after dropping duplicates', tweets_df.shape)
16

```

```

tweet shape before dropping duplicates (169761, 13)
Duplicates removed from the cleaned tweets data.
tweet shape after dropping duplicates (168685, 13)

```

```

1 import sys
2
3 # Get the size of the DataFrame in bytes
4 size_bytes = sys.getsizeof(tweets_df)
5
6 # Convert the size to megabytes
7 size_mb = size_bytes / (1024 * 1024)
8
9 # Print the size in MB
10 print(f"Size of tweets_df: {size_mb} MB")
11

```

```
Size of tweets_df: 137.04691982269287 MB

1 print(tweets_df.shape)

(168685, 13)

1 tweets_df.head()

  user_name user_location user_description user_created user_followers user_friends u:
0      Lrk  Vancouver, WA  Lrk started investing in the stock market in 1...  2018-08-11 03:17:00      116.0      8.0
1    Xiang Zhang      NaN  Professional Software Engineer  2011-01-11 01:37:00      42.0      22.0
2    Rhizoo      NaN  researcher. local maxima dunningâkruger spec...  2019-04-03 18:09:00      778.0      627.0
3    Hari Marquez  Las Vegas, NV  Donât trust, verify. #Bitcoin | El Salvador ...  2014-01-17 23:04:00      222.0      521.0
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 tweets_df.min(axis=0)

<ipython-input-33-14ceec5556a9>:1: FutureWarning: The default value of numeric_only in DataFrame.min is deprecated. In a future version
  tweets_df.min(axis=0)
user_created      2006-09-04 17:48:00
user_followers      0.0
user_friends      0.0
user_favourites      0.0
user_verified      False
date      2023-02-25 20:49:00
is_retweet      0.0
dtype: object

1 tweets_df.info()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 168685 entries, 0 to 169760
Data columns (total 13 columns):
#   Column                Non-Null Count  Dtype
---  -
0   user_name              168675 non-null object
1   user_location          83730 non-null object
2   user_description        158165 non-null object
3   user_created            168685 non-null object
4   user_followers          168685 non-null float64
5   user_friends            168685 non-null float64
6   user_favourites          168685 non-null float64
7   user_verified           168685 non-null bool
8   date                   168685 non-null object
9   text                    168673 non-null object
10  hashtags                167953 non-null object
11  source                  167953 non-null object
12  is_retweet              167953 non-null float64
dtypes: bool(1), float64(4), object(8)
memory usage: 16.9+ MB

1 #Filter Tweets and hashtag for Bitcoin only
```

```

1 import pandas as pd
2 import numpy as np
3
4 # Filter tweets containing specific keywords or hashtags related to Bitcoin
5 bitcoin_keywords = ['bitcoin', 'btc', 'crypto']
6 bitcoin_hashtags = ['#bitcoin', '#btc', '#cryptocurrency']
7
8 # Drop rows with missing values in the 'text' column
9 tweets_df.dropna(subset=['text'], inplace=True)
10
11 # Apply the filtering condition
12 filtered_tweets_df = tweets_df[tweets_df['text'].str.contains('|'.join(bitcoin_keywords + bitcoin_hashtags), case=False)]
13
14 # Print the filtered tweets DataFrame
15 print(filtered_tweets_df)
16

```

		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	
...		
169756	L3X - 00000000000000000000000000000000 v5.3		NaN	
169757		FunFacts.AI	NaN	
169758		Word On Crypto0000	SocialMedia	
169759		BTC Status Alert 0000	Japan	
169760		Jon Padilha	NaN	

		user_description	
0		Irk started investing in the stock market in 1...	
1		Professional Software Engineer 00000000Crypto ...	
2		researcher. local maxima dunning&0000kruger spec...	
3		Donâ0000t trust, verify. #Bitcoin El Salvador ...	
4		Robot that posts the closure of the bitcoin da...	
...		...	
169756		#Bitcoin & #Crypto #Trading #Strategy\nFollow ...	
169757		I post a Fun Fact every minute\n\n#OpenAI #Cha...	
169758		0000 UNDERGROUND	
169759		Tweet data that will help you consider #BTC pr...	
169760		Daytrade Institucional\nForex - BTC - BinÃ0000rias	

	user_created	user_followers	user_friends	user_favourites	
0	2018-08-11 03:17:00	116.0	8.0	4580.0	
1	2011-01-11 01:37:00	42.0	22.0	5.0	
2	2019-04-03 18:09:00	778.0	627.0	32005.0	
3	2014-01-17 23:04:00	222.0	521.0	13052.0	
4	2021-01-06 01:36:00	40.0	4.0	1.0	
...	
169756	2022-04-05 22:51:00	182.0	14.0	2978.0	
169757	2022-05-24 12:26:00	24.0	18.0	1.0	
169758	2021-06-18 16:30:00	41.0	15.0	316.0	
169759	2019-07-21 11:28:00	45666.0	2.0	15.0	
169760	2019-09-22 13:54:00	101.0	3.0	35.0	

	user_verified	date	
0	False	2023-02-25 23:59:00	
1	False	2023-02-25 23:59:00	
2	False	2023-02-25 23:59:00	
3	False	2023-02-25 23:59:00	
4	False	2023-02-25 23:59:00	
...	
169756	False	2023-03-05 18:52:00	
169757	False	2023-03-05 18:52:00	
169758	False	2023-03-05 18:52:00	
169759	False	2023-03-05 18:51:00	
169760	False	2023-03-05 18:51:00	

	text	
0	bitcoin btc rest crypto ye bitcoin cryptocurr ...	
1	retriev invest fund current ongoing tidexcoin kic...	
2	bull save monthli thread today good shit bitco...	
3	el salvador shape futur bitcoin membvk32cn	
4	candl day 25022023 close open 2319406 high 232...	

```

1 import sys
2
3 # Get the size of the DataFrame in bytes
4 size_bytes = sys.getsizeof(filtered_tweets_df)
5
6 # Convert the size to megabytes

```



```

7 size_mb = size_bytes / (1024 * 1024)
8
9 # Print the size in MB
10 print(f"Size of tweets_df: {size_mb} MB")
11
    Size of tweets_df: 136.2565097808838 MB

```

▼ Sentiment analysis with Vader

```

1 !!pip install vaderSentiment
2

['Looking in indexes: https://pypi.org/simple, https://us-python.pkg.dev/colab-wheels/public/simple/',
 'Collecting vaderSentiment',
 '  Downloading vaderSentiment-3.3.2-py2.py3-none-any.whl (125 kB)',
 '\x1b[?25l   \x1b[90m-----\x1b[0m \x1b[32m0.0/126.0 kB\x1b[0m \x1b[31m?\x1b[0m eta \x1b[36m-:--\x1b[0m',
 '-:--\x1b[0m',
 '\x1b[2K   \x1b[90m-----\x1b[0m \x1b[32m126.0/126.0 kB\x1b[0m \x1b[31m3.5 MB/s\x1b[0m eta \x1b[36m0:00:00\x1b[0m',
 '\x1b[?25hRequirement already satisfied: requests in /usr/local/lib/python3.10/dist-packages (from vaderSentiment) (2.27.1)',
 'Requirement already satisfied: urllib3<1.27,>=1.21.1 in /usr/local/lib/python3.10/dist-packages (from requests->vaderSentiment) (1.26.16)',
 'Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.10/dist-packages (from requests->vaderSentiment) (2023.5.7)',
 'Requirement already satisfied: charset-normalizer~=2.0.0 in /usr/local/lib/python3.10/dist-packages (from requests->vaderSentiment) (2.0.12)',
 'Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.10/dist-packages (from requests->vaderSentiment) (3.4)',
 'Installing collected packages: vaderSentiment',
 'Successfully installed vaderSentiment-3.3.2']

1 import pandas as pd
2 from vaderSentiment.vaderSentiment import SentimentIntensityAnalyzer
3 from tqdm import tqdm
4
5 # Assuming you have already cleaned the tweets and stored them in the 'text' column of `df_clean`
6 df_clean = filtered_tweets_df.copy()
7
8 # Handle NaN values in the 'text' column
9 df_clean['text'].fillna('', inplace=True)
10
11 # Initialize the SentimentIntensityAnalyzer
12 analyzer = SentimentIntensityAnalyzer()
13
14 # Perform sentiment analysis and store the compound scores
15 compound = []
16 for i, s in tqdm(enumerate(df_clean['text']), total=len(df_clean['text'])):
17     if isinstance(s, str): # Check if the tweet is a string
18         vs = analyzer.polarity_scores(s)
19         compound.append(vs["compound"])
20     else:
21         compound.append(0.0) # Assign a neutral score for non-string tweets
22
23 # Add the compound scores to the dataframe
24 df_clean["compound"] = compound
25
26 # Save the updated dataframe to a new CSV file
27 df_clean.to_csv("bitcoinTweets_sentiment.csv", index=False)
28
29 print("Sentiment analysis completed and data saved to 'bitcoinTweets_sentiment.csv'.")
30

100%|██████████| 167655/167655 [00:29<00:00, 5594.51it/s]
Sentiment analysis completed and data saved to 'bitcoinTweets_sentiment.csv'.

1 df_clean.head()

```

	user_name	user_location	user_description	user_created	user_followers	user_friends	u:
0	Irk	Vancouver, WA	Irk started investing in the stock market in 1...	2018-08-11 03:17:00	116.0	8.0	
1	Xiang Zhang	NaN	Professional Software Engineer ðŸ’€»ðŸ’€Crypto ...	2011-01-11 01:37:00	42.0	22.0	
2	Rhizoo	NaN	researcher. local maxima dunningâŸ’kruger spec...	2019-04-03 18:09:00	778.0	627.0	
3	Hari Marquez	Las Vegas, NV	DonâŸ’t trust, verify. #Bitcoin El Salvador ...	2014-01-17 23:04:00	222.0	521.0	

1 #Calculate a score for each tweet

```

1 # Perform sentiment analysis and calculate the scores
2 scores = []
3 for i, s in tqdm(df_clean.iterrows(), total=df_clean.shape[0]):
4     vs = analyzer.polarity_scores(s["text"])
5     score = vs["compound"] * (s["user_followers"] + 1) * (s["user_favourites"] + 1)
6     scores.append(score)
7
8 # Add the scores to the dataframe
9 df_clean["score"] = scores
10
11 # Display the first two rows of the updated dataframe
12 print(df_clean.head(2))

```

```

100%|██████████| 167655/167655 [00:47<00:00, 3543.34it/s]
0      Irk  Vancouver, WA
1  Xiang Zhang      NaN

      user_description      user_created \
0  Irk started investing in the stock market in 1...  2018-08-11 03:17:00
1  Professional Software Engineer ðŸ’€»ðŸ’€Crypto ...  2011-01-11 01:37:00

      user_followers  user_friends  user_favourites  user_verified \
0              116.0           8.0           4580.0          False
1              42.0          22.0           5.0          False

      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...

      hashtags      source \
0  ['Bitcoin', 'crypto', 'NeedsMoreCrash']  Twitter Web App
1  ['Tidexcoin', 'Kicurrency', 'LMY', 'GMK', 'SYR...  Twitter for iPhone

      is_retweet  compound      score
0              0.0    -0.4019 -215409.1563
1              0.0     0.0000  0.0000

```

1 df_clean.head(2)

	user_name	user_location	user_description	user_created	user_followers	user_friends	u:
0	lrk	Vancouver, WA	lrk started investing in the stock market in 1...	2018-08-11 03:17:00	116.0	8.0	

```

1 import sys
2
3 # Get the size of the DataFrame in bytes
4 size_bytes = sys.getsizeof(df_clean)
5
6 # Convert the size to megabytes
7 size_mb = size_bytes / (1024 * 1024)
8
9 # Print the size in MB
10 print(f"Size of tweets_df: {size_mb} MB")
11
    Size of tweets_df: 138.81472206115723 MB

1 # Check for unique values in the 'user_name' column
2 unique_user_names = df_clean['user_name'].unique()
3
4 # Check if all user names are unique
5 if len(unique_user_names) == len(df_clean):
6     print("All user names are unique.")
7 else:
8     print("There are duplicate user names.")
9 # Check the number of unique values in the 'user_name' column
10 num_unique_user_names = df_clean['user_name'].nunique()
11 print(f"The number of unique user names is: {num_unique_user_names}")

    There are duplicate user names.
    The number of unique user names is: 36054

1 #Split dataframe and save it into multiple files

1 date_column_format = df_clean['date'].dtypes
2
3 print(f"The format of the 'date' column is: {date_column_format}")

    The format of the 'date' column is: object

1 #The purpose of dividing the DataFrame into smaller chunks and then merging them back together is to handle large datasets efficiently.

1
2 import pandas as pd
3 from datetime import datetime
4
5 n = 20000 # chunk row size
6 chunks_df = [df_clean[i:i+n] for i in range(0, df_clean.shape[0], n)]
7
8 sep_char = '~'
9 path = '/content' # Specify the path where you want to save the CSV files
10
11 # Concatenate all the chunked DataFrames into a single DataFrame
12 concatenated_df = pd.concat(chunks_df)
13
14 # Get the minimum and maximum dates from the concatenated DataFrame
15 date_from = datetime.strptime(concatenated_df['date'].min(), '%Y-%m-%d %H:%M:%S').strftime('%Y-%m-%d %H:%M:%S')
16 date_to = datetime.strptime(concatenated_df['date'].max(), '%Y-%m-%d %H:%M:%S').strftime('%Y-%m-%d %H:%M:%S')
17
18 # Write the concatenated DataFrame into CSV
19 concatenated_df.to_csv(f"{path}/Chunk_{date_from}-{sep_char}-{date_to}.csv", header=True, index=False)

1 import sys
2
3 # Get the size of the DataFrame in bytes
4 size_bytes = sys.getsizeof(concatenated_df)
5
6 # Convert the size to megabytes
7 size_mb = size_bytes / (1024 * 1024)

```

```
8
9 # Print the size in MB
10 print(f"Size of tweets_df: {size_mb} MB")
11

Size of tweets_df: 138.81472206115723 MB
```

```
1 concatenated_df.info()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 167655 entries, 0 to 169760
Data columns (total 15 columns):
 #   Column                Non-Null Count  Dtype
---  -
 0   user_name             167645 non-null object
 1   user_location         83351 non-null object
 2   user_description      157326 non-null object
 3   user_created          167655 non-null object
 4   user_followers        167655 non-null float64
 5   user_friends          167655 non-null float64
 6   user_favourites       167655 non-null float64
 7   user_verified         167655 non-null bool
 8   date                  167655 non-null object
 9   text                  167655 non-null object
10   hashtags              167145 non-null object
11   source                167145 non-null object
12   is_retweet            167145 non-null float64
13   compound              167655 non-null float64
14   score                 167655 non-null float64
dtypes: bool(1), float64(6), object(8)
memory usage: 19.3+ MB
```

```
1 concatenated_df.head()
```

	user_name	user_location	user_description	user_created	user_followers	user_friends	u:
0	Irk	Vancouver, WA	Irk started investing in the stock market in 1...	2018-08-11 03:17:00	116.0	8.0	
1	Xiang Zhang	NaN	Professional Software Engineer δ□□»δ□□□Crypto ...	2011-01-11 01:37:00	42.0	22.0	
2	Rhizoo	NaN	researcher. local maxima dunningâ□□kruger spec...	2019-04-03 18:09:00	778.0	627.0	
3	Hari Marquez	Las Vegas, NV	Donâ□□t trust, verify. #Bitcoin El Salvador ...	2014-01-17 23:04:00	222.0	521.0	
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 #Sentiment Level

1 #The compound score represents the overall sentiment intensity of the text, taking into account the positive, negative, and neutral scores

1 import nltk
2 nltk.download('vader_lexicon')
3 from nltk.sentiment import SentimentIntensityAnalyzer
4 import pandas as pd
5 import numpy as np
6
7 # Initialize VADER sentiment analyzer
```

```

8 sia = SentimentIntensityAnalyzer()
9
10 # Calculate compound scores
11 concatenated_df['compound'] = concatenated_df['text'].apply(lambda x: sia.polarity_scores(x)['compound'])
12
13 # Find the range of compound values
14 compound_min = concatenated_df['compound'].min()
15 compound_max = concatenated_df['compound'].max()
16
17 # Define custom bin edges based on quantiles
18 bin_edges = np.linspace(compound_min, compound_max, num=6) # Adjust the 'num' parameter as needed
19
20 # Define labels
21 labels = ['Extreme Negative', 'Negative', 'Neutral', 'Positive', 'Extreme Positive']
22
23 # Assign sentiment levels based on custom bins
24 concatenated_df['sentiment_level'] = pd.cut(concatenated_df['compound'], bins=bin_edges, labels=labels, include_lowest=True)
25
26 # Save the updated dataframe as a new CSV file
27 concatenated_df.to_csv('updated_sentiment_data.csv', index=False)
28

```

[nltk_data] Downloading package vader_lexicon to /root/nltk_data...

```

1 concatenated_df_2 = pd.read_csv("updated_sentiment_data.csv")
2 concatenated_df_2.head()

```

	user_name	user_location	user_description	user_created	user_followers	user_friends	u:
0	Irk	Vancouver, WA	Irk started investing in the stock market in 1...	2018-08-11 03:17:00	116.0	8.0	
1	Xiang Zhang	NaN	Professional Software Engineer ðïï»ðïïCrypto ...	2011-01-11 01:37:00	42.0	22.0	
2	Rhizoo	NaN	researcher. local maxima dunningâïïkruger spec...	2019-04-03 18:09:00	778.0	627.0	
3	Hari Marquez	Las Vegas, NV	Donâïït trust, verify. #Bitcoin El Salvador ...	2014-01-17 23:04:00	222.0	521.0	
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 print(concatenated_df_2.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'],
      dtype='object')

```

```

1 label_counts = concatenated_df_2['sentiment_level'].value_counts()
2 print(label_counts)

```

```

Neutral      93170
Positive     35921
Extreme Positive 17344
Negative      15904
Extreme Negative 5316
Name: sentiment_level, dtype: int64

```

```

1 import sys
2
3 # Get the size of the DataFrame in bytes
4 size_bytes = sys.getsizeof(concatenated_df_2)
5
6 # Convert the size to megabytes
7 size_mb = size_bytes / (1024 * 1024)
8
9 # Print the size in MB
10 print(f"Size of tweets_df: {size_mb} MB")
11

```

Size of tweets_df: 148.0125036239624 MB

▼ Sentiment Analysis with TextBlob

```

1 #Calculate Polarity and Subjectivity

1 import pandas as pd
2 from textblob import TextBlob
3
4
5 polarity = []
6 subjectivity = []
7
8 # Perform sentiment analysis on each tweet
9 for tweet in concatenated_df_2['text']:
10     try:
11         analysis = TextBlob(tweet)
12         polarity.append(analysis.sentiment.polarity)
13         subjectivity.append(analysis.sentiment.subjectivity)
14     except:
15         polarity.append(0)
16         subjectivity.append(0)
17
18 # Add sentiment polarity and subjectivity columns to the dataframe
19 concatenated_df_2['polarity'] = polarity
20 concatenated_df_2['subjectivity'] = subjectivity
21
22 # Display the updated dataframe
23 print(concatenated_df_2.head())
24

```

	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 the stock market in 1...	2018-08-11 03:17:00
1	Professional Software Engineer 000000Crypto ...	2011-01-11 01:37:00
2	researcher. local maxima dunningâ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 membv32cn
4	2023-02-25 23:59:00	candl 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 concatenated_df_2.to_csv('BitcoinPriceTweets.csv', index=False)
2

```

▼ Data For Further Analysis

```

1 import pandas as pd
2
3 # Assuming you have a DataFrame named crypto_usd
4
5 # Save the DataFrame as CSV
6 crypto_usd = pd.read_csv('BitcoinPrice.csv')
7
8

```

```

1 #Crypto - Bitcoin
2 crypto_usd.head(2)
3

```

	time	close	high	low	open	volumefrom	volumeto	Date	Time
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 22

```

1 crypto_usd.info()

```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 337 entries, 0 to 336
Data columns (total 12 columns):
#   Column      Non-Null Count  Dtype
---  -
0   time        337 non-null   object
1   close       337 non-null   float64
2   high       337 non-null   float64
3   low        337 non-null   float64
4   open       337 non-null   float64
5   volumefrom  337 non-null   float64
6   volumeto   337 non-null   float64
7   Date       337 non-null   object
8   Time       337 non-null   object
9   volume     337 non-null   float64
10  marketcap  337 non-null   float64
11  price_delta 336 non-null   float64
dtypes: float64(9), object(3)
memory usage: 31.7+ KB

```

```

1 #Tweets-Bitcoin
2 tweets = pd.read_csv('/content/BitcoinPriceTweets.csv')
3 tweets.head(1)

```

user_name	user_location	user_description	user_created	user_followers	user_friends	u:
-----------	---------------	------------------	--------------	----------------	--------------	----

1 #Dividing the Data in smaller chunks and save it as size of over all data is apporx 140 MB

```

2
3
4 # Read the dataset from a CSV file
5 dataset = tweets
6 # Calculate the desired size of each subset in bytes
7 desired_size_per_subset = 23 * 1024 * 1024 # Convert 23 MB to bytes
8
9 # Calculate the total number of subsets needed
10 total_subsets = int(round(140 / 23)) # Round up to the nearest integer
11
12 # Calculate the number of rows per subset
13 rows_per_subset = int(round(len(dataset) / total_subsets)) # Round up to the nearest integer
14
15 # Create a directory to save the subsets if it doesn't exist
16 save_directory = '/content'
17 if not os.path.exists(save_directory):
18     os.makedirs(save_directory)
19
20 # Split the dataset into subsets and save each subset as a separate CSV file
21 for subset_index in range(total_subsets):
22     start_index = subset_index * rows_per_subset
23     end_index = (subset_index + 1) * rows_per_subset
24     subset = dataset.iloc[start_index:end_index]
25
26     # Generate the subset file name
27     subset_filename = f"BitcoinTweetsPreprocessed_{subset_index + 1}.csv"
28
29     # Generate the full file path
30     file_path = os.path.join(save_directory, subset_filename)
31
32     # Save the subset as a CSV file
33     subset.to_csv(file_path, index=False)
34
35     # Calculate the size of the subset file
36     subset_size = os.path.getsize(file_path)
37
38     # Get the number of rows in the subset
39     num_rows = len(subset)
40
41     # Print the name, size, and number of rows of each subset file
42     print(f"Subset file: {subset_filename}")
43     print(f"Size: {subset_size / (1024 * 1024)} MB")
44     print(f"Number of rows: {num_rows}")
45     print()
46

```

Subset file: BitcoinTweetsPreprocessed_1.csv
 Size: 11.216121673583984 MB
 Number of rows: 27942

Subset file: BitcoinTweetsPreprocessed_2.csv
 Size: 11.801884651184082 MB
 Number of rows: 27942

Subset file: BitcoinTweetsPreprocessed_3.csv
 Size: 11.73094654083252 MB
 Number of rows: 27942

Subset file: BitcoinTweetsPreprocessed_4.csv
 Size: 11.688798904418945 MB
 Number of rows: 27942

Subset file: BitcoinTweetsPreprocessed_5.csv
 Size: 11.499346733093262 MB
 Number of rows: 27942

Subset file: BitcoinTweetsPreprocessed_6.csv
 Size: 11.43162727355957 MB
 Number of rows: 27942

▼ Datasets for further Anaysis - Preprocessed

```
1 #Saved the processed file o Github and extracting again for 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	volume	from \
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 the stock market in 1...	2018-08-11 03:17:00
1	Professional Software Engineer 00000000Crypto ...	2011-01-11 01:37:00
2	researcher. local maxima dunning&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 membv32cn
4  2023-02-25 23:59:00 candl 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()
```

	user_name	user_location	user_description	user_created	user_followers	user_friends	u:
0	Irk	Vancouver, WA	Irk started investing in the stock market in 1...	2018-08-11 03:17:00	116.0	8.0	
1	Xiang Zhang	NaN	Professional Software Engineer ðŸŒŠ»ðŸŒŠCrypto ...	2011-01-11 01:37:00	42.0	22.0	
2	Rhizoo	NaN	researcher. local maxima dunningâŸŒkruger spec...	2019-04-03 18:09:00	778.0	627.0	
3	Hari Marquez	Las Vegas, NV	DonâŸŒt trust, verify. #Bitcoin El Salvador ...	2014-01-17 23:04:00	222.0	521.0	
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 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
```

