## This notebook uses BERTweet to compute the daily sentiment of non-commercial, non-spam accounts in the MIT Vape and JUUL tweet Dataset from 2015 to 2022

Each tweet is weighted by 1/(the number of tweets the user made that day), and the daily proportion of Negative, Neutral, and Positive sentiment are then computed and plotted by date

### Pulling dataset from Kaggle

```
 1 import zipfile
 2 import os
 3 from google.colab import files
 4
 5 files.upload()
 6
 7 !mkdir ~/.kaggle
 8 !cp kaggle.json ~/.kaggle/
 9 ! chmod 600 ~/.kaggle/kaggle.json
10 !kaggle datasets download -d bestgirl/vape-and-juul-tweet-dataset
11
12
13 with zipfile.ZipFile('vape-and-juul-tweet-dataset.zip', 'r') as zip_ref:
14     zip_ref.extractall('vape_and_juul')
15
16 files = []
17 for dirname, _, filenames in os.walk('vape_and_juul'):
18     for filename in filenames:
19         files.append(os.path.join(dirname, filename))
20
```

Choose Files  kaggle.json
- **kaggle.json**(application/json) - 64 bytes, last modified: 8/28/2025 - 100% done
Saving kaggle.json to kaggle (2).json
mkdir: cannot create directory '/root/.kaggle': File exists
Dataset URL: https://www.kaggle.com/datasets/bestgirl/vape-and-juul-tweet-dataset
License(s): MIT
vape-and-juul-tweet-dataset.zip: Skipping, found more recently modified local copy (use --force to force download)

### Below is a duplicate file being removed and justification for why it is a duplicate

```
1 xlsx_files = [f for f in files if f.lower().endswith(".xlsx")]
2 xlsx_files # this caused the program to crash earlier, I have a csv version in there as well so I am removing it from files since it is
```

['vape_and_juul/2020-08/OTS-2020-08-04.xlsx']

```
1 for f in files:
2     if f[-14:-4] == '2020-08-04':
3         print(f)
```

vape_and_juul/2020-08/OTS-2020-08-04.csv

```
1 # removing the duplicate xlsx
2 files.remove(xlsx_files[0])
```

### Only selecting OTS files

```
1 OTS_only = []
2 for file in files:
3   if file[22:25] =='OTS':
4     OTS_only.append(file)
5 files = OTS_only
6
```

## Instantiating BERTweet

```python
1 from transformers import AutoTokenizer, AutoModelForSequenceClassification, pipeline
2
3 model_name = "finiteautomata/bertweet-base-sentiment-analysis"
4
5 tokenizer = AutoTokenizer.from_pretrained(model_name, use_fast = True)
6 BERTweet = AutoModelForSequenceClassification.from_pretrained(model_name)
7
8 sentiment_pipeline = pipeline(
9     "sentiment-analysis",
10     model = BERTweet,
11     tokenizer = tokenizer,
12     device = 0
13 )
```

```
/usr/local/lib/python3.12/dist-packages/huggingface_hub/utils/_auth.py:94: UserWarning:
The secret `HF_TOKEN` does not exist in your Colab secrets.
To authenticate with the Hugging Face Hub, create a token in your settings tab (https://huggingface.co/settings/tokens), set it as secre
You will be able to reuse this secret in all of your notebooks.
Please note that authentication is recommended but still optional to access public models or datasets.
  warnings.warn(
emoji is not installed, thus not converting emoticons or emojis into text. Install emoji: pip3 install emoji==0.6.0
Device set to use cuda:0
```

## Computing Sentiment Proportions by Date

```python
1 from tqdm import tqdm
2 import pandas as pd
3
4 tweets_today = []
5 dates = []
6 NEGS = []
7 NEUS = []
8 POSS = []
9 files_with_issues = []
10 empty_after_spam = []
11 for file in tqdm(files, desc = "Processing files"):
12
13     try:
14         # Taking in new file and creating df
15         df = pd.read_csv(file, on_bad_lines='skip',engine = 'python', encoding = 'utf-8')
16     except Exception:
17         print(f"there is an error with file: {file}, do it by itself down below")
18         files_with_issues.append(file)
19         continue
20     df = df.dropna()
21
22     # Removing potential retweets and stores
23     df = df[~df['TweetContent'].str.contains("Repost|RT|% off|buy|coupon|sale|discount|free|shop|store|deal|get_repost|get yours today|g
24
25     # Making a dictionary of tweet counts and usernames
26     tweet_counts = df.groupby('TweeterUsername').size()
27     tweet_counts = tweet_counts.to_dict()
28
29     # Spam handling
30     # Since I dont have a bot detecter, I am going to set a hard rule that more than 20 posts in a day is spam and they are either a bot
31     spam_users = set()
32     for k, v in tweet_counts.items():
33         if v > 20:
34             spam_users.add(k)
35
36     df = df[~df['TweeterUsername'].isin(spam_users)]
37
38     # Some more spam handling / ensuring no issues with BERTweet token capacity
39
40     if df['TweetContent'].tolist() == []:
41         # this happend with one of them. Not sure why but will check out in the individual basis ones
42         empty_after_spam.append(file)
43
44         continue
45
46     encodings = tokenizer(df['TweetContent'].tolist(), add_special_tokens = True, truncation = False, padding = False)
```

```
47     df['token_length'] = [len(ids) for ids in encodings['input_ids']]
48
49     df = df[df['token_length'] <=128]
50
51
52    # Running BERTweet on all the remaining valid tweets
53     input = df['TweetContent'].to_list()
54     results = sentiment_pipeline(input, batch_size = 128)
55
56
57     # Weighting sentiment per post based on number of posts made. Each user gets an allocation of 1 unit
58     sentiment_weight = {}
59     for k, v in tweet_counts.items():
60         if k not in spam_users:
61             sentiment_weight[k] = float(1/v)
62
63     # Adding Sentiment and Sentiment weight to the df, the weight is determined by how many total posts they made
64     df['Sentiment'] = [r['label'] for r in results]
65     df['Sentiment_weight'] = df['TweeterUsername'].map(sentiment_weight)
66
67
68     # Computing daily tweet count weighted proportions
69     t = len(set(df['TweeterUsername']))
70
71     df_POS = df[df['Sentiment'] == 'POS']
72     POS_prop = df_POS['Sentiment_weight'].sum() / t
73
74     df_NEU = df[df['Sentiment'] == 'NEU']
75     NEU_prop = df_NEU['Sentiment_weight'].sum() / t
76
77     df_NEG = df[df['Sentiment'] == 'NEG']
78     NEG_prop = df_NEG['Sentiment_weight'].sum() / t
79
80     # Adding Proportions to respective lists
81     NEGS.append(NEG_prop)
82     NEUS.append(NEU_prop)
83     POSS.append(POS_prop)
84     tweets_today.append(len(df))
85
86     # Adding the date to dates
87     date = pd.to_datetime(file[-14:-4])
88     dates.append(date)
89
90
```

```
Processing files: 100%|██████████| 2891/2891 [2:31:46<00:00,  3.15s/it]
```

```
1 files_with_issues
```

```
[]
```

```
 1
 2 empty_after_spam
```

```
['vape_and_juul/2022-11/OTS-2022-11-29.csv',
 'vape_and_juul/2022-11/OTS-2022-11-30.csv']
```

## ⌄ Creating the Data Frame

```
 1 master_df = pd.DataFrame({
 2     'date' : dates,
 3     'NEG_prop': NEGS,
 4     'NEU_prop': NEUS,
 5     'POS_prop': POSS,
 6     'tweets_today': tweets_today
 7 })
 8 master_df['date'] = pd.to_datetime(master_df['date'])
 9 master_df = master_df.sort_values(by = 'date')
10 master_df
```

|  | date | NEG_prop | NEU_prop | POS_prop | tweets_today |
|---|---|---|---|---|---|
| **520** | 2015-01-01 | 0.139855 | 0.437129 | 0.423016 | 1187 |
| **503** | 2015-01-02 | 0.143684 | 0.500984 | 0.355332 | 1374 |
| **506** | 2015-01-03 | 0.202525 | 0.455398 | 0.342078 | 1346 |
| **498** | 2015-01-04 | 0.184096 | 0.472082 | 0.343821 | 1354 |
| **497** | 2015-01-05 | 0.183777 | 0.480495 | 0.335728 | 1652 |
| **...** | ... | ... | ... | ... | ... |
| **2254** | 2022-11-24 | 0.402320 | 0.417159 | 0.180493 | 2522 |
| **2265** | 2022-11-25 | 0.383537 | 0.422500 | 0.193963 | 2220 |
| **2274** | 2022-11-26 | 0.443789 | 0.403812 | 0.152399 | 2208 |
| **2275** | 2022-11-27 | 0.414121 | 0.436821 | 0.149058 | 2064 |
| **2280** | 2022-11-28 | 0.395215 | 0.440867 | 0.162505 | 408 |

2889 rows × 5 columns

Next steps: ( Generate code with `master_df` ) ( ⬤ View recommended plots ) ( New interactive sheet )

## ⌄ Saving master_df to a file in case anything goes wrong:

```
1 from google.colab import files
2 master_df.to_csv("MIT_Vape_and_JUUL_sentiment.csv", index=False)
3 files.download("MIT_Vape_and_JUUL_sentiment.csv")
```

## Inspecting: 'vape_and_juul/2022-11/OTS-2022-11-29.csv', and 'vape_and_juul/2022-11/OTS-2022-11-30.csv'

The actual files are empty and these are the very last date so they will just be ommited from the set.

## ⌄ Plotting Sentiment Proportions

```
 1 import matplotlib.pyplot as plt
 2
 3
 4 plt.figure(figsize=(12, 6))
 5
 6 plt.plot(master_df['date'], master_df['NEG_prop'], label='NEG_prop', color = 'red')
 7 plt.plot(master_df['date'], master_df['NEU_prop'], label='NEU_prop',  color = 'grey')
 8 plt.plot(master_df['date'], master_df['POS_prop'], label='POS_prop',  color = 'green')
 9
10 plt.title("Sentiment Proportions Over Time")
11 plt.xlabel("Date")
12 plt.ylabel("Proportion")
13 plt.legend()
14 plt.grid(True)
15 plt.show()
16
17 plt.figure(figsize=(12, 6))
18
19 plt.plot(master_df['date'], master_df['tweets_today'], label='tweets_today')
20
21 plt.title("Number of Tweets Over Time")
22 plt.xlabel("Date")
23 plt.ylabel("Number of Tweets")
24 plt.legend()
25 plt.grid(True)
26 plt.show()
```

Sentiment Proportions Over Time

Number of Tweets Over Time