

Data Collection

Install Tweepy

To build the dataset for the **'Twitter Sentiment Analysis'** problem, we must first extract the tweets relevant to our topic – **Expo 2020**.

For this, we would be using the built-in python library **'Tweepy'** to extract the tweets related to Expo 2020.

Since Google Colab installs Version 3.10.0 of Tweepy by default, we will be downloading Version 4.4.0 instead to access the latest functions and features offered by Tweepy and ensure that we do not run into any problems while extracting the tweets.

A link to Tweepy's documentation - <https://docs.tweepy.org/en/stable/>
(<https://docs.tweepy.org/en/stable/>).

```
In [ ]: 1 # pip install tweepy==4.4.0
```

Imports

```
In [1]: 1 import tweepy
        2 import configparser
        3 import pandas as pd
        4 import os
        5 from os import access
        6 from re import search
```

Accessing the API

For extracting the tweets using the Twitter API, we need to authenticate and provide access to the Twitter account. This can be done using the keys and tokens provided by Twitter.

We will be storing the keys in a local config file since they should not be shared with anyone but the developer. From this config file, we will be accessing our keys, which can be used to authenticate the Twitter account.

```
In [3]: 1 config = configparser.ConfigParser()           #(Spectrum, 2021)
        2 config.read('../config.ini')
```

```
Out[3]: ['../config.ini']
```

```
In [2]: 1 consumerKey = config['twitter']['apiKey']
        2 consumerSecret = config['twitter']['apiKeySecret']
        3 accessToken = config['twitter']['accessToken']
        4 accessTokenSecret = config['twitter']['accessTokenSecret']
```

Authenticating the Twitter Account

Using the functions offered by Tweepy and the keys provided by Twitter, we establish a connection with the Twitter API and provide access to our account.

After the connection and authentication are successful, we can begin collecting tweets related to Expo 2020.

To collect a large enough number of tweets, we needed to bypass the limit rate on the number of tweets allowed to be retrieved back. By default, the rate limit when searching for tweets matching a specific query is 450 tweets per 15 minutes (Twitter 2022b).

Thus, the `wait_on_rate_limit` has been set to `True`. This will allow the API initialized to wait automatically when a rate limit occur for a time period of 15 minutes (Tweenv 2022)

```
In [4]: 1 auth = tweepy.OAuthHandler(consumerKey, consumerSecret)
2 auth.set_access_token(accessToken, accessSecret)
3 api = tweepy.API(auth, wait_on_rate_limit=True) #provides access to t
```

Initial collector

```
In [ ]: 1 #themes = ["Innovation", "Entertainment", "Sustainability", "Golden Jubil
2 #themes = ["Entertainment"]
3
4 #termToSearch = "Expo 2020"
5 numOfTweets = 1000
6 engLanguage = "en"
7 expoTweets = []
8
9 #for theme in themes:
10 #termToSearch = '#Expo2020+{}'.format(theme),
11 termToSearch = '(#Expo2020) AND (Entertainment OR Sustainability OR "Go
12 #termToSearch = "#Expo2020"
13 tweetCursor = tweepy.Cursor(
14     api.search_tweets,
15     #plenty of the duplicates come from the RT, hence we filter them ou
16     q='{} -filter:retweets'.format(termToSearch),  #(Morrissey, Wasser
17     lang = engLanguage,
18     tweet_mode = "extended",  #(manoelhortaribei, 2017)
19     result_type = "mixed",
20     count = 100,
21 ).items(numOfTweets)
22
23 for tweet in tweetCursor:  #(W3schools.com, 2022)
24     tweetBody = tweet.full_text  #(manoelhortaribei, 2017)
25     expoTweets.append(tweetBody)
26
27 tweetDataFrame = pd.DataFrame(expoTweets, columns=['Tweet Body'])
28 tweetDataFrame = tweetDataFrame.drop_duplicates()  #(Lungu, 2019)
29 tweetDataFrame.to_csv('extractedTweets.csv')
```

Collection of the Tweets

- As there is a possibility of duplicate tweets appearing when collecting the tweets, we decided to use a set. All the elements in a set are unique, as a set cannot contain two elements with the same value. Python's `set()` constructor was used to create the set to store the tweets.

- To collect a wide, diverse collection of tweets, a query was built to include various keywords related to Expo 2020. We started with the term "Expo" and its possible alternate forms that could be present in a tweet, such as a hashtag or a text, and stored all stored of them in the variable `expo2020_str` . Additionally, when people tweet about Expo 2020, sometimes they tag the official Expo 2020 Dubai twitter account (@expo2020dubai). Therefore, we have added `@expo` and `@Expo` as another possible other form of the term "Expo". As we aim to extract tweets about Expo 2020, we have added "2020" to `expo2020_str` to ensure we only get the tweets that are relevant to Expo 2020.
- Using Twitter (2022a) API documentation, the idea behind building the queries was as follows. Since the main keyword to be included in all tweets is `#Expo2020` , the `AND` operation was used to include the optional keywords we want to search for. These words range from `Entertainment` to `Culture` , thus, diversifying the collection of tweets to cover all the topics that were actively discussed in relation to Expo 2020. Any keyword in quotation marks such as "Golden Jubilee" when queried will return an exact match of the keyword used. In this instance, all tweets containing exactly `#Expo` and "Golden Jubilee" will be returned. This was how the query `'(#Expo2020) AND (Entertainment OR Sustainability OR "Golden Jubilee" OR Sports OR Tourism OR Arts OR Culture)'` was created.
- From there, we expanded the queries from that simple concept to include more keywords. We have chosen 43 different keywords, such as "Mobility", "Sustainability", "Opportunity", etc., and mentioned them in alternative forms, inspired by the various ways people can mention a single word in a tweet. E.g. `(opportunity OR Opportunity OR #opportunity OR #Opportunity OR @opportunity OR @Opportunity)` . That gets concatenated with `expo2020_str` to return back specific tweets. According to Twitter (2022a), broad queries are not efficient, hence why we are not collecting only tweets where Expo 2020 appears but are being more specific.
- Retweeted content showed up multiple times in the dataset, since the original tweet was retweeted by multiple users. To avoid these duplicates, all retweets were filtered out. This was done using `-filter:retweets` in the query. Below we can see an instance of how retweets created duplicates within the dataset. The retweet does not differ from the original tweet hence why they were viewed as duplicates and were removed. (Twitter 2022a)

RT @expo2020dubai: Attention! See the preparations for the Union Fortress 8 parade on 4, 11, 18, and 25 February, along Ghaf Avenue and sto... Attention! See the preparations for the Union Fortress 8 parade on 4, 11, 18, and 25 February, along Ghaf Avenue and stopping along Jubilee Park and Al Forsan Park at Expo 2020 Dubai. Get a taste of the parade itself, scheduled for March. Join us! https://t.co/lgx2OUyY1A
--

- Since some of the Expo 2020 tweets may be written in a different language, labeling these tweets by code or manually would not produce accurate results for the sentiment analysis problem. To fix this, we have set the `lang` parameter of `api.search_tweets` to `en` . By setting the language as English, we will be able to collect only tweets posted in English.
- To get the best results for the Twitter Sentiment Analysis problem, we aimed to obtain a mixture of tweets that included the most popular and recent Expo 2020 tweets posted. As a result, we set the `result_type` parameter to `mixed` to get a variety of tweets.
- After getting the tweets, we add them to a set to collect only the unique ones. We noticed that Tweepy had truncated the tweet body to around 140 characters while collecting the tweets. Since we needed the entire tweet body for performing sentiment analysis, we fixed the issue of truncated tweets by setting the `tweet_mode` parameter as `extended` and used the `.full_text` when adding the tweets to the set.

In [8]:

```
1 '''  
2 Initializing a set to store the tweets  
3 '''  
4 expoTweets = set()
```

In [9]:

```
1 themes = ["Innovation","Entertainment","Sustainability", "Golden Jubile
2 termToSearch = "#Expo2020"
3 for theme in themes:
4     termToSearch = '#Expo2020+{}'.format(theme)
5
6 #more on building a query: https://developer.twitter.com/en/docs/twitte
7 # themes = ["Innovation", "Entertainment", "Sustainability", "Golden Ju
8 expo2020_str = "(#Expo OR Expo OR expo OR #expo OR @Expo OR @expo) AND
9 termsToSearch = [
10     "#Expo2020+Innovation",
11     "#Expo2020+Entertainment",
12     "#Expo2020+Sustainability",
13     "#Expo2020+Golden Jubilee",
14     "#Expo2020+Arts and Culture",
15     "#Expo2020+Sports",
16     "#Expo2020+Tourism",
17     '#Expo2020',
18     ' (#Expo2020) AND (Entertainment OR Sustainability OR "Golden Jubile
19     "@expo2020Dubai",
20     ' (#Expo OR Expo OR expo OR #expo OR @Expo OR @expo) AND 2020',
21     expo2020_str + ' AND (Dubai OR dubai OR #dubai OR #Dubai OR @Dubai
22     expo2020_str + ' AND (uae OR UAE OR #uae OR #UAE OR @uae OR @UAE)',
23     expo2020_str + ' AND (entertainment OR Entertainment OR #entertainm
24     expo2020_str + ' AND (sustainability OR Sustainability OR #sustaina
25     expo2020_str + ' AND (golden jubilee OR Golden Jubilee OR #golden j
26     expo2020_str + ' AND (sports OR Sports OR #sports OR #Sports OR @sp
27     expo2020_str + ' AND (Tourism OR tourism OR #Tourism OR #tourism OR
28     expo2020_str + ' AND (art OR Art OR #art OR #Art OR @art OR @Art)',
29     expo2020_str + ' AND (culture OR Culture OR #culture OR #Culture OR
30     expo2020_str + ' AND (heritage OR Heritage OR #heritage OR #Heritag
31     expo2020_str + ' AND (concert OR Concert OR #concert OR #Concert OR
32     expo2020_str + ' AND (festival OR Festival OR #festival OR #Festiva
33     expo2020_str + ' AND (innovation OR Innovation OR #innovation OR #I
34     expo2020_str + ' AND (virtual OR Virtual OR #virtual OR #Virtual OR
35     expo2020_str + ' AND (pavilion OR Pavilion OR #pavilion OR #Pavilio
36     expo2020_str + ' AND (architecture OR Architecture OR #architecture
37     expo2020_str + ' AND (dxb OR Dxb OR #dxb OR #Dxb OR @dxb OR @Dxb OR
38     expo2020_str + ' AND (design OR Design OR #design OR #Design OR @de
39     expo2020_str + ' AND (food OR Food OR #food OR #Food OR @food OR @F
40     expo2020_str + ' AND (technology OR Technology OR #technology OR #T
41     expo2020_str + ' AND (future OR Future OR #future OR #Future OR @fu
42     expo2020_str + ' AND (mobility OR Mobility OR #mobility OR #Mobilit
43     expo2020_str + ' AND (experience OR Experience OR #experience OR #E
44     expo2020_str + ' AND (exhibition OR Exhibition OR #Exhibition OR #e
45     expo2020_str + ' AND (passport OR Passport OR #passport OR #Passpor
46     expo2020_str + ' AND (covid OR Covid OR #Covid OR #Covid OR @covid
47     expo2020_str + ' AND (Cristiano OR cristiano OR #Cristiano OR #cris
48     expo2020_str + ' AND (ronaldo OR Ronaldo OR #ronaldo OR #Ronaldo OR
49     expo2020_str + ' AND (opportunity OR Opportunity OR #opportunity OR
50     expo2020_str + ' AND (forsan OR Forsan OR #forsan OR #Forsan OR @fo
51     expo2020_str + ' AND (opti OR OPTI OR Opti OR #opti OR #OPTI OR #Op
52     expo2020_str + ' AND (emirates OR Emirates OR #emirates OR #Emirate
53     expo2020_str + ' AND (robot OR Robot OR #robot OR #Robot OR @robot
54     expo2020_str + ' AND (sustainable OR Sustainable OR #sustainable OR
55     expo2020_str + ' AND (wasl OR Wasl OR #wasl OR #Wasl OR @wasl OR @W
56     expo2020_str + ' AND (terra OR Terra OR #terra OR #Terra OR @terra
57     expo2020_str + ' AND (music OR Music OR #music OR #Music OR @music
58     expo2020_str + ' AND (water OR Water OR #water OR #Water OR @water
59     expo2020_str + ' AND (digital OR Digital OR #digital OR #Digital OR
```

```

60     expo2020_str + ' AND (climate OR Climate OR #climate OR #Climate OR
61     expo2020_str + ' AND (tourists OR Tourists OR #tourists OR #Tourist
62     expo2020_str + ' AND (country OR Country OR #country OR #Country OR
63     expo2020_str + ' AND (countries OR Countries OR #countries OR #Coun
64 ]
65
66 '''
67 Setting the maximum number of tweets to search as 10000.
68 '''
69 numOfTweets = 10000
70 engLanguage = "en"
71
72 '''
73 Initializing a counter to keep track of the number of tweets we were ab
74 '''
75 tweet_count = 0
76
77 # for termToSearch in termsToSearch:
78 for i in range(len(termsToSearch)):
79     termToSearch = termsToSearch[i]
80     tweetCursor = tweepy.Cursor(
81         api.search_tweets,
82         '''
83         As plenty of duplicates appear from the Retweets (RT), we filt
84         '''
85         q='{} -filter:retweets'.format(termToSearch), #(Morrissey, Was
86         lang = engLanguage,
87         tweet_mode = "extended", #(manoeLhortaribeiRO, 2017)
88         result_type = "mixed",
89         count = 100, '''We try to retrieve 100 tweets pe
90         include_entities = False
91     ).items(numOfTweets)
92
93     for tweet in tweetCursor: #(W3schools, 2022)
94         expoTweets.add(tweet.full_text) #(manoeLhortaribeiRO, 2017)
95         tweet_count += 1 '''Incrementing the counte
96     print(termToSearch)
97     print("no. of tweets retrieved = {}, no. of unique tweets = {}".for
98 # 6613

```

#Expo2020+Innovation

no. of tweets retrieved = 79, no. of unique tweets = 77

#Expo2020+Entertainment

no. of tweets retrieved = 84, no. of unique tweets = 82

#Expo2020+Sustainability

no. of tweets retrieved = 114, no. of unique tweets = 108

#Expo2020+Golden Jubilee

no. of tweets retrieved = 116, no. of unique tweets = 110

#Expo2020+Arts and Culture

no. of tweets retrieved = 117, no. of unique tweets = 111

#Expo2020+Sports

no. of tweets retrieved = 131, no. of unique tweets = 124

#Expo2020+Tourism

no. of tweets retrieved = 154, no. of unique tweets = 146

"#Expo2020

no. of tweets retrieved = 2037, no. of unique tweets = 1884

(#Expo2020) AND (Entertainment OR Sustainability OR "Golden Jubilee" OR Sp
orts OR Tourism OR Arts OR Culture)

no. of tweets retrieved = 2149, no. of unique tweets = 1886

Q=expo2020+Innovation

Saving the tweets collected

Initially, we created a Pandas DataFrame to store the set of unique tweets. We chose to use a DataFrame for the tweets as Pandas offers a function `to_csv()`, which helps to convert the DataFrame into a csv format.

After converting it to a csv format, we can then store the tweets in a .csv file, which will serve as our dataset.

```
In [10]: 1 '''
2 Initializing the DataFrame for the tweets, and setting the column of th
3 '''
4 tweetDataFrame = pd.DataFrame(expoTweets, columns=['Tweet Body'])
5 print(tweetDataFrame.shape)

(6613, 1)
```

```
In [11]: 1 '''
2 Converting the DataFrame into a csv format so that we can store it in a
3 '''
4 tweetDataFrame.to_csv('./data/tweets.csv', index=False)
```

References

L

- Lungu, C. (2019). Exploring Twitter data using Python: An Intro to NLP and Sentiment Analysis. Medium. [online]
Available at: <https://medium.com/analytics-vidhya/exploring-twitter-data-using-python-af1287ee65f1> (<https://medium.com/analytics-vidhya/exploring-twitter-data-using-python-af1287ee65f1>).

M

- Manoelhortaribeiro (2017). Tweepy not getting full text · Issue #935 · tweepy/tweepy. GitHub. [online]
Available at: <https://github.com/tweepy/tweepy/issues/935> (<https://github.com/tweepy/tweepy/issues/935>).
- Morrissey, M., Wasser, L. and Farmer, C., 2020. *Automate Getting Twitter Data in Python Using Tweepy and API Access*. [online] Earth Data Science.
Available at: <https://www.earthdatascience.org/courses/use-data-open-source-python/intro-to-apis/twitter-data-in-python/> (<https://www.earthdatascience.org/courses/use-data-open-source-python/intro-to-apis/twitter-data-in-python/>).

S

- Spectrum, A. (2021). How to get TWEETS by Python | Twitter API 2022. YouTube. [online]
Available at: <https://www.youtube.com/watch?v=Lu1nskBkPJU&t=601s> (<https://www.youtube.com/watch?v=Lu1nskBkPJU&t=601s>).

T

- Tweepy. (2022). *Examples — tweepy 4.5.0 documentation*. [online]
Available at: https://docs.tweepy.org/en/stable/examples.html?highlight=wait_on_rate_limit (https://docs.tweepy.org/en/stable/examples.html?highlight=wait_on_rate_limit)
- Twitter. (2022a). *Search Tweets - How to build a query*. [online]
Available at: <https://developer.twitter.com/en/docs/twitter-api/tweets/search/integrate/build-a-query> (<https://developer.twitter.com/en/docs/twitter-api/tweets/search/integrate/build-a-query>).
- Twitter. (2022b). *Standard search API*. [online]
Available at: <https://developer.twitter.com/en/docs/twitter-api/v1/tweets/search/api-reference/get-search-tweets> (<https://developer.twitter.com/en/docs/twitter-api/v1/tweets/search/api-reference/get-search-tweets>).

W

- W3schools. (2022). *Python Iterators*. [online]
Available at: https://www.w3schools.com/python/python_iterators.asp (https://www.w3schools.com/python/python_iterators.asp)