```python
1 import os
2 import datetime
3 import json
4 import twython
5 import json
```

```python
1 import pandas as pd
2 import numpy as np
```

```python
1 from twython import Twython
2 from contextlib import suppress
3 from requests_oauthlib import OAuth1Session
4 from apscheduler.schedulers.background import BackgroundScheduler as Scheduler
```

```python
 1 # Enter your keys/secrets as strings in the following fields
 2 credentials = {}
 3 credentials['CONSUMER_KEY'] =
 4 credentials['CONSUMER_SECRET'] =
 5 credentials['ACCESS_KEY'] =
 6 credentials['ACCESS_SECRET'] =
 7
 8 # Save the credentials object to file
 9 with open("data/twitter_credentials.json", "w") as file:
10     json.dump(credentials, file)
```

```python
1 # Load credentials from json file
2 with open("data/twitter_credentials.json", "r") as file:
3     creds = json.load(file)
4
5 # Instantiate an object
6 python_tweets = Twython(creds['CONSUMER_KEY'], creds['CONSUMER_SECRET'])
```

```python
1 # getting the header for our dataframe so we have something to append to
2 startdf = pd.read_csv('df')
3 # when our dataframe is created the index is unnamed, when you export
4 # it is renamed so we need to fix it so our df realizes it's the same
5 startdf = startdf.head(0).rename(columns={'Unnamed: 0':''})
6 #assign id as most recent tweet id to  so that we get everything after
7 id_ = python_tweets.search(**{'q': 'RT', 'result_type':
8                                'recent', 'count': 2})['statuses'][0]['id']
```

```python
1 def gather_tweets(q):
2     # import global variables we will be using this function to alter
3     global id_
4     global startdf
5     # Create our query
```

```python
    query = {'q': q,
             'result_type': 'recent',
             'count': 100,
             'lang': 'en',
             'max_id': id_,
             'tweet_mode' : 'extend',
             'entities': {
                 "hashtags": [],
                 "symbols": [],
                 "user_mentions": []
             }
    }
    #build our query into a dictionary to easily turn into dataframe
    dict_ = {'user': [], 'user_id':[], 'post_id': [], 'text': [],
             'favorite_count': [], 'hashtags':[], 'symbols':[],
             'user_mentions':[], 'retweet_count':[]}
    for status in python_tweets.search(**query)['statuses']:
        dict_['user'].append(status['user']['screen_name'])
        dict_['user_id'].append(status['user']['id'])
        dict_['text'].append(status['text'])
        dict_['hashtags'].append(status['entities']['hashtags'])
        dict_['symbols'].append(status['entities']['symbols'])
        dict_['user_mentions'].append(status['entities']['user_mentions'])
        dict_['favorite_count'].append(status['favorite_count'])
        dict_['retweet_count'].append(status['retweet_count'])
        dict_['post_id'].append(status['id'])
    # put data in a DataFrame to work with it easier
    df = pd.DataFrame(dict_)
    #removing our vairables stuck in dictionaries within our
    #dataframe and give them their own columns
    df['mentions'] = ' '
    df['hashtag'] = ' '
    for i in range(len(df)):
        for j in range(len(df.at[i, 'user_mentions'])):
            try:
                df.at[i, 'mentions'] = str(df.at[i, 'mentions'] ,
                    ' ' , df.at[i, 'user_mentions'][j]['id'])
            except:
                df.at[i, 'mentions'] = str(df.at[i, 'mentions'])
        for k in range(len(df.at[i, 'hashtags'])):
            try:
                df.at[i, 'hashtag'] = str(df.at[i, 'hashtag']),
                    ' ' , str(df.at[i, 'hashtags'][k]['text'])
            except:
                df.at[i, 'hashtag'] = str(df.at[i, 'hashtag'])
    # get rid of columns that we took everything we need from
    df = df.drop(['user_mentions', 'hashtags'], 1)
    startdf = startdf.append(df)
    id_ = None
    # get oldest id_ (lowest number) and subtract one so
    # that that you don't reinclude it in the next search reset
    #  id_ to be oldest
```

```python
57      #     id_ to be oldest
58      try:
59          id_ = df['post_id'].sort_values(ascending=True)
60          id_ = int(id_.to_frame().reset_index()['post_id'][0] - 1)
61      except:
62          id_ = python_tweets.search(**{'q': 'RT', 'result_type': 'recent',
63                                      'count': 2})['statuses'][0]['id']
64      return startdf, id_
```

```python
1  # keeping a list of all of our used hashtags
2  #  'astronomy', 'Starship', 'mars', 'curiosityrover',
3  #  'oppertunityrover', 'starlink', 'falconheavy',
4  #  'sls','ESA', 'NASA', 'spacex', 'virgingalactic',
5  #  'virginorbit', 'JAXA', 'Roscosmos', 'artemis',
6  #  'starliner', 'blueorigin', 'spacetravel', 'marswebcam',
7  #  'falcon9', 'nasa_app', 'universe', 'cosmos',
8  #  'iss', 'climate', 'internationalspacestation', 'futurism',
9  #  'starliner', 'blueorigin', 'spacetravel','starlink', 'falconheavy',
10 #  'astronomy', 'mars', 'curiosityrover', 'oppertunityrover',
11 #  'NASA', 'spacex', 'virgingalactic', 'virginorbit', 'JAXA', 'areospace'
```

```python
1 query_list = []
```

```python
1  # eventual goal is to use those loop to do all of our searches
2  end_item = datetime.datetime.now()
3  id_ = python_tweets.search(**{'q': 'RT', 'result_type': 'recent',
4                              'count': 2})['statuses'][0]['id']
5
6  for item in query_list:
7      #reset id_ to be most recent tweet for each item
8      id_ = python_tweets.search(**{'q': 'RT', 'result_type': 'recent',
9                                  'count': 2})['statuses'][0]['id']
10     # for each item, reset the start time to the end of the last item
11     start_time = end_item + datetime.timedelta(seconds=1)
12     # set the end time to be 90 min after you start, gathering 50
13     #  tweets every 100 seconds for 2,700 tweets from each item in list
14     end_item = end_item + datetime.timedelta(seconds=1500)
15     #print start and end times so I know when the program
16     # will be finished and where it should be at
17     print(f'starting {item} time : {start_time.strftime("%H:%M:%S")},'+
18            'ending time: {end_item.strftime("%H:%M:%S")}')
19     # occassionally getting key/value errors of 0 that I cannot find the
20     # cause of, but do not affect anything in how our dataframe
21     # is made, I just don't want them printing
22     with suppress(KeyError, ValueError):
23         #initiate scheduler
24         sch = Scheduler()
25         #add our function and items, we're doing 50/10 seconds
26         sch.add_job(myfn(q), 'interval', (item, startdf), seconds= 5,
27                       start_date=start_time, end_date=end_item)
```

```
28         #start scheduler
29         sch.start()

    starting marswebcam time : 22:21:02, ending time: 22:46:01
    starting falcon9 time : 22:46:02, ending time: 23:11:01
    starting nasa_app time : 23:11:02, ending time: 23:36:01
    starting universe time : 23:36:02, ending time: 00:01:01
    starting cosmos time : 00:01:02, ending time: 00:26:01
    starting iss time : 00:26:02, ending time: 00:51:01
    starting climate time : 00:51:02, ending time: 01:16:01
    starting internationalspacestation time : 01:16:02, ending time: 01:41:01
    starting futurism time : 01:41:02, ending time: 02:06:01
    starting starliner time : 02:06:02, ending time: 02:31:01
    starting blueorigin time : 02:31:02, ending time: 02:56:01
    starting spacetravel time : 02:56:02, ending time: 03:21:01
    starting astronomy time : 03:21:02, ending time: 03:46:01
    starting mars time : 03:46:02, ending time: 04:11:01
    starting curiosityrover time : 04:11:02, ending time: 04:36:01
    starting oppertunityrover time : 04:36:02, ending time: 05:01:01
    starting starlink time : 05:01:02, ending time: 05:26:01
    starting falconheavy time : 05:26:02, ending time: 05:51:01
    starting sls time : 05:51:02, ending time: 06:16:01
    starting ESA time : 06:16:02, ending time: 06:41:01
    starting NASA time : 06:41:02, ending time: 07:06:01
    starting spacex time : 07:06:02, ending time: 07:31:01
    starting virgingalactic time : 07:31:02, ending time: 07:56:01
    starting virginorbit time : 07:56:02, ending time: 08:21:01
    starting JAXA time : 08:21:02, ending time: 08:46:01
    starting Roscosmos time : 08:46:02, ending time: 09:11:01
    starting areospace time : 09:11:02, ending time: 09:36:01
```

```
1 startdf = startdf[['user', 'user_id', 'text',
2                      'favorite_count', 'retweet_count',
3                      'mentions', 'hashtag', 'post_id']]
```

```
1 # don't need these columns anymore
2 startdf = startdf.drop(['index', '', 'symbols'], 1)
```

```
1 # in order to add our strings for our hashtags and mentions together we
2 # had to make the entire column a string, sinlucding empty cells
3 # here we go through and replace empty cells with nan values so pandas
4 # will read them as being empty instead of a string
5 for i in range(len(startdf)):
6     if startdf.at[i, 'mentions'] == ' ':
7         startdf.at[i, 'mentions'] = np.nan
8     else:
9         pass
10    if startdf.at[i, 'hashtag'] == ' ':
11        startdf.at[i, 'hashtag'] = np.nan
12    else:
13        pass
```

```
1 startdf.to_csv('tweets.4.7')
```

✓ 0s    completed at 12:35 AM