

Example from

- https://github.com/RodolfoFerro/pandas_twitter/blob/master/01-extracting-data.md (https://github.com/RodolfoFerro/pandas_twitter/blob/master/01-extracting-data.md) **Added some notes and explanation**

Extracting twitter data (tweepy + pandas)

```
In [1]: 1 # General:
2 import tweepy          # To consume Twitter's API
3 import pandas as pd    # To handle data
4 import numpy as np     # For number computing
5
6 # For plotting and visualization:
7 from IPython.display import display
8 import matplotlib.pyplot as plt
9 import seaborn as sns
10 %matplotlib inline
```

```
In [3]: 1 # Twitter App access keys for @user
2 #required keys and tokens
3
4 ACCESS_TOKEN = '1219125691028930560-CZyXhFlgCpMM8rG11KwuYJaMoX7uNa'
5 ACCESS_SECRET = '0qSAXeOLmH9pKDPVFy2pQzbloaGxRcbb0JJgnkhq5F2d4'
6 CONSUMER_KEY = 'HpU6B5BVTuwAfa5nYX1vAVxgD'
7 CONSUMER_SECRET = 'Cgcs5YIHpIp5Pu5US3N0XAX8N4jlJgPQrE4aK8LYkM89nqeTQa'
8
```

```
In [5]: 1 # We import our access keys:
2
3
4 # API's setup:
5 def twitter_setup():
6     """
7     Utility function to setup the Twitter's API
8     with our access keys provided.
9     """
10    # Authentication and access using keys:
11    auth = tweepy.OAuthHandler(CONSUMER_KEY, CONSUMER_SECRET)
12    auth.set_access_token(ACCESS_TOKEN, ACCESS_SECRET)
13
14    # Return API with authentication:
15    api = tweepy.API(auth)
16    return api
```

Tweets Extraction

```
In [6]: ▶ 1 # We create an extractor object:
2 extractor = twitter_setup()
3
4 # We create a tweet list as follows:
5 tweets = extractor.user_timeline(screen_name="realDonaldTrump", count=20
6 print("Number of tweets extracted: {}".format(len(tweets)))
7
8 # We print the most recent 5 tweets:
9 print("5 recent tweets:\n")
10 for tweet in tweets[:5]:
11     print(tweet.text)
12     print()
```

Number of tweets extracted: 200.

5 recent tweets:

RT @MeatInstitute: "We are grateful to @realDonaldTrump for protecting our nation's food supply," said @MeatInstitute Pres. & CEO Julie Ann...

No, I think Amash would make a wonderful candidate, especially since he is way behind in his district and has no ch... <https://t.co/RMbhUCax4n> (<https://t.co/RMbhUCax4n>)

THANK YOU @MarkLevinShow!
<https://t.co/GJSNM0LUQH> (<https://t.co/GJSNM0LUQH>)

RT @TrumpWarRoom: Yahoo reporter apologizes to President Trump after making false coronavirus test claim in Oval Office meeting
<https://t.co/...> (<https://t.co/...>)

At least they admit it. The Failing @nytimes & @washingtonpost never correct their Fake Reporting! <https://t.co/3aUKA4826K> (<https://t.co/3aUKA4826K>)

Creating a (pandas) DataFrame

```
In [7]: 1 # We create a pandas dataframe as follows:
2 data = pd.DataFrame(data=[tweet.text for tweet in tweets], columns=['Twe
3
4 # We display the first 10 elements of the dataframe:
5 display(data.head(10))
```

	Tweets
0	RT @MeatInstitute: "We are grateful to @realDo...
1	No, I think Amash would make a wonderful candi...
2	THANK YOU @MarkLevinShow! \nhttps://t.co/GJSNM...
3	RT @TrumpWarRoom: Yahoo reporter apologizes to...
4	At least they admit it. The Failing @nytimes &...
5	RT @ErinMPerrine: 📺 TUNE IN NOW 📺 \n\nTonight...
6	RT @TeamTrump: WATCH: Team Trump Online with @...
7	RT @TeamTrump: Tisa Clark, CEO Of J.D. Clark P...
8	RT @EquipoTrump: As Trump campaign senior advi...
9	RT @TrumpWarRoom: Joe Biden rambles and strugg...

```
In [8]: 1 # Internal methods of a single tweet object:
2 print(dir(tweets[0]))
```

```
['__class__', '__delattr__', '__dict__', '__dir__', '__doc__', '__eq__', '_
_format__', '__ge__', '__getattr__', '__getstate__', '__gt__', '__hash
__', '__init__', '__init_subclass__', '__le__', '__lt__', '__module__', '_
ne__', '__new__', '__reduce__', '__reduce_ex__', '__repr__', '__setattr__',
'__sizeof__', '__str__', '__subclasshook__', '__weakref__', '_api', '_jso
n', 'author', 'contributors', 'coordinates', 'created_at', 'destroy', 'enti
ties', 'favorite', 'favorite_count', 'favorited', 'geo', 'id', 'id_str', 'i
n_reply_to_screen_name', 'in_reply_to_status_id', 'in_reply_to_status_id_st
r', 'in_reply_to_user_id', 'in_reply_to_user_id_str', 'is_quote_status', 'l
ang', 'parse', 'parse_list', 'place', 'retweet', 'retweet_count', 'retweete
d', 'retweeted_status', 'retweets', 'source', 'source_url', 'text', 'trunca
ted', 'user']
```

In [9]: ▶

```

1 # We print info from the first tweet:
2 print(tweets[0].id)
3 print(tweets[0].created_at)
4 print(tweets[0].source)
5 print(tweets[0].favorite_count)
6 print(tweets[0].retweet_count)
7 print(tweets[0].geo)
8 print(tweets[0].coordinates)
9 print(tweets[0].entities)

```

1255558366455365634

2020-04-29 18:03:36

Twitter for iPhone

0

3173

None

None

```

{'hashtags': [], 'symbols': [], 'user_mentions': [{'screen_name': 'MeatInst
itute', 'name': 'North American Meat Institute', 'id': 30944238, 'id_str':
'30944238', 'indices': [3, 17]}, {'screen_name': 'realDonaldTrump', 'name':
'Donald J. Trump', 'id': 25073877, 'id_str': '25073877', 'indices': [39, 5
5]}, {'screen_name': 'MeatInstitute', 'name': 'North American Meat Institut
e', 'id': 30944238, 'id_str': '30944238', 'indices': [103, 117]}], 'urls':
[]}]

```

Adding relevant info to our dataframe

In [10]: ▶

```

1 # We add relevant data:
2 data['len'] = np.array([len(tweet.text) for tweet in tweets])
3 data['ID'] = np.array([tweet.id for tweet in tweets])
4 data['Date'] = np.array([tweet.created_at for tweet in tweets])
5 data['Source'] = np.array([tweet.source for tweet in tweets])
6 data['Likes'] = np.array([tweet.favorite_count for tweet in tweets])
7 data['RTs'] = np.array([tweet.retweet_count for tweet in tweets])

```

```
In [11]: 1 # Display of first 10 elements from dataframe:
2 display(data.head(10))
```

	Tweets	len	ID	Date	Source	Likes	RTs
0	RT @MeatInstitute: "We are grateful to @realDo...	144	1255558366455365634	2020-04-29 18:03:36	Twitter for iPhone	0	3173
1	No, I think Amash would make a wonderful candi...	140	1255510996623527936	2020-04-29 14:55:22	Twitter for iPhone	41646	10442
2	THANK YOU @MarkLevinShow! \nhttps://t.co/GJSNM...	50	1255507868280856581	2020-04-29 14:42:56	Twitter for iPhone	45413	14846
3	RT @TrumpWarRoom: Yahoo reporter apologizes to...	140	1255485332339990528	2020-04-29 13:13:23	Twitter for iPhone	0	7670
4	At least they admit it. The Failing @nytimes &...	126	1255484391364665351	2020-04-29 13:09:39	Twitter for iPhone	54994	17476
5	RT @ErinMPerrine: 📺 TUNE IN NOW 📺 \n\nTonight...	147	1255482641245188108	2020-04-29 13:02:42	Twitter for iPhone	0	3231
6	RT @TeamTrump: WATCH: Team Trump Online with @...	140	1255482613927739393	2020-04-29 13:02:35	Twitter for iPhone	0	3346
7	RT @TeamTrump: Tisa Clark, CEO Of J.D. Clark P...	140	1255482503349047298	2020-04-29 13:02:09	Twitter for iPhone	0	3542
8	RT @EquipoTrump: As Trump campaign senior advi...	140	1255482443647385602	2020-04-29 13:01:54	Twitter for iPhone	0	5175
9	RT @TrumpWarRoom: Joe Biden rambles and strugg...	124	1255481563195850754	2020-04-29 12:58:25	Twitter for iPhone	0	4471

Visualization and basic statistics

Averages and popularity

```
In [12]: 1 # We extract the mean of lenghts:
2 mean = np.mean(data['len'])
3
4 print("The lenght's average in tweets: {}".format(mean))
```

The lenght's average in tweets: 117.915

- To extract more data, we will use some pandas' functionalities:

```

In [13]: ▶ 1 # We extract the tweet with more FAVs and more RTs:
2
3 fav_max = np.max(data['Likes'])
4 rt_max = np.max(data['RTs'])
5
6 fav = data[data.Likes == fav_max].index[0]
7 rt = data[data.RTs == rt_max].index[0]
8
9 # Max FAVs:
10 print("The tweet with more likes is: \n{}".format(data['Tweets'][fav]))
11 print("Number of likes: {}".format(fav_max))
12 print("{} characters.\n".format(data['len'][fav]))
13
14 # Max RTs:
15 print("The tweet with more retweets is: \n{}".format(data['Tweets'][rt]))
16 print("Number of retweets: {}".format(rt_max))
17 print("{} characters.\n".format(data['len'][rt]))

```

The tweet with more likes is:
 Happy Birthday to Melania, our great First Lady!
 Number of likes: 574952
 48 characters.

The tweet with more retweets is:
 RT @realDonaldTrump: BRILLIANT, A MUST WATCH! @RepDanCrenshaw <https://t.co/W6pGsJQ2ua>
 Number of retweets: 93830
 85 characters.

Time series

- Pandas has its own object for time series. Since we have a whole vector with creation dates, we can construct time series respect tweets lengths, likes and retweets.
- The way we do it is:

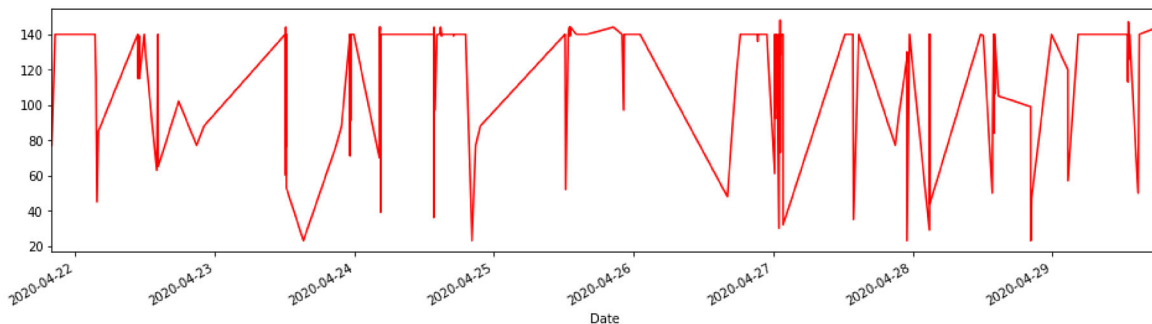
```

In [14]: ▶ 1 # We create time series for data:
2 tlen = pd.Series(data=data['len'].values, index=data['Date'])
3 tfav = pd.Series(data=data['Likes'].values, index=data['Date'])
4 tret = pd.Series(data=data['RTs'].values, index=data['Date'])

```

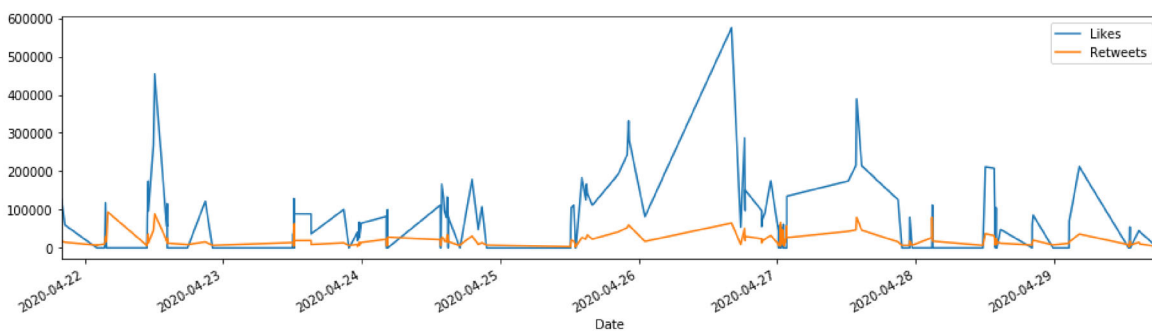
- And if we want to plot the time series, pandas already has its own method in the object. We can plot a time series as follows:

```
In [15]: 1 # Lenghts along time:
2 tlen.plot(figsize=(16,4), color='r');
```



- And to plot the likes versus the retweets in the same chart:

```
In [16]: 1 # Likes vs retweets visualization:
2 tfav.plot(figsize=(16,4), label="Likes", legend=True)
3 tret.plot(figsize=(16,4), label="Retweets", legend=True);
```



- Pie charts of sources

```
In [17]: ▶ 1 # We obtain all possible sources:
2 sources = []
3 for source in data['Source']:
4     if source not in sources:
5         sources.append(source)
6
7 # We print sources list:
8 print("Creation of content sources:")
9 for source in sources:
10     print("* {}".format(source))
```

Creation of content sources:

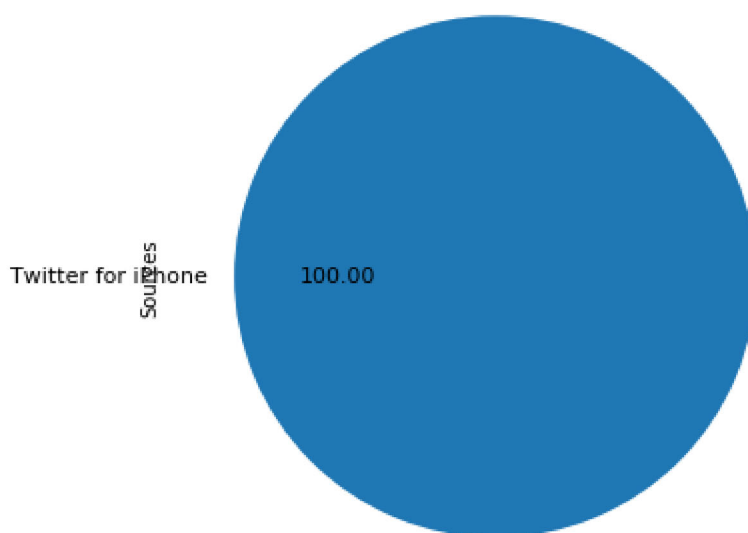
* Twitter for iPhone

- With the following output, we realize that basically this twitter account has two sources:
Creation of content sources: * Twitter for iPhone * Media Studio
- We now count the number of each source and create a pie chart. You'll notice that this code cell is not the most optimized one.


```

In [18]: ▶ 1 # We create a numpy vector mapped to Labels:
2 percent = np.zeros(len(sources))
3
4 for source in data['Source']:
5     for index in range(len(sources)):
6         if source == sources[index]:
7             percent[index] += 1
8         pass
9
10 percent /= 100
11
12 # Pie chart:
13 pie_chart = pd.Series(percent, index=sources, name='Sources')
14 pie_chart.plot.pie(fontsize=11, autopct='%.2f', figsize=(6, 6));

```



Sentiment analysis

- **Importing textblob**
- textblob will allow us to do sentiment analysis in a very simple way. We will also use the re library from Python, which is used to work with regular expressions.
- For this, I'll provide you two utility functions to: a) clean text (which means that any symbol distinct to an alphanumeric value will be remapped into a new one that satisfies this condition), and b) create a classifier to analyze the polarity of each tweet after cleaning the text in it. I won't explain the specific way in which the function that cleans works, since it would be extended and it might be better understood in the official redocumentation.

```
In [19]: ▶ 1 from textblob import TextBlob
2 import re
3
4 def clean_tweet(tweet):
5     '''
6     Utility function to clean the text in a tweet by removing
7     links and special characters using regex.
8     '''
9     return ' '.join(re.sub("(@[A-Za-z0-9]+)|([^0-9A-Za-z \t])|(\w+:\/\/\w+)",
10
11
12 def analyze_sentiment(tweet):
13     '''
14     Utility function to classify the polarity of a tweet
15     using textblob.
16     '''
17     analysis = TextBlob(clean_tweet(tweet))
18     if analysis.sentiment.polarity > 0:
19         return 1
20     elif analysis.sentiment.polarity == 0:
21         return 0
22     else:
23         return -1
```

```
In [20]: 1 # We create a column with the result of the analysis:
2 data['SA'] = np.array([ analyze_sentiment(tweet) for tweet in data['Twee
3
4 # We display the updated dataframe with the new column:
5 display(data.head(10))
```

	Tweets	len	ID	Date	Source	Likes	RTs	SA
0	RT @MeatInstitute: "We are grateful to @realDo...	144	1255558366455365634	2020-04-29 18:03:36	Twitter for iPhone	0	3173	0
1	No, I think Amash would make a wonderful candi...	140	1255510996623527936	2020-04-29 14:55:22	Twitter for iPhone	41646	10442	1
2	THANK YOU @MarkLevinShow! \nhttps://t.co/GJSNM...	50	1255507868280856581	2020-04-29 14:42:56	Twitter for iPhone	45413	14846	0
3	RT @TrumpWarRoom: Yahoo reporter apologizes to...	140	1255485332339990528	2020-04-29 13:13:23	Twitter for iPhone	0	7670	-1
4	At least they admit it. The Failing @nytimes &...	126	1255484391364665351	2020-04-29 13:09:39	Twitter for iPhone	54994	17476	-1
5	RT @ErinMPerrine: 🤖🤖 TUNE IN NOW 🤖🤖 \n\nTonight...	147	1255482641245188108	2020-04-29 13:02:42	Twitter for iPhone	0	3231	1
6	RT @TeamTrump: WATCH: Team Trump Online with @...	140	1255482613927739393	2020-04-29 13:02:35	Twitter for iPhone	0	3346	0
7	RT @TeamTrump: Tisa Clark, CEO Of J.D. Clark P...	140	1255482503349047298	2020-04-29 13:02:09	Twitter for iPhone	0	3542	1
8	RT @EquipoTrump: As Trump campaign senior advi...	140	1255482443647385602	2020-04-29 13:01:54	Twitter for iPhone	0	5175	0
9	RT @TrumpWarRoom: Joe Biden rambles and strugg...	124	1255481563195850754	2020-04-29 12:58:25	Twitter for iPhone	0	4471	0

Analyzing the results

- To have a simple way to verify the results, we will count the number of neutral, positive and negative tweets and extract the percentages.

```
In [21]: 1 # We construct lists with classified tweets:
2 pos_tweets = [ tweet for index, tweet in enumerate(data['Tweets']) if da
3 neu_tweets = [ tweet for index, tweet in enumerate(data['Tweets']) if da
4 neg_tweets = [ tweet for index, tweet in enumerate(data['Tweets']) if da
```

- Now that we have the lists, we just print the percentages:

```
In [22]: ▶ 1 # We print percentages:
          2
          3 print("Percentage of positive tweets: {}".format(len(pos_tweets)*100/len
          4 print("Percentage of neutral tweets: {}".format(len(neu_tweets)*100/len
          5 print("Percentage de negative tweets: {}".format(len(neg_tweets)*100/len
```

Percentage of positive tweets: 48.0%

Percentage of neutral tweets: 34.0%

Percentage de negative tweets: 18.0%

```
In [ ]: ▶ 1
```