



Store, Analyze and Visualize “tweets” from Twitter

Team members:

Lakshmi Chamala
Gulnoza Khakimova
Khushbu Kolhe



Instructor: Dr. Praveen Rao
CS5540
Principles of Big Data Management

Introduction

Twitter is an online social network where users communicate with each other through messages called “tweets”. Tweets are visible to everyone but communicator can restrict message delivery to just his/her followers. Users can group “tweets” by topic or hashtags. The most posted topics are called “trending topic”. Those topics help users to identify what is happening in the world. On average there are 1.6 billion of messages posted every day. In this project we will collect “tweets”, analyze them based on specific criteria and visualize results.

Our research consists of three phases: The first step will be to collect tweets on the following topics: “Bitcoin”, “Forbes” and “Winter Olympics”. In the second step we will analyze “tweets” by writing queries. The last step will be demonstration of our project.

Phase 1

In this phase we need to collect “tweets” for the following topics: “Bitcoin”, “Forbes”, “Winter Olympics”. First we had to create and register our application on <http://apps.twitter.com> in order to access Twitter data and interact with Twitter API. After registration we were given a consumer key, consumer secret and access tokens which are used to get access to Twitter data.

Application Settings

Keep the "Consumer Secret" a secret. This key should never be human-readable in your application.

Consumer Key (API Key) LBIJTcvUFWxjG2cC5xYaRLaLQ

Consumer Secret (API Secret) ExuqqOHAW9C5s3GiNQpZLBclwdxWJKgb7CBioMaqFa45dpujP

We tried different ways to get tweets using Twitter streaming API . After collecting good amount of tweets, which we stored into google drive (link to collected tweets: https://drive.google.com/open?id=1EM9e3iK-OxoyAe2XetQRdU_LWNjyCBn6) we were able to extract hashtags and URLs for each “tweet” so we could run word count in Apache Hadoop and Apache Spark.

Apache Hadoop

Collect tweets

Below is a code on Python which were used to collect “tweets”:

```
*collecttweets.py (-/Desktop)
File Edit View Search Tools Documents Help
from tweepy.streaming import StreamListener
from tweepy import OAuthHandler
from tweepy import Stream

consumer_key="LBIJTcvUFWxjG2cC5xYaRLaLQ"
consumer_secret="ExuqqOHAW9C5s3GiNQpZLBclwdxWJKgb7CBioMaqFa45dpujP"
access_token="926091370363793413-BXu0ZBqH58kUDWs54XVWgH8AEWm8s"
access_token_secret="FskYTykFTJnqplMod1HPZpxAKnr2Dn2ExKl0vdxJ4DTqP"

class StdOutListener(StreamListener):

    def on_data(self, data):
        with open('tweets_output.json','a') as tf:
            tf.write(data)
            print(data)
            return True

if __name__ == '__main__':
    auth = OAuthHandler(consumer_key, consumer_secret)
    auth.set_access_token(access_token, access_token_secret)
    stream = Stream(auth, StdOutListener())
    stream.filter(track=['Forbes','Bitcoin','Winter Olympics'])
```

Word count on Apache Hadoop

Following command was used to perform word count on extracted hashtags and URL's from "tweets" on Apache Hadoop:

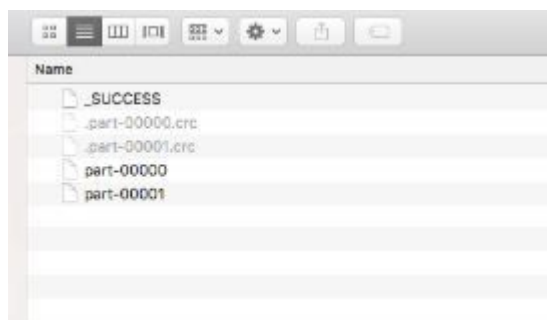
```
bin/hadoop jar share/hadoop/mapreduce/hadoop-mapreduce-examples-2.8.1.jar wordcount /Project/output.txt /Project/out
```

Apache Spark

Below is a code for Apache Spark which was used to perform "wordcount" on extracted hashtags and URLs

```
WordCount.scala x
1
2 import org.apache.spark.{SparkContext, SparkConf}
3
4 object WordCount {
5
6   def main(args: Array[String]) {
7
8     val sparkConf = new SparkConf().setAppName("WordCount").setMaster("local[*]")
9
10    val sc=new SparkContext(sparkConf)
11
12    val input=sc.textFile("/Volumes/Data/PrinciplesOfBigData/WordCountHashTags/hashtagsAndUrl.txt")
13
14    val wc=input.flatMap(line=>{line.split(" ")}).map(word=>{word,1}).cache()
15
16    val output=wc.reduceByKey(_+_ )
17
18    output.saveAsTextFile("sparkWordCountOutput")
19
20    val o=output.collect()
21
22    var s:String="Words:Count \n"
23    o.foreach{case (word,count)=>{
24
25      s+=word+" : "+count+"\n"
26
27    }}
28
29   }
30
31 }
32
33
```

Output:



```

WordCount.scala  part-00000  part-00001
1  [[('http://ift.tt/2CmzYCW',,1)
2  (('http://ift.tt/2CpB0HV',,1)
3  (('https://twitter.com/i/web/status/9684404528667148290']],1)
4  ('AméricaLatina',,1)
5  (('https://fb.me/WvEX7g4i']],1)
6  (('https://twitter.com/i/web/status/968443369562234880']],1)
7  (('Easy',,1)
8  (('http://www.gizmodo.co.uk/2018/02/50-cent-no-longer-bragging-about-a-bitcoin-fortune-now-that-the-us-government-is-interested/']],2)
9  (('https://twitter.com/i/web/status/968440201134641152']],1)
10 (('https://twitter.com/i/web/status/968390561445261313']],1)
11 ('firm',,1)
12 (('https://is.gd/H5gB6J']],1)
13 (('https://twitter.com/i/web/status/968466000118198272']],1)
14 (('btcnews',,3)
15 (('https://twitter.com/i/web/status/968446470005309312']],1)
16 (('http://ift.tt/2FAjSUL',,1)
17 (('https://www.reddit.com/r/Bitcoin/comments/90cw5/i_built_this_lightning_network_searchanalysis/']],1)
18 (('https://twitter.com/i/web/status/968398381964107776']],1)
19 (('https://twitter.com/i/web/status/968452375957585920']],1)
20 (('https://www.youtube.com/watch?v=w2sGn_dzRoo']],,9)
21 (('https://twitter.com/i/web/status/968415308401782786']],1)
22 (('https://twitter.com/i/web/status/968390634400965000']],1)
23 (('https://twitter.com/i/web/status/968428042724036608']],1)
24 (('https://twitter.com/i/web/status/968489170553790464']],1)
25 ('SEO',,5)
26 (('https://www.betmoose.com/bet/fifa-world-cup-2018-champion-3624?ref=scacco']],1)
27 (('employee recharge',,1)
28 ('Grocain',,1)
29 (('https://forum.bitcoin.com/viewtopic.php?t=719556p=148927#p148927']],1)
30 (('http://dlvr.it/QJ0Gwd']],1)
31 (('https://twitter.com/i/web/status/968396732856848384']],1)
32 ('Motivation',,3)
33 ('etc',,18)
34 (('https://twitter.com/deex_exchange/status/967933559775350784']],1)
35 (('https://twitter.com/i/web/status/968420140743300993']],1)

```

Hashtags and URLs extraction

Code which was used to get Hashtags and URLs from collected tweets:

```
collecttweets.py x HashtagAndURlExtraction.py x hashtagAndURLoutput.txt x
1 import codecs
2 from datetime import datetime
3 import json
4 # import requests
5 import os
6 import string
7 import sys
8 import time
9
10
11 def parse_json_tweet(line):
12     tweet = json.loads(line)
13
14     hashtags = [hashtag['text'] for hashtag in tweet['entities']['hashtags']]
15     urls = [url['expanded_url'] for url in tweet['entities']['urls']]
16
17
18
19     return [hashtags, urls]
20
21
22
23 '''start main'''
24 if __name__ == "__main__":
25     file_timeordered_json_tweets = codecs.open("tweets_output.json", 'r', 'utf-8')
26     fout = codecs.open("hashtagAndURLoutput.txt", 'w', 'utf-8')
27
28     for line in file_timeordered_json_tweets:
29         try:
30             [hashtags, urls] = parse_json_tweet(line)
31             fout.write(str([hashtags, urls]) + "\n")
32             # fout.write("Hashtag" + str([hashtags]) + "URLs" + str([urls]) + "\n")
33         except:
34             pass
35     file_timeordered_json_tweets.close()
36     fout.close()
37
38
```

```
collecttweets.py x HashtagAndURlExtraction.py x
1 import codecs
2 from datetime import datetime
3 import json
4 # import requests
5 import os
6 import string
7 import sys
8 import time
9
10
11 def parse_json_tweet(line):
12     tweet = json.loads(line)
13
14     hashtags = [hashtag['text'] for hashtag in tweet['entities']['hashtags']]
15     urls = [url['expanded_url'] for url in tweet['entities']['urls']]
16
17
18
19     return [hashtags, urls]
20
21
22
23 '''start main'''
24 if __name__ == "__main__":
25     file_timeordered_json_tweets = codecs.open("tweets_output.json", 'r', 'utf-8')
26     fout = codecs.open("hashtagAndURLoutput.txt", 'w', 'utf-8')
27
28     for line in file_timeordered_json_tweets:
29         try:
30             [hashtags, urls] = parse_json_tweet(line)
31             fout.write(str([hashtags, urls]) + "\n")
32             # fout.write("Hashtag" + str([hashtags]) + "URLs" + str([urls]) + "\n")
33         except:
34             pass
35     file_timeordered_json_tweets.close()
36     fout.close()
37
38
```

Output:

```

[[], []]
[[], []]
[['nachrichten'], ['http://ift.tt/2ov0hhz']]
[['SpendMoneyWisely'], ['http://bit.ly/2ASnUbe']]
[[], []]
[['EXO'], []]
[['BestFanArmy', 'iHeartAwards', 'EXOL'], ['https://twitter.com/weareoneexo/status/968823850129620992']]
[[], []]
[[], []]
[[], []]
[['Olympics', 'PyeongChang2018'], []]
[['EXO'], []]
[[], []]
[['PyeongChang2018', 'Olympics'], []]
[['Verge', 'XVG', 'Wraith'], ['https://twitter.com/crypto_sarah9/status/967485388497842177']]
[[], []]
[['GetThatMoney'], ['http://bit.ly/2ARFXU0']]
[[], ['http://ift.tt/2H0OT8r', 'https://twitter.com/i/web/status/968373716097921024']]
[[], []]
[[], ['http://ift.tt/2ELkxv5', 'https://twitter.com/i/web/status/9683737178000008449']]
[['BTSARMY', 'BestFanArmy'], []]
[[], ['https://twitter.com/businessinsider/status/968347478381666384']]
[[], []]
[['EXO'], []]
[[], []]
[['Bitcoin', 'btc', 'bitcoin'], ['http://www.infocryptos.com/bitcoin/ethereum-price-technical-analysis-eth-usd-gains-traction/?utm_source=tw']]
[[], ['http://ift.tt/2CmxyzX']]
[[], []]
[['GoGetTheMoney'], ['http://bit.ly/2ARgfU4']]
[[], ['https://trib.al/WhytNih']]
[['BankingOnBitcoin'], ['https://twitter.com/readbitcoins/status/968342099501469698']]
[[], ['https://twitter.com/witnmbc/status/967763124261568512']]
[['Bitcoin', 'btc', 'bitcoin'], ['http://www.infocryptos.com/bitcoin/bitcoin-buyers-tapping-currency-into-cars-goods/?utm_source=tw']]
[[], []]
[[], ['http://www.forbes.com/sites/tanaherman/2018/02/26/videos-from-adele-ed-sheeran-cl-exo-others-spike-on-youtube-following-pyeongchang-olympics/#79286aef21fd']]
[[], ['https://twitter.com/i/web/status/968373726705340418']]
[[], ['http://youtu.be/e_9QouZMjvc']]
[[], []]
[['Bermuda', 'BVI', 'Caribbean', 'CentralAmerica', 'ico'], ['https://www.youtube.com/watch?v=yxuF1mFggLA']]
[[], ['https://twitter.com/witnmbc/status/967763124261568512']]
[['bitcoin'], ['http://ift.tt/2oNw0B']]
[['EXO', 'Olympics2018'], []]
[[], ['https://www.thenational.ae/arts-culture/music/k-pop-stars-cl-and-exo-close-out-winter-olympics-in-pyeongchang-1.707875']]
[[], []]
[[], ['http://www.forbes.com/sites/andrewarnold/2018/02/16/how-ar-and-vr-are-revolutionizing-the-supply-chain/#7c26c2997f33']]
[['iHeartAwards', 'BestFanArmy', 'EXOL'], []]
[['EXO'], []]
[['MoneyMature'], ['http://bit.ly/2ATDF1I']]
[['EXO'], []]
[['EXO'], []]
[['Kpop', 'EXO', 'EXOL', 'KOR', 'FirstLady'], []]
[[], ['https://trib.al/WhytNih']]
[['cryptocurrency'], []]
[[], ['https://www.forbes.com/sites/shelliekarabell/2018/02/14/executive-compensation-is-out-of-control-what-now/#46e7ee36431f']]
[[], []]
[[], []]

```

Link to Github repository which contains code, input, output and logs:
<https://github.com/Gnkhakimova/CS5540-Twitter>