# Principles of Big Data Management

# Project Phase - 1

**Team Members:**
**Ankitha Wankhede - 16233344**
**Rajeswari Devi Namana -18135299**
**Sravya Reddy Bokka - 16240386**

## Goal:

- To Collect Tweets using Twitter's Streaming APIs (e.g., 100K Tweets)
  (https://dev.twitter.com/docs/streaming-apis)

- Extract all the hashtags and URLs in the tweets
- Run the WordCount example in Apache Hadoop and Apache Spark on the extracted hashtags/URLs and collect the output and log files from Hadoop. Add a README file.

## Step 1: Collection of tweets from Twitter Streaming API:

Executed code in python using tweepy library.

## Code Used:

```python
#Import the necessary methods from tweepy library
from tweepy.streaming import StreamListener
from tweepy import OAuthHandler
from tweepy import Stream

#Variables that contains the user credentials to access Twitter API
access_token = "888958069098635264-AsvKqThYHQL02oVnfa148kEU9ooPks2"
access_token_secret = "rgWc4aerpJruAwdZHpjAZYHn55qD92TJnUUpaNd8G1j7I"
consumer_key = "zKhTdLsdjIG8HlfRa0Fhxoofo"
consumer_secret = "vorTGnycUlmTr6RhjocPmYJ0A3yttUTpXNBIyjB7bNw2YOSGvC"

#This is a basic listener that just prints received tweets to stdout.
class StdOutListener(StreamListener):

    def on_data(self, data):
        print(data)
        return True

    def on_error(self, status):
        print(status)

if __name__ == '__main__':

    #This handles Twitter authentication and the connection to Twitter Streaming API
    l = StdOutListener()
    auth = OAuthHandler(consumer_key, consumer_secret)
    auth.set_access_token(access_token, access_token_secret)
    stream = Stream(auth, l)
```

```
#This line filter Twitter Streams to capture data by the keywords: 'mom'
#stream.filter(track=['hashtags', "(?P<url>https?://[^\s]+)"])
stream.filter(track=["mom"])
```

## Code Screenshots :

## Output Screenshot: Collected Tweets:



## ***Step 2: Extraction of URLs and hashtags from collected tweets:***

## Code:

```python
import codecs
from datetime import datetime
import json
import os
import string
import sys
import time


def parse_json_tweet(line):
    tweet = json.loads(line)
    # print line
    if tweet['lang'] != 'en':
        # print "non-english tweet:", tweet['lang'], tweet
        return ['', '', '', [], [], []]

    date = tweet['created_at']
    id = tweet['id']
    nfollowers = tweet['user']['followers_count']
    nfriends = tweet['user']['friends_count']

    if 'retweeted_status' in tweet:
        text = tweet['retweeted_status']['text']
    else:
```

```python
        text = tweet['text']

    hashtags = [hashtag['text'] for hashtag in tweet['entities']['hashtags']]
    users = [user_mention['screen_name'] for user_mention in
tweet['entities']['user_mentions']]
    urls = [url['expanded_url'] for url in tweet['entities']['urls']]

    media_urls = []
    if 'media' in tweet['entities']:
        media_urls = [media['media_url'] for media in tweet['entities']['media']]

    return [hashtags, urls]


'''start main'''
if __name__ == "__main__":
    file_timeordered_json_tweets = codecs.open(sys.argv[1], 'r', 'utf-8')
    fout = codecs.open(sys.argv[2], 'w', 'utf-8')

    # efficient line-by-line read of big files
    for line in file_timeordered_json_tweets:
        try:
            [tweet_gmttime, tweet_id, text, hashtags, users, urls, media_urls,
nfollowers, nfriends] = parse_json_tweet(
                line)
            #    if not tweet_gmttime: continue
            #    fout.write(line)
            # "created_at":"Mon Feb 17 14:14:44 +0000 2014"
            try:
                c = time.strptime(tweet_gmttime.replace("+0000", ''), '%a %b %d
%H:%M:%S %Y')
            except:
                print("pb with tweet_gmttime", tweet_gmttime, line)
                pass
            tweet_unixtime = int(time.mktime(c))
            #        fout.write(line)
            fout.write(str(
                [tweet_unixtime, tweet_gmttime, tweet_id, text, hashtags, users, urls,
media_urls, nfollowers,
                 nfriends]) + "\n")
        except:
            # print "pb with tweet:", line
            #        print sys.exc_info()[0], line
            pass
    file_timeordered_json_tweets.close()
    fout.close()
```
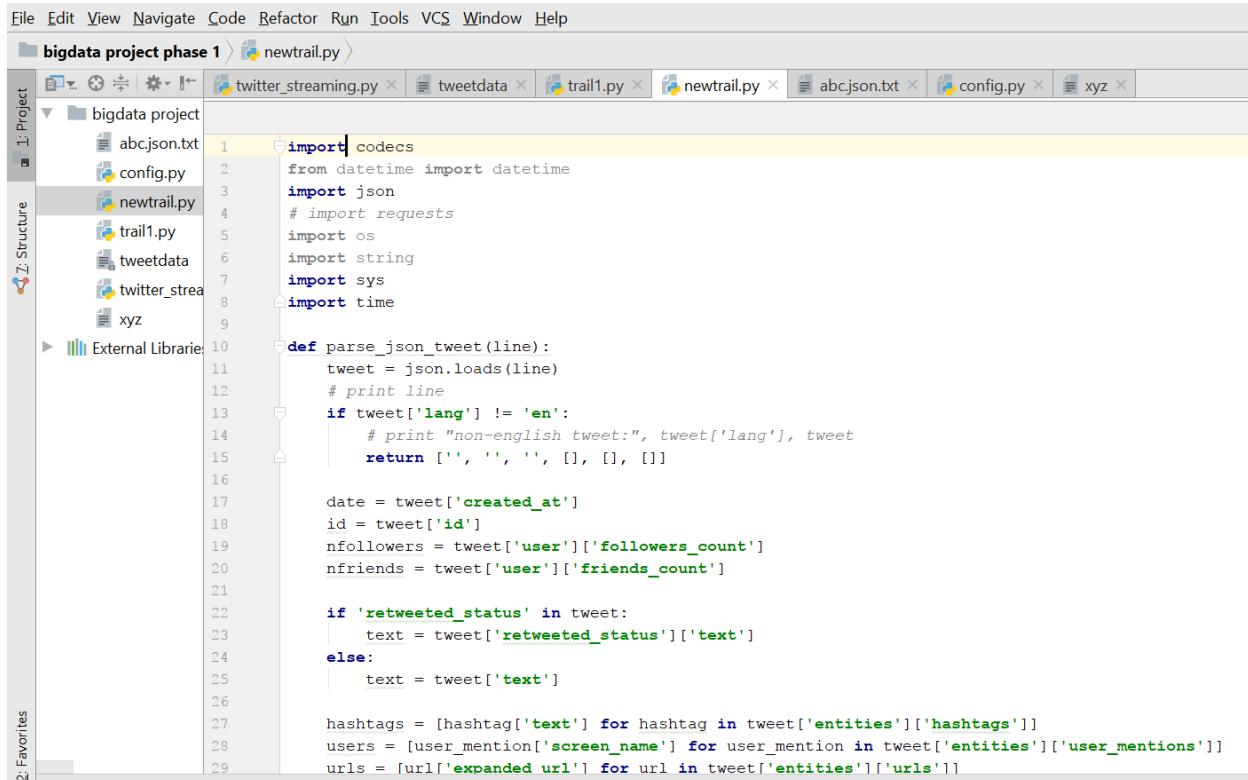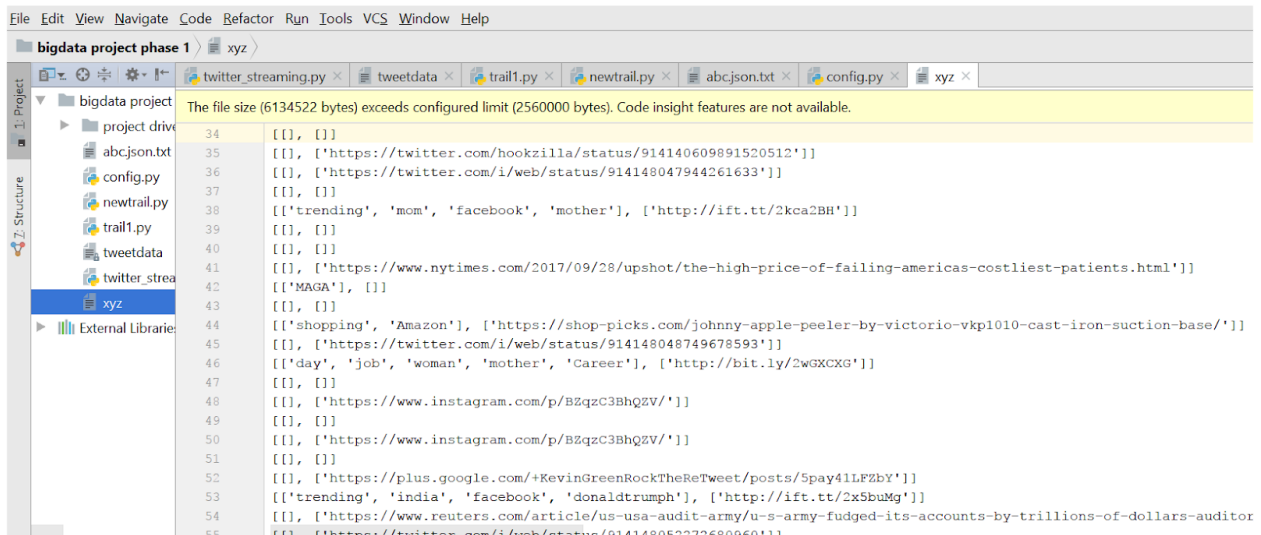
## Code Screenshots :

bigdata project phase 1 - [C:\Users\rajin\Desktop\Masters subjects\Principles of big data\bigdata project phase 1] - ...\newtrail.py - PyCharm Community Edition 2017.1.3

File  Edit  View  Navigate  Code  Refactor  Run  Tools  VCS  Window  Help

bigdata project phase 1 > newtrail.py

twitter_streaming.py  |  tweetdata  |  trail1.py  |  newtrail.py  |  abc.json.txt  |  config.py  |  xyz

bigdata project
- abc.json.txt
- config.py
- newtrail.py
- trail1.py
- tweetdata
- twitter_strea
- xyz
- External Librarie

```python
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  def parse_json_tweet(line):
11      tweet = json.loads(line)
12      # print line
13      if tweet['lang'] != 'en':
14          # print "non-english tweet:", tweet['lang'], tweet
15          return ['', '', '', [], [], []]
16
17      date = tweet['created_at']
18      id = tweet['id']
19      nfollowers = tweet['user']['followers_count']
20      nfriends = tweet['user']['friends_count']
21
22      if 'retweeted_status' in tweet:
23          text = tweet['retweeted_status']['text']
24      else:
25          text = tweet['text']
26
27      hashtags = [hashtag['text'] for hashtag in tweet['entities']['hashtags']]
28      users = [user_mention['screen_name'] for user_mention in tweet['entities']['user_mentions']]
29      urls = [url['expanded_url'] for url in tweet['entities']['urls']]
```

## Output Screenshot: Extracted URLs and Hashtags:

bigdata project phase 1 - [C:\Users\rajin\Desktop\Masters subjects\Principles of big data\bigdata project phase 1] - ...\xyz - PyCharm Community Edition 2017.1.3

File  Edit  View  Navigate  Code  Refactor  Run  Tools  VCS  Window  Help

bigdata project phase 1 > xyz

twitter_streaming.py  |  tweetdata  |  trail1.py  |  newtrail.py  |  abc.json.txt  |  config.py  |  xyz

bigdata project
- project drive
- abc.json.txt
- config.py
- newtrail.py
- trail1.py
- tweetdata
- twitter_strea
- xyz
- External Librarie

The file size (6134522 bytes) exceeds configured limit (2560000 bytes). Code insight features are not available.

```
34   [[], []]
35   [[], ['https://twitter.com/hookzilla/status/914140609891520512']]
36   [[], ['https://twitter.com/i/web/status/914148047944261633']]
37   [[], []]
38   [['trending', 'mom', 'facebook', 'mother'], ['http://ift.tt/2kca2BH']]
39   [[], []]
40   [[], []]
41   [[], ['https://www.nytimes.com/2017/09/28/upshot/the-high-price-of-failing-americas-costliest-patients.html']]
42   [['MAGA'], []]
43   [[], []]
44   [['shopping', 'Amazon'], ['https://shop-picks.com/johnny-apple-peeler-by-victorio-vkp1010-cast-iron-suction-base/']]
45   [[], ['https://twitter.com/i/web/status/914148048749678593']]
46   [['day', 'job', 'woman', 'mother', 'Career'], ['http://bit.ly/2wGXCXG']]
47   [[], []]
48   [[], ['https://www.instagram.com/p/BZqzC3BhQZV/']]
49   [[], []]
50   [[], ['https://www.instagram.com/p/BZqzC3BhQZV/']]
51   [[], []]
52   [[], ['https://plus.google.com/+KevinGreenRockTheReTweet/posts/5pay41LFZbY']]
53   [['trending', 'india', 'facebook', 'donaldtrumph'], ['http://ift.tt/2x5buMg']]
54   [[], ['https://www.reuters.com/article/us-usa-audit-army/u-s-army-fudged-its-accounts-by-trillions-of-dollars-auditor
55   [[], ['https://twitter.com/i/web/status/914148052272680960']]
```

## Step 3: Running Word Count Program in Hadoop:

## Screenshots of commands executed while running word count program:

## Word Count of URLs:



```
['https://www.google.com/amp/amp.sacbee.com/opinion/california-forum/article175434686.html']]	1
['https://www.google.com/amp/amp.slate.com/articles/health_and_science/science/2017/09/puerto_rico_needs_long_term_support_not_just_short_term_aid.html
']]	2
['https://www.google.com/amp/amp.thedailybeast.com/the-time-donald-trump-turned-away-in-disgust-while-a-man-bled-to-death-in-front-of-him']]	1
['https://www.google.com/amp/amp.timeinc.net/fortune/2014/03/20/bill-clinton-on-leadership']]	1
['https://www.google.com/amp/amp.timeinc.net/fortune/2017/09/29/colin-kaepernick-nfl-protest']]	1
['https://www.google.com/amp/amp.timeinc.net/si/tech-media/2017/09/26/nfl-ratings-increase-donald-trump-comments-monday-night-football%3Fsource=dam']]1
['https://www.google.com/amp/amp.timeinc.net/time/4606082/jfk-assassination-secrets/%3fsource=dam']]	1
['https://www.google.com/amp/amp.usatoday.com/story/470245001/']]	1
['https://www.google.com/amp/amp.usatoday.com/story/706196001/']]	1
['https://www.google.com/amp/amp.usatoday.com/story/716464001/']]	1
['https://www.google.com/amp/blog.organizer.com/what-is-deep-canvassing%3Fhs_amp=true']]	1
['https://www.google.com/amp/deadline.com/2017/09/donald-trump-twitter-puerto-rico-big-decisions-rebuilding-1202179269/amp/']]	1
['https://www.google.com/amp/denver.cbslocal.com/2017/09/16/rocky-mountain-arsenal-lawsuit/amp/']]	1
['https://www.google.com/amp/ew.com/gallery/stars-against-trump/amp/']]	4
['https://www.google.com/amp/fox2now.com/2017/09/29/illinois-high-school-football-team-honors-first-responders-while-taking-field/amp/']]	1
['https://www.google.com/amp/globalnation.inquirer.net/72279/philippines-a-jewish-refuge-from-the-holocaust/amp',	1
['https://www.google.com/amp/m.huffpost.com/us/entry/us_596cfc31e4b0376db8b659e1/amp']]	1
['https://www.google.com/amp/mobile.reuters.com/article/amp/idUSKCN1C435K']]	1
['https://www.google.com/amp/mobile.reuters.com/article/amp/idUSKCN1C50IH']]	4
['https://www.google.com/amp/mondoweiss.net/2016/12/sentenced-community-palestinian/amp/']]	1
['https://www.google.com/amp/nypost.com/2014/07/09/mayor-during-katrina-gets-10-years-for-corruption/amp/']]	5
['https://www.google.com/amp/nypost.com/2017/09/27/puerto-rico-should-have-been-ready-for-maria/amp/#ampshare=http://nypost.com/2017/09/27/puerto-rico-
should-have-been-ready-for-maria/']]	1
['https://www.google.com/amp/nypost.com/2017/09/27/puerto-rico-should-have-been-ready-for-maria/amp/']]	1
['https://www.google.com/amp/people.com/politics/lin-manuel-miranda-donald-trump-straight-to-hell-puerto-rico/amp/']]	1
['https://www.google.com/amp/pix11.com/2017/09/16/instagram-video-of-queens-students-singing-popular-song-with-racial-slur-causes-controversy/amp/']]	1
['https://www.google.com/amp/profootballtalk.nbcsports.com/2017/09/24/report-raiders-offensive-line-plans-to-kneel-as-a-group/amp/']]	1
['https://www.google.com/amp/s/amp.businessinsider.com/how-does-moviepass-make-money-2017-8']]	1
['https://www.google.com/amp/s/amp.businessinsider.com/salary-after-taxes-us-cities-2017-9']]	1
['https://www.google.com/amp/s/amp.cincinnati.com/amp/716329001']]	1
['https://www.google.com/amp/s/amp.cnn.com/cnn/2016/04/29/politics/donald-trump-tweets-daniel-scavino/index.html']]	1
['https://www.google.com/amp/s/amp.cnn.com/cnn/2017/09/26/politics/saudi-arabia-woman-drive/index.html']]	6
['https://www.google.com/amp/s/amp.cnn.com/cnn/2017/09/27/us/puerto-rico-aid-problem/index.html']]	3
['https://www.google.com/amp/s/amp.cnn.com/cnn/2017/09/30/politics/trump-tweets-puerto-rico-mayor/index.html#ampshare=http://www.cnn.com/2017/09/30/pol
itics/trump-tweets-puerto-rico-mayor/index.html']]	1
['https://www.google.com/amp/s/amp.cnn.com/cnn/2017/09/30/politics/trump-tweets-puerto-rico-mayor/index.html']]	1
['https://www.google.com/amp/s/amp.cnn.com/cnn/2017/09/30/us/puerto-rico-hurricane-recovery/index.html']]	2
['https://www.google.com/amp/s/amp.freep.com/amp/708496001#ampshare=http://www.freep.com/story/news/nation/2017/09/30/power-puerto-rico-detroit-solar-p
anels-hurricane-maria/708496001/']]	1
['https://www.google.com/amp/s/amp.theguardian.com/us-news/2016/sep/22/trump-ohio-campaign-chair-no-racism-before-obama']]	1
['https://www.google.com/amp/s/amp.theguardian.com/world/2010/sep/14/photographer-ernest-withers-fbi-informer']]	6
['https://www.google.com/amp/s/amp.theguardian.com/world/2017/sep/29/puerto-rico-crisis-supply-food-water']]	1
['https://www.google.com/amp/s/blog.bulletproof.com/womens-health-post-birth-control-syndrome-brain-injuries-dr-jolene-brighten-415/amp']]	1
['https://www.google.com/amp/s/citizensagainstequineslaughter.org/2017/07/07/congressman-chris-stewart-violates-u-s-c-title-18/amp/']]	3
['https://www.google.com/amp/s/drjengunter.wordpress.com/2016/04/12/check-your-privilege-and-your-facts-before-discussing-sex-selective-abortion/amp/']
]	3
['https://www.google.com/amp/s/hiddenremote.com/2017/09/28/wednesday-tv-ratings-seal-team-enjoys-solid-premiere-leads-night-total-viewers/amp/']]	1
['https://www.google.com/amp/s/mobile.nytimes.com/2016/07/27/us/puerto-rico-debt-mayors.amp.html',	1
['https://www.google.com/amp/s/mobile.nytimes.com/2017/09/27/nyregion/mexico-puerto-rico-diasters-aid-new-york-city.amp.html']]	1
['https://www.google.com/amp/s/nep.247sports.com/Bolt/Rex-Burkhead-returns-to-New-England-Patriots-practice-on-Friday--108159043/Amp']]	1
```
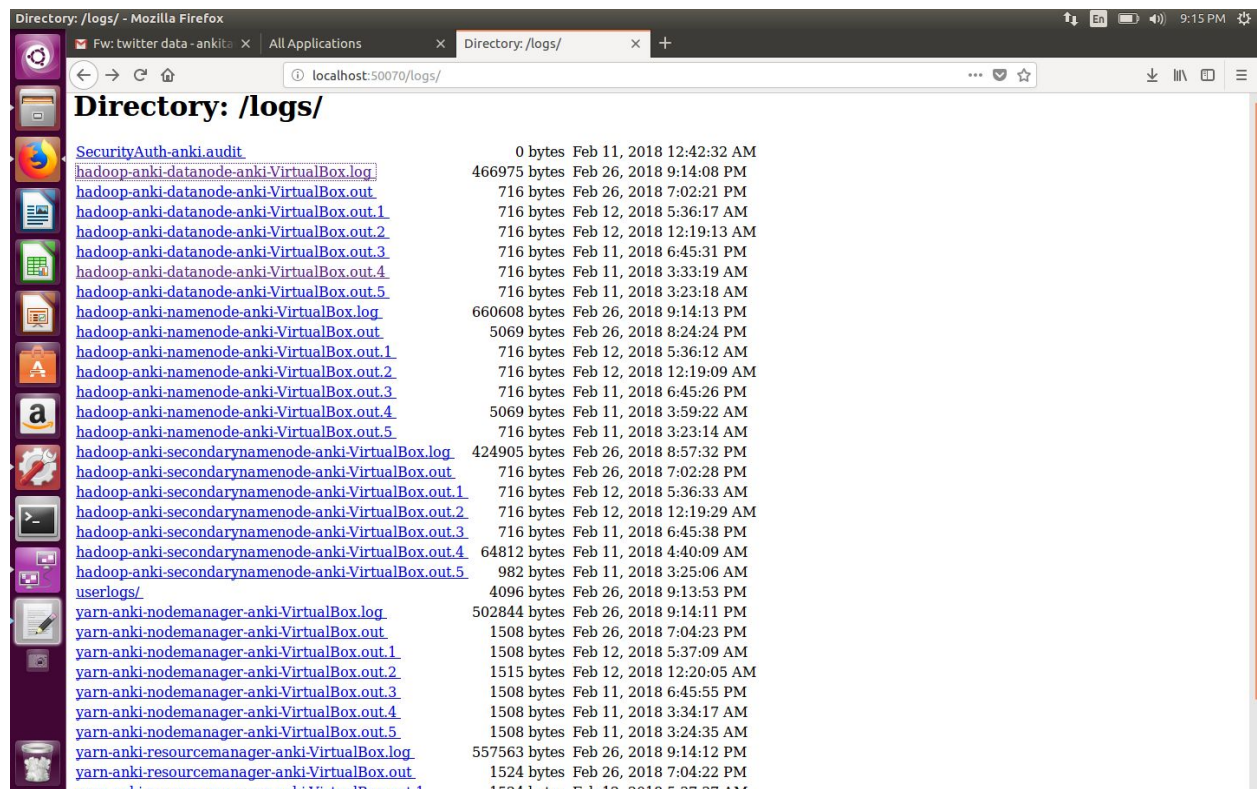
## Word Count of HashTags:



```
[['RocknRollerCoaster'],	1
[['RodmanDisgraceful'], 1
[['Rohingya',	7
[['Rohingya'], 3
[['RohingyaCrisis'],	1
[['RohingyaMuslims',	5
[['RohingyaMuslims'],	4
[['RohingyaRefugees'], 9
[['RohingyaTerrorExposed',	1
[['Rohingyas'], 8
[['Roku',	1
[['Role',	2
[['RollingStone',	1
[['RollyVortex'],	1
[['RomanticSuspense'], 1
[['Rome',	1
[['RoyMoore', 1
[['RuPaulsDragRace',	1
[['Ruan',	1
[['Ruckerfire'],	1
[['RussellWestbrook',	8
[['Russia',	8
[['Russia'], 4
[['Russian',	3
[['Russian'], 2
[['RussianBot', 1
[['RussianHacking'],	1
[['RyanGosling'],	1
[['RyderCup',	1
[['SAIGON',	2
[['SAN',	1
[['SATC',	1
[['SATC'],	1
[['SAVEJPWISHTOWN'],	2
[['SAVEMEDEDUCATION'], 1
[['SAvBAN'],	2
[['SBIR',	1
[['SC', 3
[['SCCs',	1
[['SCD2017'],	1
[['SCOTUS',	1
[['SCtop10',	1
[['SDLive',	1
[['SEAmayorLive'],	1
[['SEAtraffic'],	1
[['SEC',	2
[['SECHSKIES', 2
[['SEO',	14
[['SEO'],	5
[['SERP',	1
[['SHIITE'],	1
```

## HDFS Screenshot :



## Screenshot of log files created:

## Step 4: Running Word Count Program in Apache Spark:

## Screenshots of commands executed while running word count program:

1.Starting Spark: using bin/spark-shell command:



## 2. Loading the data:

3. Splitting content in our file as strings and applying map and reduce functions to perform word count, using the below highlighted command:



4. Finally collecting output which displays the word count in the form of an array of key value pairs:

### References:

http://adilmoujahid.com/posts/2014/07/twitter-analytics/

https://github.com/heerme/twitter-topics/blob/master/extract-json-to-text-stream.py

https://www.youtube.com/watch?v=YZnNb0BTrS4&list=PLJNKKS4iwuamrvNVahopRzi
urK7XNCc5B&index=3