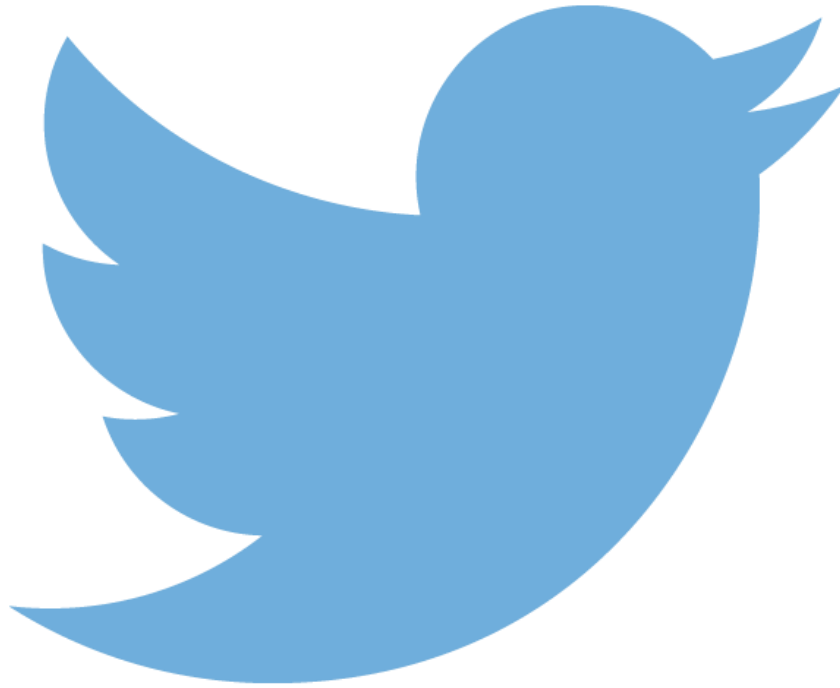


Principles of Big Data Management

Project 2



Team

1. Saketh Garuda (sg7kf)
2. Mudunuri Sri Sai Sarat Chandra Varma (smyx4)
3. Yalamanchili Sowmya (syb7c)
4. Nandanamudi Sreelakshmi (snhnc)

1.Introduction

The main objective of this project is analyzing the twitter data about 'Politics' category and analyze the data using Hadoop MapReduce. After collecting, the necessary tweets data we have used python code to retrieve text from the tweets. Based on the text we have calculated the list of words in the tweets text file which includes the list of duplicates and unique words in the tweets text file using MapReduce program and also calculated the ratio of number of unique words and to the number of duplicate words. We also implemented MapReduce program to return the top ten best times to post a tweet on twitter.

The entire document gives the walk through and scope of the working environment of the project.

2.System Requirements

Software Requirements:

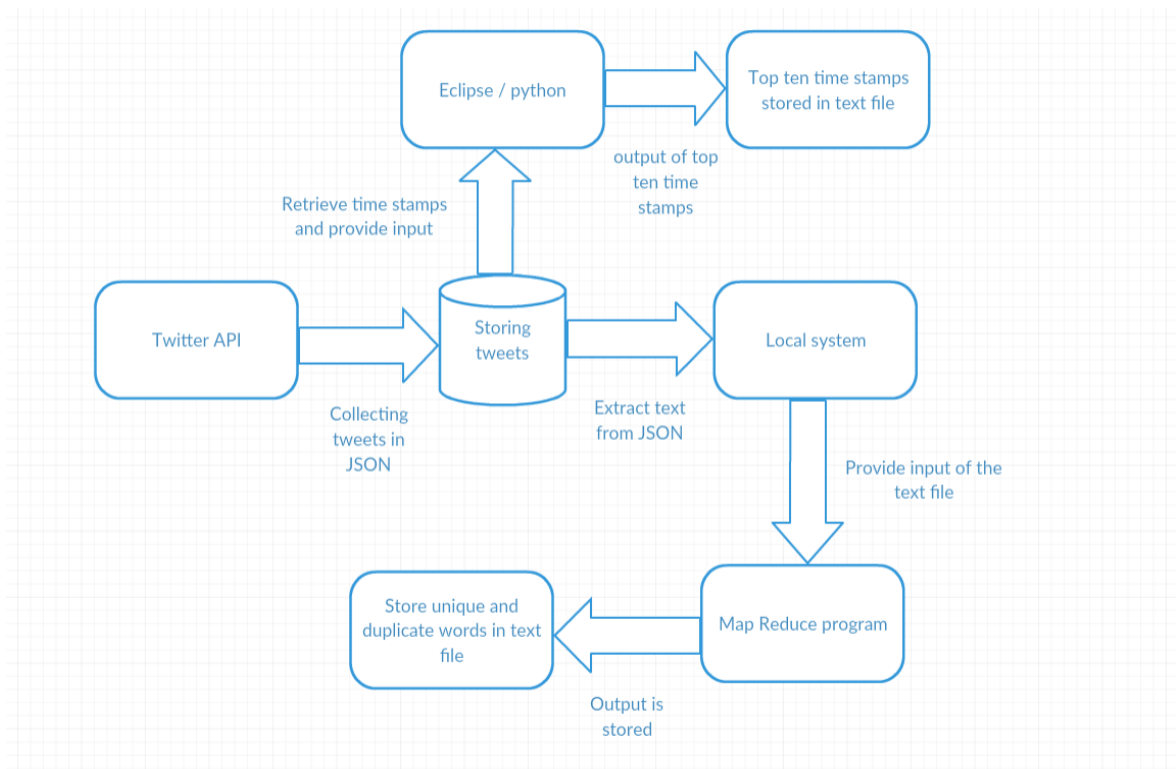
- Eclipse
- Net Beans
- Apache Hadoop V 2.6.0-cdh5.8.0
- JDK 1.8

Programming Languages:

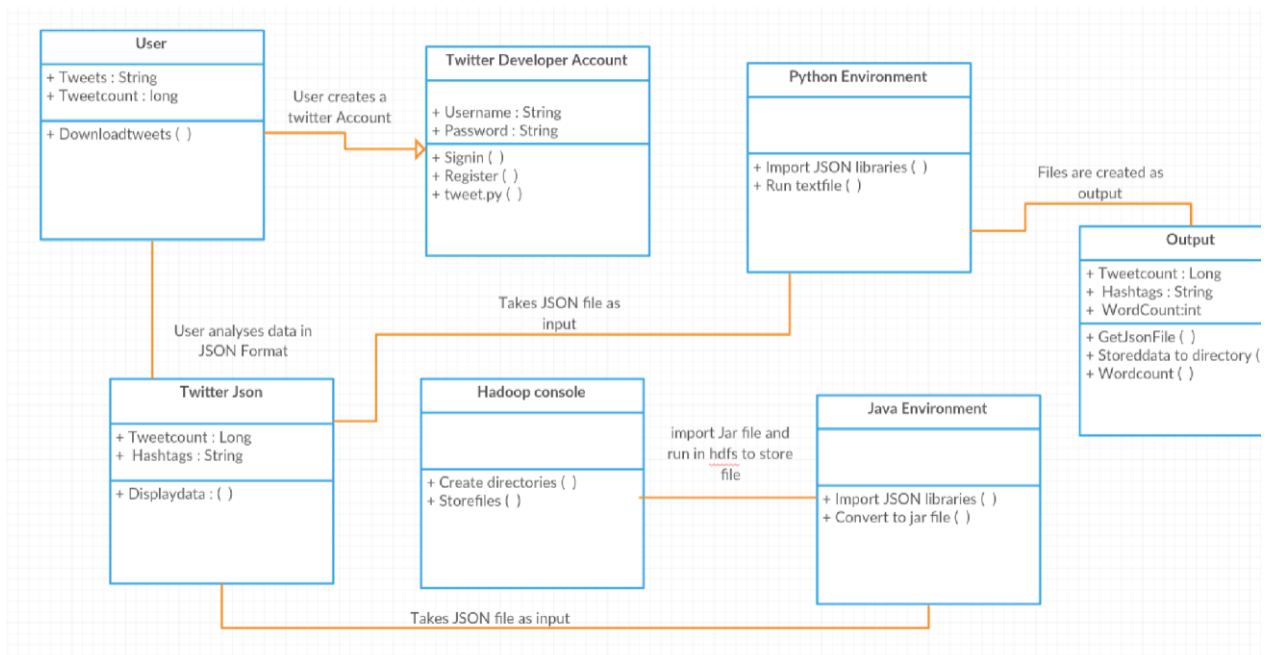
- Java
- Python

3. Software Architecture and Design

- Twitter tweets are extracted using the API by running python program using the twitter tokens. The collected tweets are stored in JSON format in the local system.



- Here is the class diagram for the project



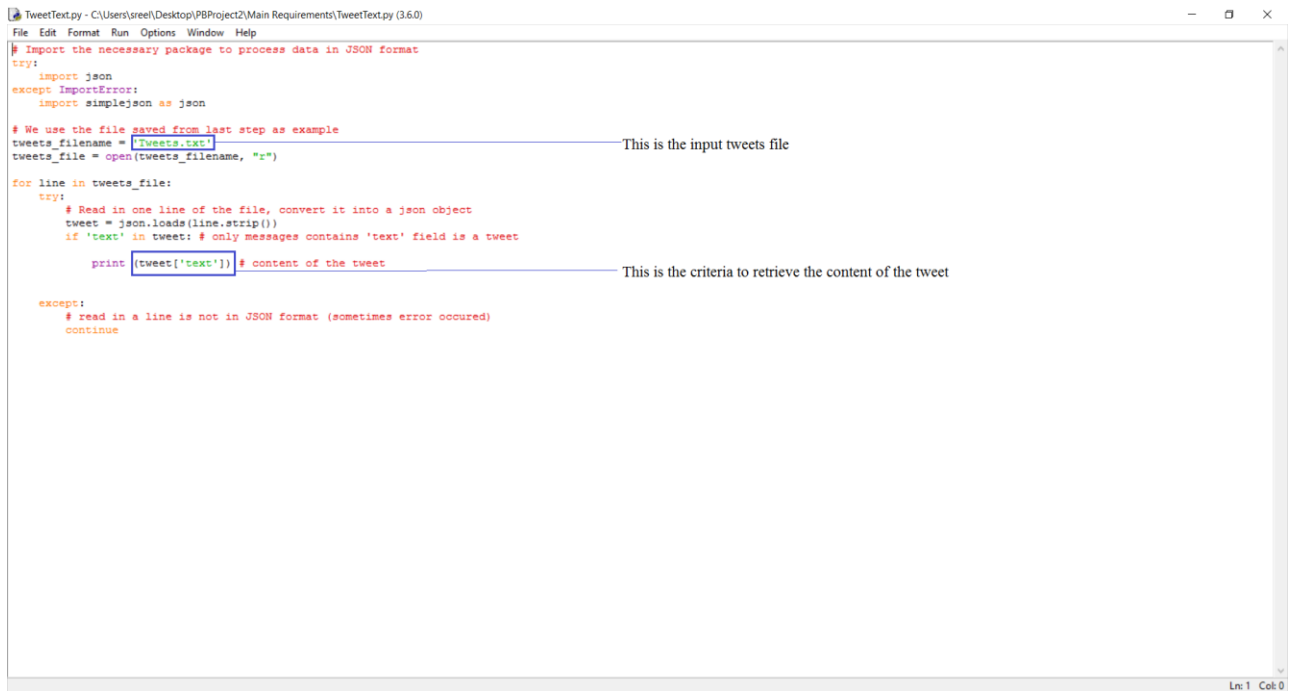
4. Main Requirements

Collect tweets in JavaScript Object Notation

- We've collected 100k+ tweets related to category 'Politics' into a JSON file using the twitter API. JSON file is used for future work in retrieving the text and top 10 best times from the entire list of tweets. Here is the screenshot of JSON file,

```
{
  "created_at": "Tue Feb 28 04:46:48 +0000 2017",
  "id": "836437472527724547",
  "id_str": "836437472527724547",
  "text": "Agreed, make it like Desi po",
  "created_at": "Tue Feb 28 04:46:48 +0000 2017",
  "id": "836437474628927488",
  "id_str": "836437474628927488",
  "text": "RT @bianconoco: \ufffc\u7c60",
  "created_at": "Tue Feb 28 04:46:48 +0000 2017",
  "id": "836437474981314560",
  "id_str": "836437474981314560",
  "text": "RT @englishpaulm: I'm donati",
  "created_at": "Tue Feb 28 04:46:49 +0000 2017",
  "id": "836437476176637953",
  "id_str": "836437476176637953",
  "text": "RT @owillis: if i have my fa",
  "created_at": "Tue Feb 28 04:46:49 +0000 2017",
  "id": "836437476206080000",
  "id_str": "836437476206080000",
  "text": "RT @clizyatin: This Beautif",
  "created_at": "Tue Feb 28 04:46:49 +0000 2017",
  "id": "836437477028159489",
  "id_str": "836437477028159489",
  "text": "RT @SitaramYechury: All effo",
  "created_at": "Tue Feb 28 04:46:50 +0000 2017",
  "id": "836437480215760896",
  "id_str": "836437480215760896",
  "text": "https://t.co/vwLw0yJYDf",
  "created_at": "Tue Feb 28 04:46:50 +0000 2017",
  "id": "836437480647831552",
  "id_str": "836437480647831552",
  "text": "RT @BernieSanders: Trump mus",
  "created_at": "Tue Feb 28 04:46:50 +0000 2017",
  "id": "836437483575353344",
  "id_str": "836437483575353344",
  "text": "RT @GabrielAlon: Even though",
  "created_at": "Tue Feb 28 04:46:51 +0000 2017",
  "id": "836437485450207232",
  "id_str": "836437485450207232",
  "text": "RT @AngelaRoss1: @thehill Di",
  "created_at": "Tue Feb 28 04:46:51 +0000 2017",
  "id": "836437485530021888",
  "id_str": "836437485530021888",
  "text": "RT @stevenjfrisch: Careless",
  "created_at": "Tue Feb 28 04:46:52 +0000 2017",
  "id": "836437488839213056",
  "id_str": "836437488839213056",
  "text": "RT @DrDinD: 8 people who owe",
  "created_at": "Tue Feb 28 04:46:52 +0000 2017",
  "id": "836437488843501568",
  "id_str": "836437488843501568",
  "text": "RT @The latest LiveWord? Politic",
  "created_at": "Tue Feb 28 04:46:52 +0000 2017",
  "id": "836437489904578561",
  "id_str": "836437489904578561",
  "text": "RT @koufukuron1925:\u3010\u051",
  "created_at": "Tue Feb 28 04:46:53 +0000 2017",
  "id": "836437493431967744",
  "id_str": "836437493431967744",
  "text": "RT @debasishdas568: https://",
  "created_at": "Tue Feb 28 04:46:53 +0000 2017",
  "id": "836437493935292417",
  "id_str": "836437493935292417",
  "text": "RT @debasishdas568: https://",
  "created_at": "Tue Feb 28 04:46:53 +0000 2017",
  "id": "836437493935304706",
  "id_str": "836437493935304706",
  "text": "RT @RealM_Zubair: All Pakist",
  "created_at": "Tue Feb 28 04:46:53 +0000 2017",
  "id": "836437493935304704",
  "id_str": "836437493935304704",
  "text": "RT @debasishdas568: https://",
  "created_at": "Tue Feb 28 04:46:53 +0000 2017",
  "id": "836437494577025024",
  "id_str": "836437494577025024",
  "text": "RT @peterbakernyt: Bush 43 i",
  "created_at": "Tue Feb 28 04:46:53 +0000 2017",
  "id": "836437494493294593",
  "id_str": "836437494493294593",
  "text": "RT @MajorPoonia: Vultures fr",
  "created_at": "Tue Feb 28 04:46:54 +0000 2017",
  "id": "836437496892370944",
  "id_str": "836437496892370944",
  "text": "RT @funder: Anti-Trump Group",
  "created_at": "Tue Feb 28 04:46:55 +0000 2017",
  "id": "836437500616929280",
  "id_str": "836437500616929280",
  "text": "RT @Live: Justice Minister confi",
  "created_at": "Tue Feb 28 04:46:54 +0000 2017",
  "id": "836437499648032770",
  "id_str": "836437499648032770",
  "text": "RT @Pres. Trump suggests Obama w",
  "created_at": "Tue Feb 28 04:46:55 +0000 2017",
  "id": "836437501090926593",
  "id_str": "836437501090926593",
  "text": "RT @DineshGhodke: Why has po",
  "created_at": "Tue Feb 28 04:46:55 +0000 2017",
  "id": "836437501455831044",
  "id_str": "836437501455831044",
  "text": "RT @The Politics of Black Mascu",
  "created_at": "Tue Feb 28 04:46:55 +0000 2017",
  "id": "836437501866729472",
  "id_str": "836437501866729472",
  "text": "RT @AnnCoulter: Trump critic",
  "created_at": "Tue Feb 28 04:46:55 +0000 2017",
  "id": "836437501384474624",
  "id_str": "836437501384474624",
  "text": "RT @He is following his father's",
  "created_at": "Tue Feb 28 04:46:55 +0000 2017",
  "id": "836437502101639168",
  "id_str": "836437502101639168",
  "text": "RT @sandeepfromvms: @mehartw",
  "created_at": "Tue Feb 28 04:46:55 +0000 2017",
  "id": "836437503984926721",
  "id_str": "836437503984926721",
  "text": "RT @bemoredev: \ud45c\ucc3d\u",
  "created_at": "Tue Feb 28 04:46:56 +0000 2017",
  "id": "836437506082103296",
  "id_str": "836437506082103296",
  "text": "RT @SivaswamiPK: @mediacrook",
  "created_at": "Tue Feb 28 04:46:56 +0000 2017",
  "id": "8364375080857612288",
  "id_str": "8364375080857612288",
  "text": "RT @There is no good reason to m",
  "created_at": "Tue Feb 28 04:46:57 +0000 2017",
  "id": "836437510628622336",
  "id_str": "836437510628622336",
  "text": "RT @malviyami: All credit t",
  "created_at": "Tue Feb 28 04:46:57 +0000 2017",
  "id": "836437511262126080",
  "id_str": "836437511262126080",
  "text": "RT @RepJoeCourtney: It would",
  "created_at": "Tue Feb 28 04:46:58 +0000 2017",
  "id": "836437514348945409",
  "id_str": "836437514348945409",
  "text": "RT @BernieSanders: Trump mus",
  "created_at": "Tue Feb 28 04:46:58 +0000 2017",
  "id": "836437515670282240",
  "id_str": "836437515670282240",
  "text": "RT @BKNight561: Muslims aren't t",
  "created_at": "Tue Feb 28 04:46:58 +0000 2017",
  "id": "836437516404195328",
  "id_str": "836437516404195328",
  "text": "RT @AddityaRajKaul: Says, Mr.",
  "created_at": "Tue Feb 28 04:46:59 +0000 2017",
  "id": "836437517226254336",
  "id_str": "836437517226254336",
  "text": "RT @oh no they did the politics",
  "created_at": "Tue Feb 28 04:46:59 +0000 2017",
  "id": "836437520963416064",
  "id_str": "836437520963416064",
  "text": "RT @TonyKrvaric: Shame on C1",
  "created_at": "Tue Feb 28 04:47:00 +0000 2017",
  "id": "836437523266220032",
  "id_str": "836437523266220032",
  "text": "RT @lany891: #TuckerCarlson",
  "created_at": "Tue Feb 28 04:47:01 +0000 2017",
  "id": "836437526101426177",
  "id_str": "836437526101426177",
  "text": "RT @Ryan Is Really Looking For H",
  "created_at": "Tue Feb 28 04:47:01 +0000 2017",
  "id": "836437526738919424",
  "id_str": "836437526738919424",
  "text": "RT @S/VX primarily on local nati",
  "created_at": "Tue Feb 28 04:47:02 +0000 2017",
  "id": "836437529955983360",
  "id_str": "836437529955983360",
  "text": "RT @SyahmiBrahim: Carefully",
  "created_at": "Tue Feb 28 04:47:02 +0000 2017",
  "id": "836437531705028608",
  "id_str": "836437531705028608",
  "text": "RT @dingos1946: Yes Abetz 5",
  "created_at": "Tue Feb 28 04:47:02 +0000 2017",
  "id": "836437531755483136",
  "id_str": "836437531755483136",
  "text": "RT @Namrataa: So she is in",
  "created_at": "Tue Feb 28 04:47:02 +0000 2017",
  "id": "836437531818328064",
  "id_str": "836437531818328064",
  "text": "RT @funder: Chaffetz abusing",
  "created_at": "Tue Feb 28 04:47:03 +0000 2017",
  "id": "836437534028652545",
  "id_str": "836437534028652545",
  "text": "RT @JohnWren1950: BREAKING:",
  "created_at": "Tue Feb 28 04:47:03 +0000 2017",
  "id": "836437536125796352",
  "id_str": "836437536125796352",
  "text": "RT @Hood_Biologist: One day
```

- After collecting the 100k tweets, we've written a python code for separating text from tweets file and it is stored in tweet text file.
- Below screenshot indicates the python code for retrieving text from the tweets file. We have given the input “**Tweets.txt**” file and parsed the JSON data with “**tweet[‘text’]**” and stored it an output file.



```

TweetText.py - C:\Users\sree\Desktop\PBProject2\Main Requirements\TweetText.py (3.6.0)
File Edit Format Run Options Window Help
# Import the necessary package to process data in JSON format
try:
    import json
except ImportError:
    import simplejson as json

# We use the file saved from last step as example
tweets_filename = 'Tweets.txt'
tweets_file = open(tweets_filename, "r")

for line in tweets_file:
    try:
        # Read in one line of the file, convert it into a json object
        tweet = json.loads(line.strip())
        if 'text' in tweet: # only messages contains 'text' field is a tweet
            print (tweet['text']) # content of the tweet

    except:
        # read in a line is not in JSON format (sometimes error occurred)
        continue

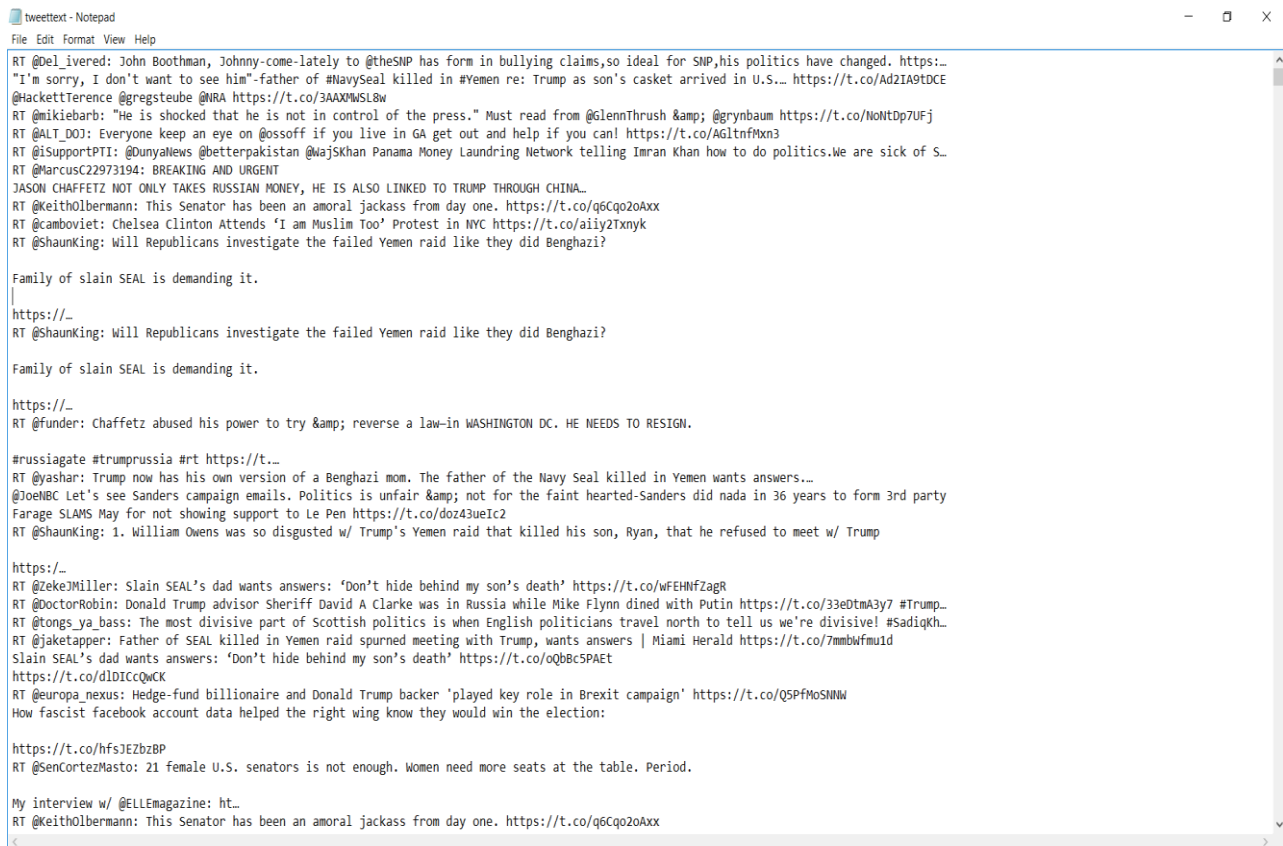
```

This is the input tweets file

This is the criteria to retrieve the content of the tweet

Ln: 1 Col: 0

- Here is the screenshot of output TweetsText file



```
tweettext - Notepad
File Edit Format View Help
RT @Del_ivered: John Boothman, Johnny-come-lately to @theSNP has form in bullying claims,so ideal for SNP,his politics have changed. https://t.co/Ad2IA9tDCE
"I'm sorry, I don't want to see him"-father of #NavySeal killed in #Yemen re: Trump as son's casket arrived in U.S... https://t.co/Ad2IA9tDCE
@HackettTerence @gregsteube @NRA https://t.co/3AAXWMSL8w
RT @mikielbarb: "He is shocked that he is not in control of the press." Must read from @GlennThrush & @grynbaum https://t.co/NoMtDp7Ufj
RT @ALT_DOJ: Everyone keep an eye on @ossoff if you live in GA get out and help if you can! https://t.co/AGltnfMxn3
RT @ISupportPTI: @DunyaNews @betterpakistan @WajSKhan Panama Money Laundering Network telling Imran Khan how to do politics.We are sick of S...
RT @MarcusC22973194: BREAKING AND URGENT
JASON CHAFFETZ NOT ONLY TAKES RUSSIAN MONEY, HE IS ALSO LINKED TO TRUMP THROUGH CHINA...
RT @KeithOlbermann: This Senator has been an amoral jackass from day one. https://t.co/q6Cqo2oAxx
RT @Camboviet: Chelsea Clinton Attends 'I am Muslim Too' Protest in NYC https://t.co/aaiy2Txnyk
RT @Shaunking: Will Republicans investigate the failed Yemen raid like they did Benghazi?

Family of slain SEAL is demanding it.

https://_
RT @Shaunking: Will Republicans investigate the failed Yemen raid like they did Benghazi?

Family of slain SEAL is demanding it.

https://_
RT @funder: Chaffetz abused his power to try & reverse a law-in WASHINGTON DC. HE NEEDS TO RESIGN.

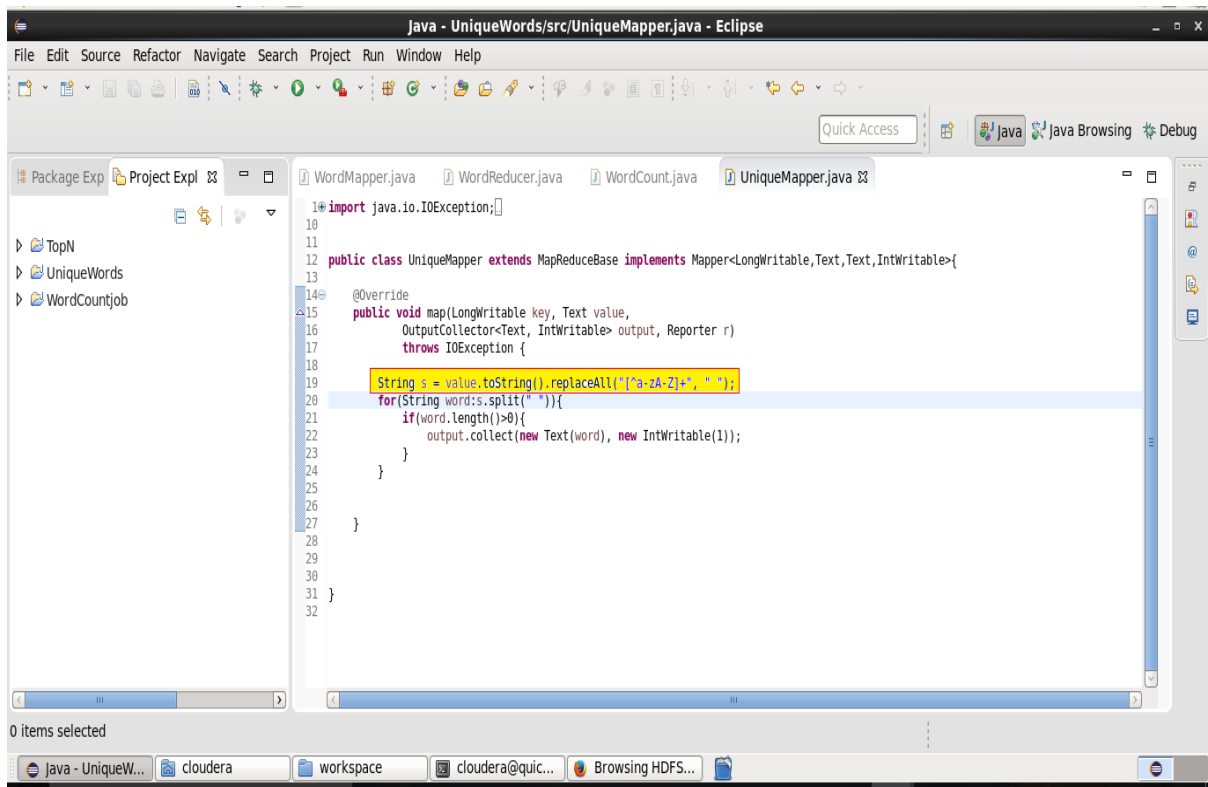
#russiagate #trumprussia #rt https://t...
RT @yashar: Trump now has his own version of a Benghazi mom. The father of the Navy Seal killed in Yemen wants answers...
@JoelHBC Let's see Sanders campaign emails. Politics is unfair & not for the faint hearted-Sanders did nada in 36 years to form 3rd party
Farage SLAMS May for not showing support to Le Pen https://t.co/doz43ueIc2
RT @Shaunking: 1. William Owens was so disgusted w/ Trump's Yemen raid that killed his son, Ryan, that he refused to meet w/ Trump

https://_
RT @ZekeJMiller: Slain SEAL's dad wants answers: 'Don't hide behind my son's death' https://t.co/wFEHnfZagR
RT @DoctorRobin: Donald Trump advisor Sheriff David A Clarke was in Russia while Mike Flynn dined with Putin https://t.co/33e0tmA3y7 #Trump...
RT @tongs_ya_bass: The most divisive part of Scottish politics is when English politicians travel north to tell us we're divisive! #SadiqKh...
RT @jaketapper: Father of SEAL killed in Yemen raid spurned meeting with Trump, wants answers | Miami Herald https://t.co/7mmbwfmui1
Slain SEAL's dad wants answers: 'Don't hide behind my son's death' https://t.co/oQbC5PAEt
https://t.co/dlDlCCQwCK
RT @europa_nexus: Hedge-fund billionaire and Donald Trump backer 'played key role in Brexit campaign' https://t.co/Q5PFMoSNNW
How fascist facebook account data helped the right wing know they would win the election:

https://t.co/hfsJEZbzBP
RT @SenCortezMasto: 21 female U.S. senators is not enough. Women need more seats at the table. Period.

My interview w/ @ELLEmagazine: ht...
RT @KeithOlbermann: This Senator has been an amoral jackass from day one. https://t.co/q6Cqo2oAxx
```

- We have implemented the Hadoop MapReduce in java code to find the list of duplicate and unique words. The words that occurred more than once are considered as duplicate words and the words which occurred once are considered as unique.
- We have filtered words in the text file using the below code which removes the special characters and split them with whitespace delimitation. The below screenshot refers the same:



```
Java - UniqueWords/src/UniqueMapper.java - Eclipse
File Edit Source Refactor Navigate Search Project Run Window Help

Package Exp Project Expl WordMapper.java WordReducer.java WordCount.java UniqueMapper.java
TopN
UniqueWords
WordCountJob

import java.io.IOException;
10
11
12 public class UniqueMapper extends MapReduceBase implements Mapper<LongWritable,Text,Text,IntWritable>{
13
14 @Override
15 public void map(LongWritable key, Text value,
16                 OutputCollector<Text, IntWritable> output, Reporter r)
17                 throws IOException {
18
19     String s = value.toString().replaceAll("[^a-zA-Z]+", "");
20     for(String word:s.split(" ")){
21         if(word.length()>0){
22             output.collect(new Text(word), new IntWritable(1));
23         }
24     }
25 }
26
27 }
28
29
30
31 }
32

0 items selected
Java - UniqueW... cloudera workspace cloudera@quic... Browsing HDFS...
```

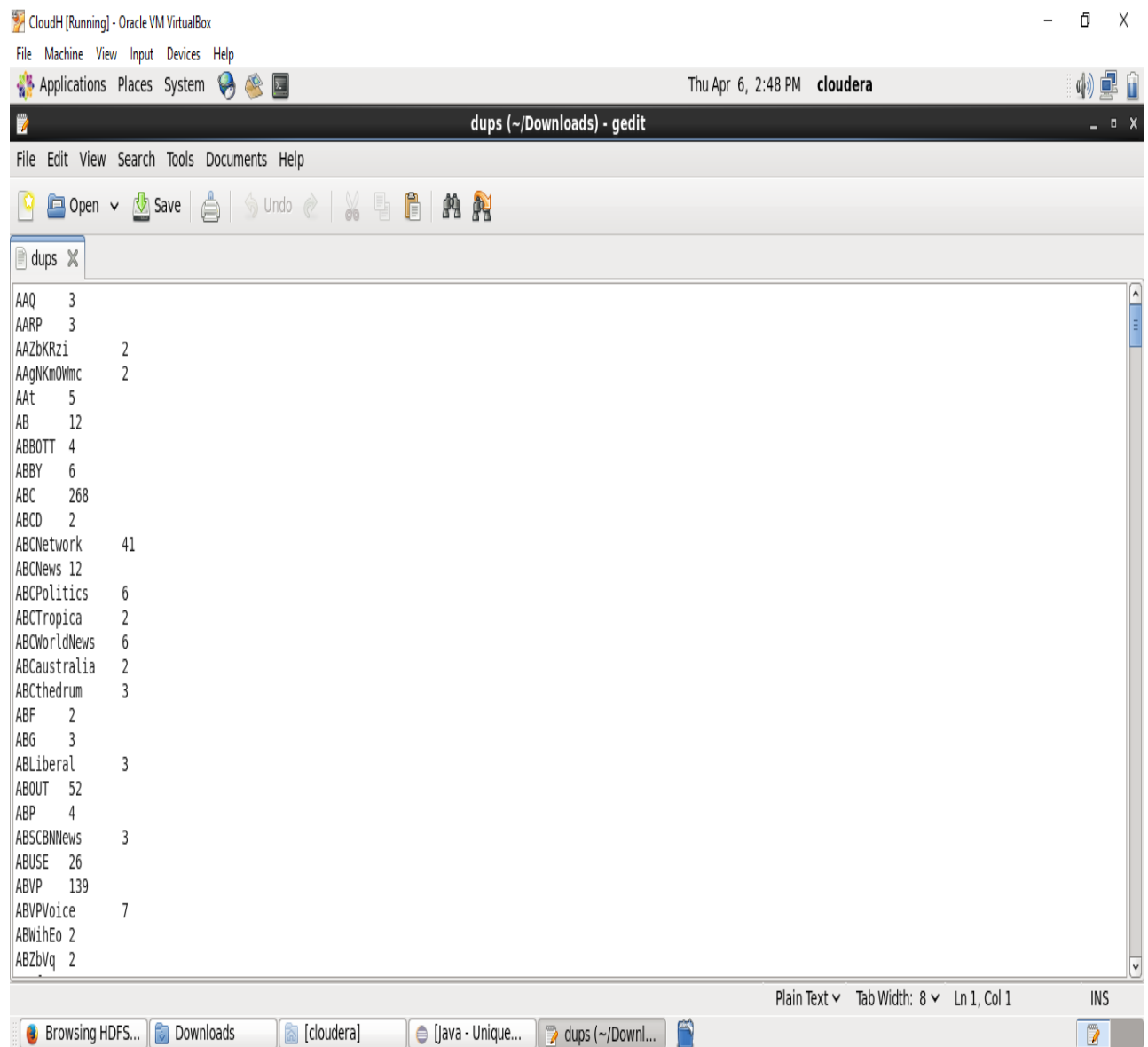
For Duplicate Words:

We have exported the java MapReduce program to jar file and run in the Hadoop environment to get the list of duplicate words. By using the command:

```
hadoop jar Duplicates.jar WordCount /user/cloudera/test/tweettext.txt /user/cloudera/test/dups.txt
```

The above hadoop command takes the Duplicates jar file with the class name that has main method i.e **WordCount** with input path “/user/cloudera/test/tweettext.txt” and output path “/user/cloudera/test/dups.txt”.

Below screenshot shows the list of duplicate words which are stored in the dups text file the count on right of the words shows that they are duplicated as they occurred more than once.

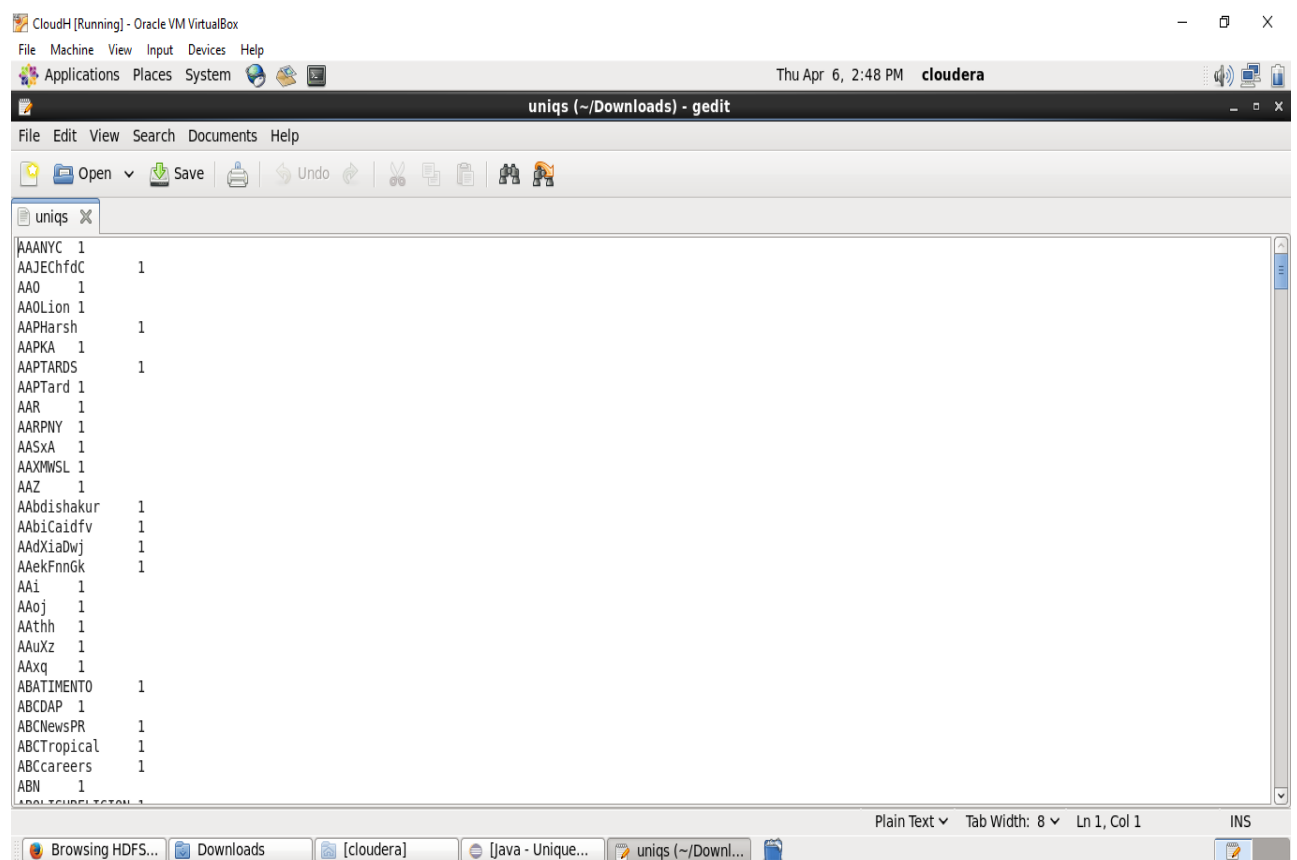


For Unique Words:

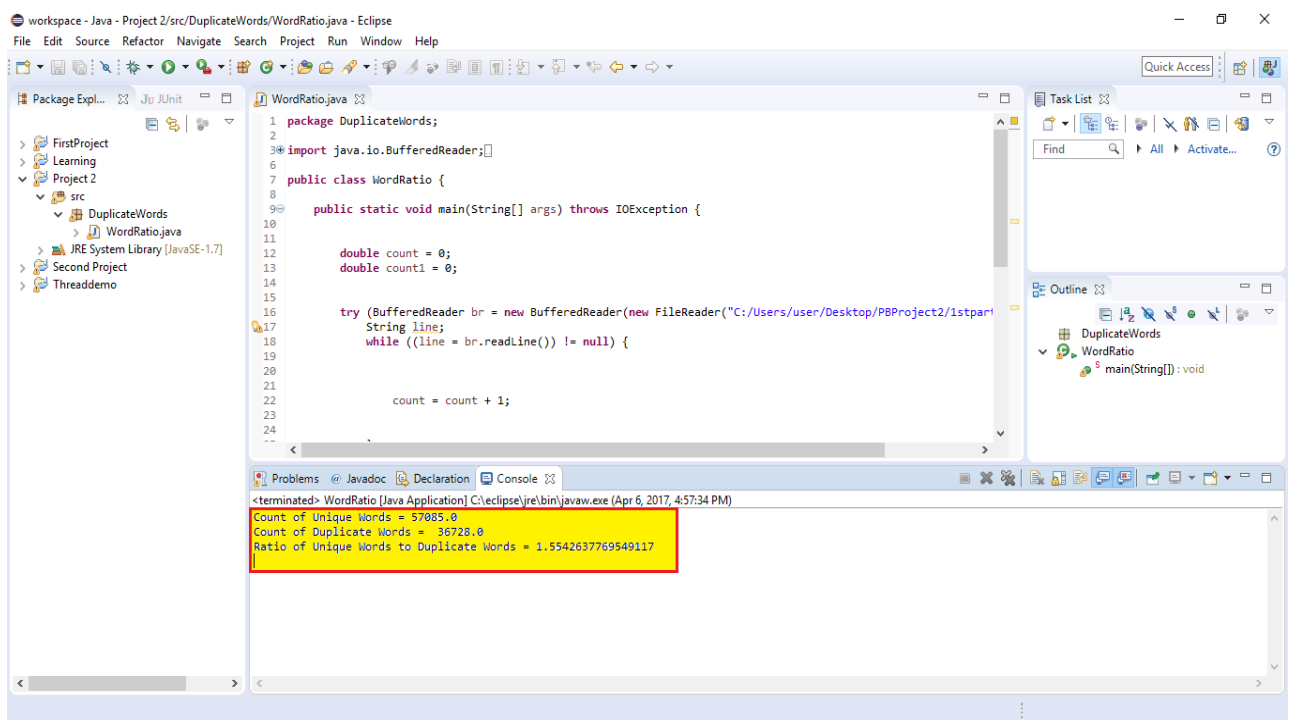
We have exported the java MapReduce program to jar file and run in the Hadoop environment to get the list of duplicate words. By using the command:

```
hadoop jar Uniques.jar UniqueCount /user/cloudera/test/tweettext.txt /user/cloudera/test/uniqs.txt
```

Below screenshot shows the list of unique words which are stored in the uniqs text file and the count on the right side of the words shows that the words are unique as they occurred only once.



We have used java code to calculate the ratio of the number of duplicate words and unique words, by taking uniqs and dups text files as input and calculated the **Unique Words: 57085.0** and **Duplicate Words 36728.0** and also the **Ratio of Unique Words to Duplicate Words 1.554** and the screenshot below depicts the same.



The screenshot displays the Eclipse IDE interface. The main editor shows the `WordRatio.java` file with the following code:

```
1 package DuplicateWords;
2
3 import java.io.BufferedReader;
4
5
6 public class WordRatio {
7
8     public static void main(String[] args) throws IOException {
9
10
11         double count = 0;
12         double count1 = 0;
13
14
15         try (BufferedReader br = new BufferedReader(new FileReader("C:/Users/user/Desktop/PBProject2/1stpart
16 String line;
17         while ((line = br.readLine()) != null) {
18
19
20
21             count = count + 1;
22
23
24
25 }
```

The console output at the bottom shows the results of the program execution:

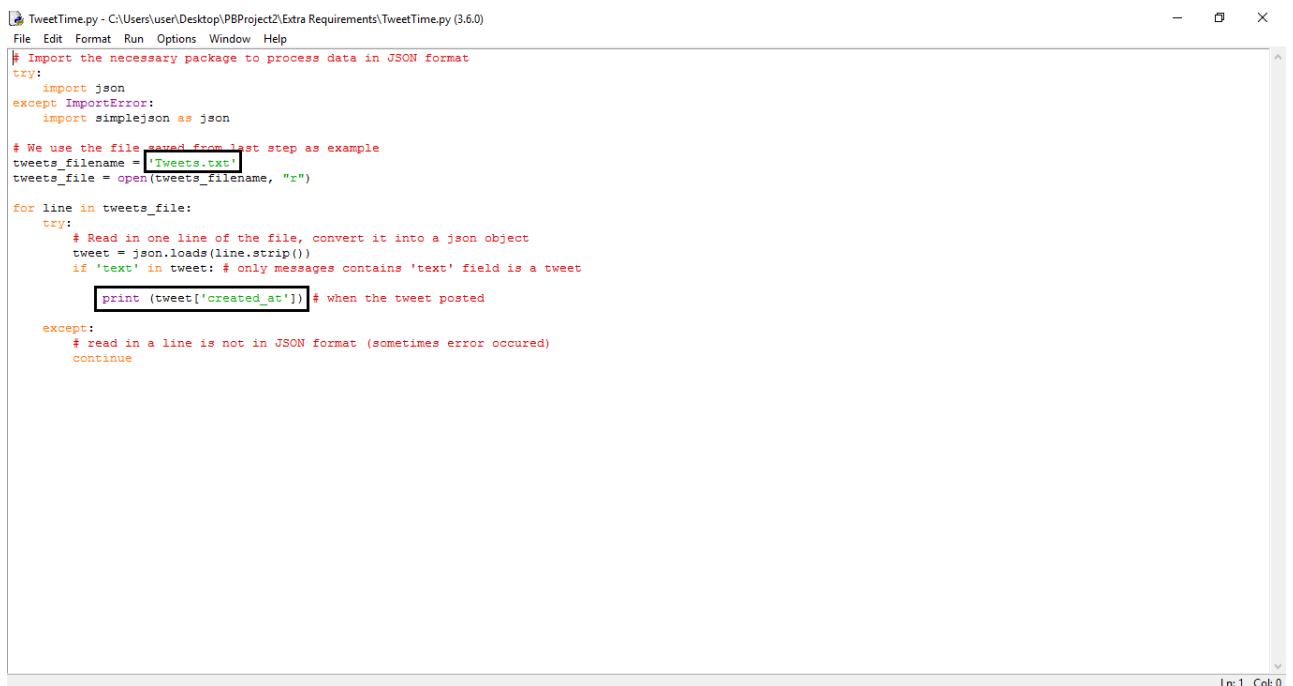
```
<terminated> WordRatio [Java Application] C:\eclipse\jre\bin\javaw.exe (Apr 6, 2017, 4:57:34 PM)
Count of Unique words = 57085.0
Count of Duplicate Words = 36728.0
Ratio of Unique Words to Duplicate Words = 1.5542637769549117
```

5. EXTRA REQUIREMENT:

We have also done the extra requirement to calculate the top ten best times to post a tweet on a twitter using java MapReduce program.

As first step we have retrieved the timestamp from tweets using python code and the output is stored in the TimeStamp text file

The below python code shows “Tweets.txt” is the input file and it prints out the timestamp of the tweets creation and stores it in the TimeStamp text file.



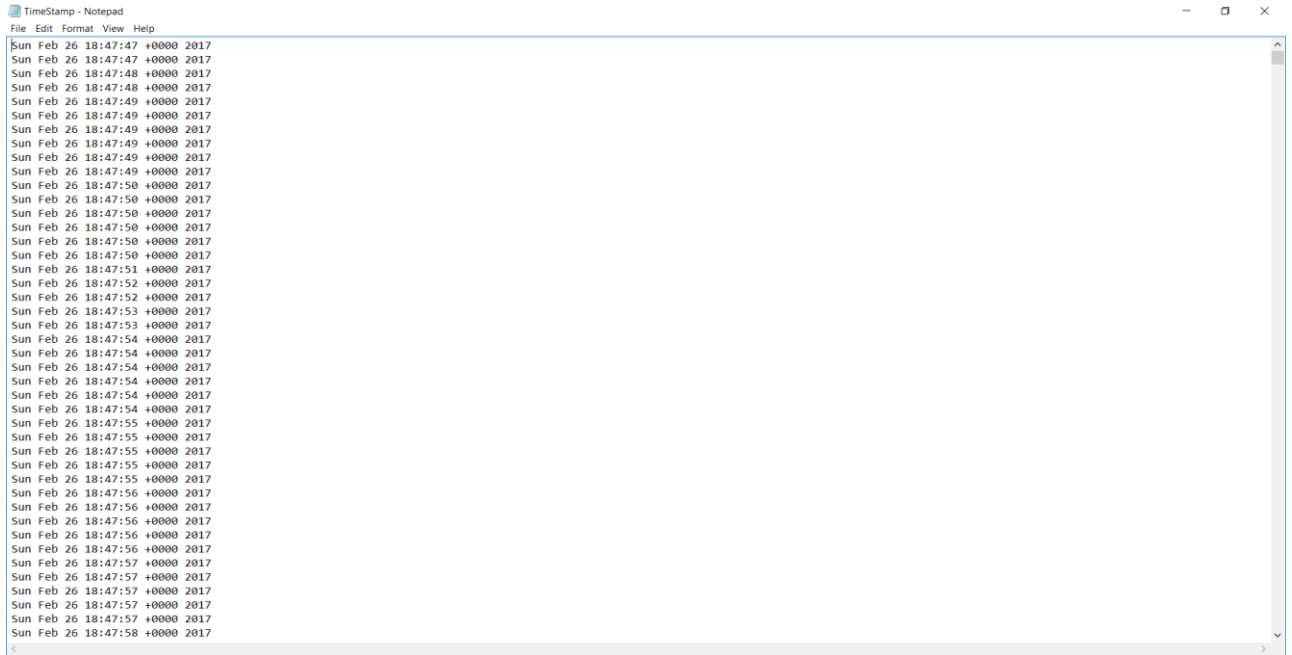
```
TweetTime.py - C:\Users\user\Desktop\PBProject2\Extra Requirements\TweetTime.py (3.6.0)
File Edit Format Run Options Window Help
# Import the necessary package to process data in JSON format
try:
    import json
except ImportError:
    import simplejson as json

# We use the file saved from last step as example
tweets_filename = 'Tweets.txt'
tweets_file = open(tweets_filename, "r")

for line in tweets_file:
    try:
        # Read in one line of the file, convert it into a json object
        tweet = json.loads(line.strip())
        if 'text' in tweet: # only messages contains 'text' field is a tweet
            print (tweet['created_at']) # when the tweet posted
    except:
        # read in a line is not in JSON format (sometimes error occurred)
        continue
```

Ln: 1 Col: 0

TimeStamp Text file:



```
TimeStamp - Notepad
File Edit Format View Help
Sun Feb 26 18:47:47 +0000 2017
Sun Feb 26 18:47:47 +0000 2017
Sun Feb 26 18:47:48 +0000 2017
Sun Feb 26 18:47:48 +0000 2017
Sun Feb 26 18:47:49 +0000 2017
Sun Feb 26 18:47:49 +0000 2017
Sun Feb 26 18:47:49 +0000 2017
Sun Feb 26 18:47:49 +0000 2017
Sun Feb 26 18:47:49 +0000 2017
Sun Feb 26 18:47:49 +0000 2017
Sun Feb 26 18:47:50 +0000 2017
Sun Feb 26 18:47:50 +0000 2017
Sun Feb 26 18:47:50 +0000 2017
Sun Feb 26 18:47:50 +0000 2017
Sun Feb 26 18:47:50 +0000 2017
Sun Feb 26 18:47:50 +0000 2017
Sun Feb 26 18:47:51 +0000 2017
Sun Feb 26 18:47:52 +0000 2017
Sun Feb 26 18:47:52 +0000 2017
Sun Feb 26 18:47:53 +0000 2017
Sun Feb 26 18:47:53 +0000 2017
Sun Feb 26 18:47:54 +0000 2017
Sun Feb 26 18:47:54 +0000 2017
Sun Feb 26 18:47:54 +0000 2017
Sun Feb 26 18:47:54 +0000 2017
Sun Feb 26 18:47:54 +0000 2017
Sun Feb 26 18:47:55 +0000 2017
Sun Feb 26 18:47:55 +0000 2017
Sun Feb 26 18:47:55 +0000 2017
Sun Feb 26 18:47:55 +0000 2017
Sun Feb 26 18:47:55 +0000 2017
Sun Feb 26 18:47:55 +0000 2017
Sun Feb 26 18:47:56 +0000 2017
Sun Feb 26 18:47:56 +0000 2017
Sun Feb 26 18:47:56 +0000 2017
Sun Feb 26 18:47:56 +0000 2017
Sun Feb 26 18:47:56 +0000 2017
Sun Feb 26 18:47:57 +0000 2017
Sun Feb 26 18:47:57 +0000 2017
Sun Feb 26 18:47:57 +0000 2017
Sun Feb 26 18:47:57 +0000 2017
Sun Feb 26 18:47:58 +0000 2017
```

To find the number of tweets created at the particular time:

We have exported the java MapReduce program to jar file and run in the Hadoop environment to get the list of the timestamps. By using the below hadoop command the java map reduce program will run and gives the below output file with count of the timestamps.

```
hadoop jar TimeSort.jar TimeCount /user/cloudera/test/TimeStamp.txt
/user/cloudera/test/TimeStampOutput1.txt
```

The below screen shots shows the output after running the above code and it gives number of tweets that occurred at that particular time. The count on the right side of the timestamp shows the same.

```
1 Mon Feb 27 19:25:57 +0000 2017 2
2 Mon Feb 27 19:25:58 +0000 2017 6
3 Mon Feb 27 19:26:00 +0000 2017 5
4 Mon Feb 27 19:26:01 +0000 2017 4
5 Mon Feb 27 19:26:02 +0000 2017 2
6 Mon Feb 27 19:26:03 +0000 2017 9
7 Mon Feb 27 19:26:04 +0000 2017 4
8 Mon Feb 27 19:26:05 +0000 2017 6
9 Mon Feb 27 19:26:06 +0000 2017 4
10 Mon Feb 27 19:26:07 +0000 2017 5
11 Mon Feb 27 19:26:09 +0000 2017 5
12 Mon Feb 27 19:26:10 +0000 2017 6
13 Mon Feb 27 19:26:11 +0000 2017 5
14 Mon Feb 27 19:26:12 +0000 2017 6
15 Mon Feb 27 19:26:13 +0000 2017 3
16 Mon Feb 27 19:26:14 +0000 2017 4
17 Mon Feb 27 19:26:15 +0000 2017 3
18 Mon Feb 27 19:26:17 +0000 2017 8
19 Mon Feb 27 19:26:18 +0000 2017 3
20 Mon Feb 27 19:26:19 +0000 2017 2
21 Mon Feb 27 19:26:20 +0000 2017 4
22 Mon Feb 27 19:26:21 +0000 2017 2
23 Mon Feb 27 19:26:22 +0000 2017 4
24 Mon Feb 27 19:26:23 +0000 2017 7
25 Mon Feb 27 19:26:24 +0000 2017 4
26 Mon Feb 27 19:26:25 +0000 2017 4
27 Mon Feb 27 19:26:26 +0000 2017 9
28 Mon Feb 27 19:26:27 +0000 2017 7
29 Mon Feb 27 19:26:28 +0000 2017 5
30 Mon Feb 27 19:26:29 +0000 2017 5
31 Mon Feb 27 19:26:30 +0000 2017 3
32 Mon Feb 27 19:26:31 +0000 2017 2
33 Mon Feb 27 19:26:32 +0000 2017 3
34 Mon Feb 27 19:26:33 +0000 2017 4
35 Mon Feb 27 19:26:34 +0000 2017 3
36 Mon Feb 27 19:26:35 +0000 2017 4
37 Mon Feb 27 19:26:36 +0000 2017 11
38 Mon Feb 27 19:26:37 +0000 2017 8
```

Best Time to post tweet:

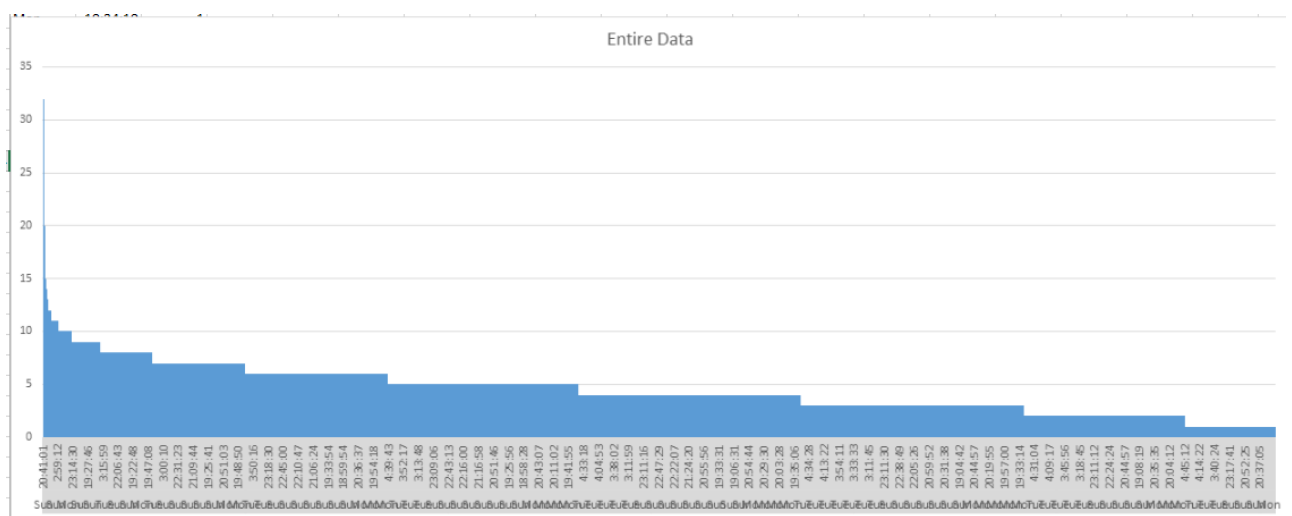
The below screenshot shows that “**Sun Feb 26 20:41:01 +0000 2017**” is the best time to post a tweet. It is the time that occurred 32 times more than any timestamp of tweet creation. This means most of the tweets are created at this time.

```
BestTime-1.txt [Read Only] (/tmp/mozilla_cloudera0) - gedit
Sun Feb 26 20:41:01 +0000 2017 32
```

Metric Criterion:

We have considered TimeStamp of the tweet to propose a metric criterion. TimeStamp refers to the time of the tweet creation. As we can see from the top 10 best times output most of the tweets are created during **Sunday** and the best time to post a tweet is also on Sunday.

Below is the graph for the entire data.



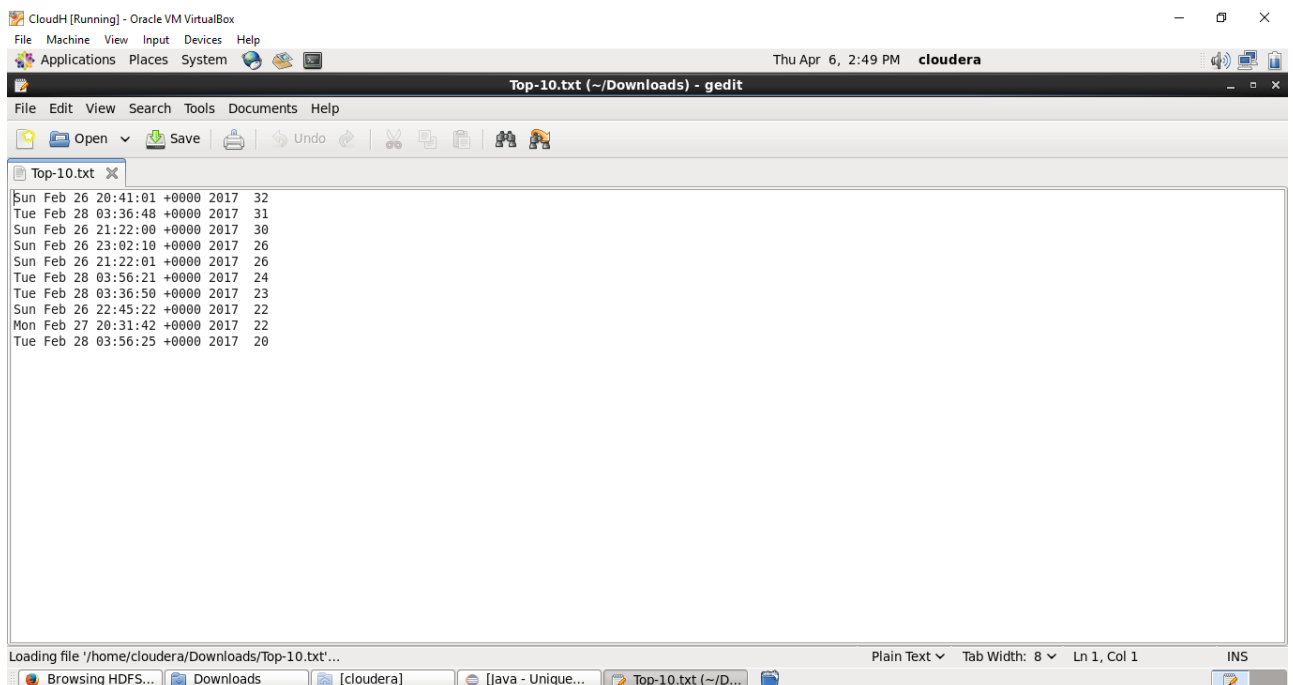
In the above graph X-axis represents the time and day the tweet was created and Y-axis shows the number of tweets created at that particular time. From the graph we can see that at timestamp “**Sun, 20:41:01**” most of the tweets are created and the lowest in this sample 20 graph is seen at “**Tue, 3:20:09**”. In this sample data out of 20 days; 11 days are Sundays i.e tweets are created or posted mostly on Sundays. From the data we can also assume that best day to post a tweet is on Sunday. And another 8 days the tweets are posted on Tuesdays and remaining one day the tweets are posted on Monday. Also, if we consider timings most of the tweets are collected during night time. From this we can conclude that best time or most tweets are posted during night time on Sundays.

Top Ten Times:

To get the **Top Ten** timestamps from the TimeStampOutput1 file from above by using the command:

```
hadoop fs -cat /user/cloudera/test/TimeStampOutput1.txt/part* | sort -n -k7 -r | head -n10 > Top-10.txt
```

The above hadoop command takes the timestamp output from the above and sorts the value in the key 7 that is the count of the timestamps when the tweets are collected and gives the top 10 timestamps and stores in the Top-10.txt file.



5.References

- <http://stackoverflow.com>
- <http://www.tweepy.org/>
- <https://www.cloudera.com/>
- <https://www.youtube.com/watch?v=TWYd0TD8Ops>
- <http://www.aegissofttech.com/Articles/how-to-get-top-n-words-count-using-big-data-hadoop-mapreduce-paradigm-with-developers-assistance.html>