

Sentimental Analysis of Twitter using Map Reduce

edureka!

edureka!

© 2014 Brain4ce Education Solutions Pvt. Ltd.

Table of Contents

Sentimental Analysis of Twitter using Map Reduce.....2

edureka!

Sentimental Analysis of Twitter using Map Reduce

Problem Statement:

Here, we have shown the sentimental analysis of twitter data which is done with the apache flume.

Important Links:

Fetching Data from Twitter using Apache Flume:

https://www.youtube.com/watch?v=723krKpwe_k

Edureka VM Installation:

https://edureka.wistia.com/medias/o1bd868187/download?media_file_id=64035017

Codes along with the Dataset:

https://edureka.wistia.com/medias/d7pjeszy5/download?media_file_id=57962991

Snippet of JSON Data:

https://edureka.wistia.com/medias/w1j2cqohlg/download?media_file_id=61275117

Dataset:

Let us consider a sample dataset as in the below screenshot.

```

/geosuperhdtv","utc_offset":null,"time_zone":null,"notifications":null,"profile_use_background_image":true,"friends_count":41,"p
bar_fill_color":"DDEEF6","screen_name":"geosuperhdtv","id_str":"1409155964","profile_image_url":"http://pbs.twimg.com/profile_in
/415372246706688000/aGx3Djx2_normal.jpeg","listed_count":1,"is_translator":false}}
{"filter_level":"medium","retweeted":false,"in_reply_to_screen_name":null,"possibly_sensitive":false,"truncated":false,"lang":"d
y_to_status_id_str":null,"id":"540586865423028225","in_reply_to_user_id_str":null,"timestamp_ms":"1417720926992","in_reply_to_stat
","created_at":"Thu Dec 04 19:22:06 +0000 2014","favorite_count":0,"place":null,"coordinates":null,"text":"RT @geosuperhdtv: wao
for #Pakistan in T20 Internationals. They have become first team to reach this milestone. #cwc15 #pakvz
#crv2026","contributors":null,"retweeted_status":{"filter_level":"low","contributors":null,"text":"wao 50th win for #Pakistan i
Internationals. They have become first team to reach this milestone. #cwc15 #pakvz
#cricket","geo":null,"retweeted":false,"in_reply_to_screen_name":null,"possibly_sensitive":false,"truncated":false,"lang":"en",
{"trends":[],"symbols":[],"urls":[],"hashtags":[{"text":"Pakistan","indices":[17,26]},{"text":"cwc15","indices":[103,109]},{
{"text":"pakvz","indices":[110,117]},{"text":"cricket","indices":[118,126]}],"user_mentions":
[{}],"in_reply_to_status_id_str":null,"id":"540586864789700608","source":"<a href='\"http://www.cloudhopper.com/\"' rel='\"nofollow
\">Cloudhopper</a>","in_reply_to_user_id_str":null,"favorited":false,"in_reply_to_status_id":null,"retweet_count":1,"created_at
04 19:22:06 +0000 2014","in_reply_to_user_id":null,"favorite_count":0,"id_str":"540586864789700608","place":null,"user":
{"location":"Dubai,
Pakistan","default_profile":true,"statuses_count":24530,"profile_background_tile":false,"lang":"en","profile_link_color":"0084
9155964","following":null,"favourites_count":153,"protected":false,"profile_text_color":"333333","verified":false,"description":
is Pakistan's 1st Sports channel.\r\nNow we launch Hd tv.\r\nFollow us for the live cricket updates and other sports
news."},"contributors_enabled":false,"profile_sidebar_border_color":"CODEED","name":"Geo Super
Hd","profile_background_color":"CODEED","created_at":"Tue May 07 01:36:41 +0000

```

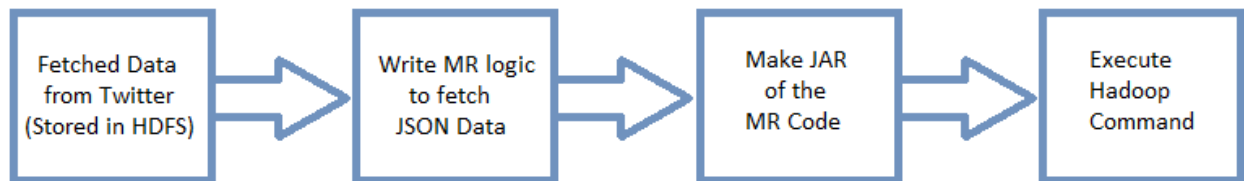
Dataset Description:

The above data represents the JSON format which is fetched from twitter. For better understanding, you can refer the snippet of JSON data which is given in previous section.

Tools and Technologies used:

- Apache Flume
- Eclipse

Dataflow Diagram:



Implementation:

To fetch the data from twitter, first by using apache flume we store the data in HDFS which is in JSON format by default. Further to read the all the fields, we used map-reduce, JSON SerDe (Serializer/Deserializer) to read JSON data.

In this, we are storing data in hdfs from twitter by using apache flume, reading the data and after map-reduce, storing the result again in hdfs.

Let us see how to do that:

First we have fetched the data from twitter using flume and data is stored in HDFS.

Wrote a map-reduce logic to fetch JSON data.

```
1 package twitter;
2
3
4 import java.io.IOException;
5
6
7
8
9
10
11
12
13 public class TwitterMapper extends Mapper<LongWritable, Text, Text, Text> {
14
15
16 @Override
17 protected void map(LongWritable key, Text value,
18 Mapper<LongWritable, Text, Text, LongWritable>.Context context,
19 Mapper<LongWritable, Text, Text, LongWritable>.Reporter reporter)
20 throws IOException, InterruptedException {
21
22 ObjectMapper mapper = new ObjectMapper();
23 mapper.configure(DeserializationConfig.Feature.FAIL_ON_UNKNOWN_PROPERTIES, false);
24
25 Tweet tweet = mapper.readValue(value.toString(), Tweet.class);
26
27 if(tweet!=null && tweet.getRetweeted_status()!=null && tweet.getRetweeted_status().getRetweeted_count()!=null)
28 {
29 String screen_name = tweet.getRetweeted_status().getUser().getScreen_name();
30 long retweeted_count = tweet.getRetweeted_status().getRetweeted_count();
31
32 if(screen_name!=null )
33 {
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
```

In Map-Reduce, firstly we read the complete data line by line. After that, we have created the objects according to our requirements. In this, we have created the objects of all fields which are mentioned in JSON file.

The format of first line of flume data (JSON) file is on Page 2.

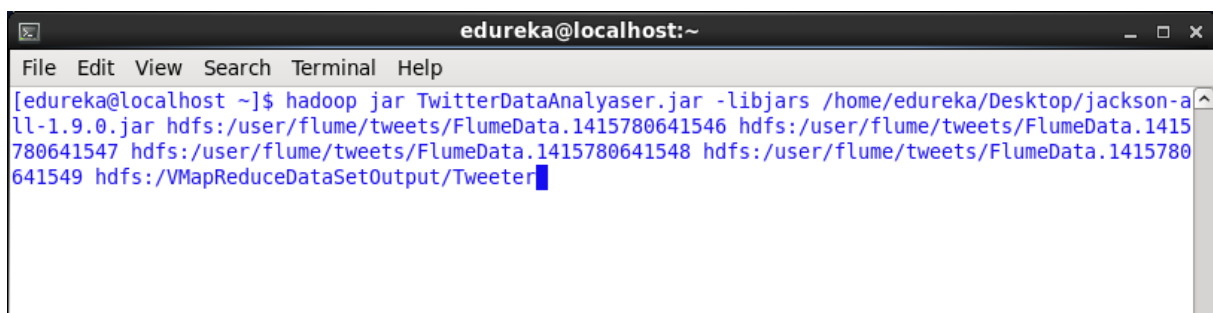
In this, we have focused in only two sections:

1. *Retweet_Count*
2. *Screen_Name, which is sub-section of User*

So, we have taken Screen_Name as Key and Retweet_Count as Value, both of these parameters is passed by mapper in the format of Key_Value pair.

Reducer receives both the parameters and put a counter for each Screen_Name and count all the Retweet_Count that is reducer output which is stored in HDFS.

Command: `hadoop jar TwitterDataAnalyaser.jar -libjars /home/cloudera/Desktop/jackson-all-1.9.0.jar hdfs:/user/flume/tweets/FlumeData.1415780641546 hdfs:/user/flume/tweets/FlumeData.1415780641547 hdfs:/user/flume/tweets/FlumeData.1415780641548 hdfs:/user/flume/tweets/FlumeData.1415780641549 hdfs:/VMapReduceDataSetOutput/Tweeter;`



```
edureka@localhost:~  
File Edit View Search Terminal Help  
[edureka@localhost ~]$ hadoop jar TwitterDataAnalyaser.jar -libjars /home/edureka/Desktop/jackson-all-1.9.0.jar hdfs:/user/flume/tweets/FlumeData.1415780641546 hdfs:/user/flume/tweets/FlumeData.1415780641547 hdfs:/user/flume/tweets/FlumeData.1415780641548 hdfs:/user/flume/tweets/FlumeData.1415780641549 hdfs:/VMapReduceDataSetOutput/Tweeter;
```

You can check the output in the HDFS, if you are using all the data files which is given with the codes. Then the output we will be as:

```
ImFaizhr      6
SRTrendulkar  1
bhaleraosarang 1
dna           3
rahulnanda86 15
tangerine_army 1
|
```

We have successfully fetched the data using MR code.