# Weblog and Social Networking Data Analysis Fostering To the Todays Business Using Hadoop Ecosystem

**Mohammed Faisal, Sanjay Gupta**

Department of Information Technology, Nizwa College of Technology

**Introduction:** *Big data has enhanced the business value, make a better advertisement, make better recommendation and get more user attraction, Big data has given a new definition to enhance business in the hi-tech era along with the lot of online data, we can have lot of offline data which can be used for analysis which could be predictive, prescription, diagnostic and descriptive. When the streaming data is becoming big, its sounds like more complicated, difficult to move, difficult to manage. But the technology simplified it that is nothing but Big Data which is really easy to handle, and from large volume of data we can do analysis by using the Hadoop ecosystem. In this paper will emphasize of predictive analysis using twitter data and weblog data using Hadoop, flume and hive. With the twitter we have discussed opinion mining, and with weblog data we have discussed how to collect the weblog data and how it is use full from business perspective. In this paper priority given to practically generated output along with justification and assumption, it's an effort to show the output practically in the relevant components of Hadoop Ecosystem.*

*Keywords: Weblog, social networking data, twitter, twitter application, big data, Hadoop, Flume, and Hive.*

## 1.1 Technological Background

### Hadoop Ecosystem



**Fig 1 – Hadoop Ecosystem**

Fig 1 is the illustration of Hadoop ecosystem, where HDFS is the bottom layer of ecosystem which managing the data in the Hadoop, just above of the HDFS we have YARN, which is providing the daemons, to do any kind of job we must start all daemons as a service, yarn has introduced in the

Hadoop 2, and all other Hadoop tools name as HBASE and Hive is running over the yarn where MapReduce is the inbuilt component of Hadoop and HIVE is the top most components of ecosystem which having all the capabilities of all components which is placed below of hive.

In ecosystem, flume is most important component which is used to get the input file from the web, such as weblog files and Twitter data. Oozie is a workflow scheduler system to manage Apache Hadoop jobs.

In this paper we will through two examples which is really the hottest buzz in the industries from business point of view [1].

## 2.0 Opinion mining using Twitter application:

Opinion mining is the analysis of people's opinions, sentiments, evaluations, appraisals, attitudes and emotions in relation to entities like individuals, products, events, services, organizations and topics by classifying the expressions as negative / positive opinions [2].

With this example of opinion mining from business perspective, we can do analysis on any product to know the opinion of the users/people of the product. If opinion if positive we can increase the production, if the opinion is negative we can do the changes into the product or we can reduce/stop the production. In this way it will help the merchandisers which will definitely solid and strait factors to predict the demand of the product into the current market.

## 2.1 Scope of work:

In this example you will be trying to do opinion mining of the opinions on twitter using

Hadoop Ecosystem. The reason we chose twitter is because the social media is gaining popularity for the customer reviews and it is also creating a good business, customer relationship in the market. The

reviews would certainly reflect the service of the company in the market.

**2.2 Objectives:**

The objective of the example is to find the sentiments of the users on Twitter. Companies can analyze the

Sentiments of the users regarding their products / services and use it for betterment of the same.

**2.3 Plan of work:**

i.   **Collect Data**: The initial step is to collect all twitter tweets.

ii.  **Preprocess Data:** We need to write MR Job to preprocess the data.

iii. **Classification of Data:** Write Hive UDF for the classification of the data into positive / negative Opinions.

iv.  **Print Data:** Final stage would be to print the desired results with the number of good and bad Tweets collected.

**2.3.1 Collect Data**

**Step 1:** To collect the data from the twitter first we need to create an application from https://apps.twitter.com/ and we must generate the 2 types of API keys and 2 kind's access token as the following:

**Consumer Key (API Key)**:
\*\*\*\*\*\*\*\*\*\*TD1sAyfQMMNIkh43Zrr

**Consumer Secret (API Secret):**
\*\*\*\*\*\*\*\*\*U5ycG5BAzlA1w4HjSZRGJaw7RuwMk5fcoW7AuEBLK

**Access Token**: \*\*\*\*\*\*\*\*\*\*\*4082176-VUAHhOS0dGVlkUljIartLy1qlqqAg2n

**Access Token Secret:**
\*\*\*\*\*\*\*8fGsJvvuJ0F47m4nmo2MlxCaPawg6dk3HjC

The above keys and access token should be updated into the flume configuration file.

**Step 2:** Create a Configuration file flume.conf file under **apache-flume-1.6.0-bin/conf** folder as following:

*TwitterAgent.sources= Twitter*

*TwitterAgent.channels= MemChannel*

*TwitterAgent.sinks=HDFS*

*TwitterAgent.sources.Twitter.type = com.cloudera.flume.source.TwitterSource*

*TwitterAgent.sources.Twitter.channels=MemChannel*

*TwitterAgent.sources.Twitter.consumerKey= \*\*\*\*\*\*\*\*\*\*\*TD1sAyfQMMNIkh43Zrr*

*TwitterAgent.sources.Twitter.consumerSecret=\*\*\*\*\*\*\*\*\*\*U5ycG5BAzlA1w4HjSZRGJaw7RuwMk5fcoW7AuEBLK*

*TwitterAgent.sources.Twitter.accessToken=\*\*\*\*\*\*\*\*\*\*\*4082176-VUAHhOS0dGVlkUljIartLy1qlqqAg2n*

*TwitterAgent.sources.Twitter.accessTokenSecret= \*\*\*\*\*\*\*\*\*8fGsJvvuJ0F47m4nmo2MlxCaPawg6dk3HjC*

*TwitterAgent.sources.Twitter.keywords= Toyota*

*TwitterAgent.sinks.HDFS.channel=MemChannel*

*TwitterAgent.sinks.HDFS.type=hdfs*

*TwitterAgent.sinks.HDFS.hdfs.path=hdfs://localhost:9000/Faisal/Toyota*

*TwitterAgent.sinks.HDFS.hdfs.fileType=DataStream*

*TwitterAgent.sinks.HDFS.hdfs.writeformat=Text*

*TwitterAgent.sinks.HDFS.hdfs.batchSize=1000*

*TwitterAgent.sinks.HDFS.hdfs.rollSize=0*

*TwitterAgent.sinks.HDFS.hdfs.rollCount=10000*

*TwitterAgent.sinks.HDFS.hdfs.rollInterval=600*

*TwitterAgent.channels.MemChannel.type=memory*

*TwitterAgent.channels.MemChannel.capacity=10000*

*TwitterAgent.channels.MemChannel.transactionCapacity=1000*

**Step 3**: Run below mentioned command to get the twitter data from local server as per given keywords, in this example we just given Toyota as keyword for the simplicity, keywords can be many as per the demand of analysis.



**Fig 2:** Snapshot to run the command to fetch the data from twitter

**flume-ng agent -n TwitterAgent –f/home/faisal/apache-flume-1.6.0-bin/conf/flume.con**

**Fig 3:** Snapshot to show the output of the snapshot Fig 2

We will get **FlumeData.1459583966469** data (Continuous similar files will get from nearest server) in the **Toyoto** directory in the **HDFS**. The contained of the files will be raw data which we need to parse by using Hive for further analysis purpose.



**Fig4:** Content of FlumeData.1459583966469

### 2.3.2 Preprocess Data

The above content is the raw data, it's in **JSON** format, and with this format we cannot analyze the data so we need to parse it with **HIVE.**

We need to create an external with the following command:

**Step 1:** CREATE EXTERNAL TABLE toyotadata (id BIGINT,entities STRUCT<hashtags:ARRAY<STRUCT<text:STRING>>>) ROW FORMAT SERDE 'com.cloudera.hive.serde.JSONSerDe' LOCATION '/Toyota';



**Fig5:** Snapshot to create external table

to view the created table we need to run the following command:

**Step 2:** Select * from hashtoyota;



**Fig5:** Snapshot to view hashtoyota table

### 2.3.3 Classification of Data

Still we need to work out to break the data in the form of key, value where will break each key and its all correspondent values individually.

By running the following steps we can simplified the data which will be easy to analyze:

**Step 1**: create table hashtoyota_word as select id as id,hashtag from hashtoyota LATERAL VIEW explode(words) w as ht;

**Step 2**: select hashtag, count(hashtag) from hashtag_toyota group by hashtag;

### 2.3.4 Print Data

Finally we will get the data as follows:



**Fig 6:** Snapshot to view hashtag_toyota counts of individual value of each key

With this format we can easily count the negative or positive verbs to get the conclusion, what is the opinion of people in that zone, because flume and twitter application will generate the data from the nearest server.

For the appropriate analysis as per the client requirements we can use various hive commands to customize our analysis.

### 3.0 Predictive analysis using weblog data of ecommerce sites:

After the emergence of ecommerce industry and rapid growth in this sector it is very important for the merchandiser to know every bit of information about the consumer like:

How many users hit every day to the websites?

From which regions more hits?

What type of product mostly they will look for?

User's feed back

And many more information.

Answers for the above questions can be answered if we able to analyses web log data efficiently.

To get complete information what is happening on the web site, the system should support gathering and processing information in **real-time**. Furthermore, as weblog data is collected and made available for analysis, it is also necessary to support **ad-hoc queries** via Apache Hive [3]

In Hadoop ecosystem there are many tools and technologies like Flume, Hive, Pig and Oozie which can be used for transporting and processing large amount of weblog data. In this section of the paper we will explore how we can transport weblogs data in HDFS using apache flume and then we will analyze it using apache Hive.
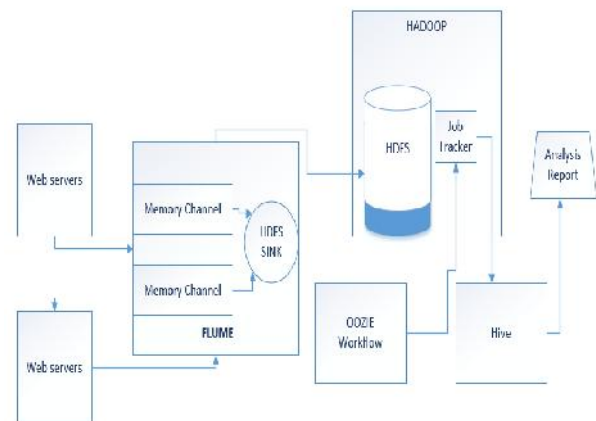


**Fig 7:** Illustration of conceptual view of weblog data analysis

**Steps:**

### 1. Transporting web logs data using Apache flume:

Apache flume allows users to stream data from multiple sources into Hadoop analysis, collect high volume weblogs in real time, guarantee data delivery and scale horizontally to handle additional data volume. Before configuring the flume we need to understand following components of flume:

```
flume-ng –n sandbox - -conf / apache-flume-1.6.0-bin /conf/flume.conf
```

**Source** – the object through which data enters into Flume. Sources either actively poll for data or passively wait for data to be delivered to them. A variety of sources allow data to be collected, such as log4j logs.

**Sink** – the object that delivers the data to the destination. There many ways to design sink. One example is the HDFS sink that writes events to HDFS.

**Channel** – This act between the Source and the Sink. Sources ingest events into the channel and the sinks drain the channel.

**Agent** – any physical Java virtual machine running Flume. It is a collection of sources, sinks and channels. [4]

Now configure and start apache flume using following commands:

```
# Flume agent config
sandbox.sources = eventlog
sandbox.channels = file_channel
sandbox.sinks = sink_to_hdfs

# Define / Configure source
sandbox.sources.eventlog.type = exec
sandbox.sources.eventlog.command = tail -F /var/log/eventlog-demo.log
sandbox.sources.eventlog.restart = true
sandbox.sources.eventlog.batchSize = 1000
#sandbox.sources.eventlog.type = seq

# HDFS sinks
sandbox.sinks.sink_to_hdfs.type = hdfs
sandbox.sinks.sink_to_hdfs.hdfs.fileType = DataStream
sandbox.sinks.sink_to_hdfs.hdfs.path = /flume/events
sandbox.sinks.sink_to_hdfs.hdfs.filePrefix = eventlog
sandbox.sinks.sink_to_hdfs.hdfs.fileSuffix = .log
sandbox.sinks.sink_to_hdfs.hdfs.batchSize = 1000

# Use a channel which buffers events in memory
sandbox.channels.file_channel.type = file
sandbox.channels.file_channel.checkpointDir = /var/flume/checkpoint
sandbox.channels.file_channel.dataDirs = /var/flume/data

# Bind the source and sink to the channel
sandbox.sources.eventlog.channels = file_channel
```

**Fig 8:** Snapshot configuration file flume.conf file under apache-flume-1.6.0-bin/conf

Next we will configure **log4j.properties** and append the following in it.

```
flume, root, logger=INFO, LOGFILE
flume.log.dir=/var/log/flume
flume.log.file=flume.log
```

Now start flume using following command:

Now generate the server log data

```
cd home/sanjay/apache-flume-1.6.0-bin
apache-flume-1.6.0-bin/conf/nano flume.conf
```

```
python generate_logs.py
```

When the log file has been generated, a timestamp will appear, and the command prompt will return to normal. It may take several seconds to generate the log file. Sample data will be as shown below:



**Fig 9:** Snapshot of sample weblog data

## 2. Creating tables from weblog data

```
CREATE EXTERNAL TABLE weblog_data
(
remoteIP STRING,
remoteLogname STRING,
user STRING,
time STRING,
request STRING,
statusCode STRING,
byteString STRING,
referral STRING,
browser STRING
)
ROW FORMAT SERDE
'org.apache.hadoop.hive.contrib.serde2.R
egexSerDe'
WITH SERDEPROPERTIES (
"input.regex" = '^(\\S+) (\\S+) (\\S+)
\\[(.+?)\\] \"([^\"]*)\" (\\S+) (\\S+)
\"([^\"]*)\" \"([^\"]*)\"',
"output.format.string" = "%1$s %2$s %3$s
%4$s %5$s %6$s %7$s %8$s %9$s")
LOCATION '/user/sanjay/weblogs';
```

Now we will start hive and use following command to create external table:

Now we can see our data from the table using following hive command

Now data will be in the desired format as you can see below:

```
hive>select * from weblog_data;
```

**Fig 9:** Snapshot of parsed output data

Similarly we can use different hive commands for analysis which could be predictive, prescriptive, diagnostic and descriptive.

## Conclusion:

It was an effort to emphasizing the use of new technology in today's business, in today's world social media is drastically changing the way of communication, and peoples are habituated to express his opinion, views on the social media.By using the Big Data technology we can utilize these social media data for the analysis purpose, which can be linked with the business to know the feedback of the users and people from that zone about a particular product,because big data will fetch the data from the server of that zone. Still variety of analysis can be done by using Hadoop ecosystem.

## Reference:

[1] Hadoop: The Definitive Guide by tom white, O'REILLY publication.

[2] Apache Cloudera-http://www.cloudera.com/documentation/other/tutorial/CDH5/Hadoop-Tutorial/ht_example_4_sentiment_analysis.html

[3] http://hortonworks.com/hadoop/data-science/

[4] http://hortonworks.com/hadoop-tutorial/how-to-refine-and-visualize-server-log-data/