

# **Project report for CSE 3025**

## **Stock Market Sentiment Analysis using various Big Data tools and their Comparative Analysis.**

<b>Vedang Sawarkar</b>	<b>19BCE1303</b>
<b>Divyanshi Thapa</b>	<b>19BCE1367</b>
<b>Prakriti Sharma</b>	<b>19BCE1655</b>

# **Abstract**

The main goal of this project is to analyze the stock market with the help of sentiment analysis done on stock market data collected from twitter using the Hadoop ecosystem. In a financially volatile market, as the stock market, it is important to have a very precise prediction of a future trend. Because of the financial crisis and scoring profits, it is mandatory to have a secure prediction of the values of the stocks. The successful prediction of a stock's future price could yield significant profit. Investors are constantly aware of the behaviour of stock markets. This affects their emotions and motivates them to buy or sell shares. Financial sentiment analysis allows us to understand the effect of social media reactions and emotions on the stock market and vice versa. Twitter, which has large audience potential, currently attracts an estimated average of 271 million users every month. So twitter-based models can then be built to aggregate the opinions of the collective population. They can be used to predict future trends while gaining useful insights into individual behavior. Keeping this in mind, we will be collecting the twitter data using NLP techniques and processing the data. Now there are a number of big data tools available so we will be using different tools in the Hadoop ecosystem to analyze the data and help in predicting the movement of the stock. Then the ability and efficiency will also be used to compare the tools.

# Motivation

Stock Market is the financial ground on which a high amount of data is released at every single point of time, which is very complex and nonlinear in nature. These data sets are used to predict the high profits and risks that a user can have in investing in a company. Manual predictive ways like studying the stock market with graphs, news, important dates, and many more practices failed many times which led to loss of profits because trading on the stock market required accurate and timely inputs. The magnitude of data that is generated within the stock market on a daily basis is impossible to be managed, analysed and made sense of by human beings due to the sheer volume of data generated and speed at which this financial data is being generated from various sources. Moreover, analyzing this much amount of data takes a lot of time. So, to overcome this challenge, we need some big data tools to analyze and give us filtered results. In the current day and age, financial analysis alone is no longer adequate for examining share prices and share price behaviour. These financial analyses need to be integrated with external factors such as social and economic trends within the economy, political environment, consumer behaviour and preferences etc. Social media platforms like Twitter are one of the biggest resources of this data. So we came up with an idea of generating common reviews of people regarding different stocks using sentiment analysis and then to predict stock from there. As the data is of a very huge amount, big data tools are required but there are so many tools available with their own speciality, it becomes difficult to choose. So we decided to use various tools and perform a comparative analysis on efficiency and accuracy of the tools.

# Introduction

Forecasting of stock market returns is a challenging research activity that is now expanding with the availability of new data sources, markets, financial instruments, and algorithms. At its core, the predictability of prices still raises important questions. Stock markets are driven by volatile factors such as microblogs and news that make it hard to predict stock market index based on merely the historical data. The enormous stock market volatility emphasizes the need to effectively assess the role of external factors in stock prediction.

Big data is pushing the financial industry, and also has an impact on investing. Huge amounts of data are generated each day since online trading has simplified the job and it's easier to view the market from your mobile by using an online trading platform or various stock trading applications. Social networks offering microblogging services enable the rapid spread of user generated content (UGC) from a handful of individuals to millions of people around the world in the form of short text, images or videos. Microblogging platforms have grown so exponentially, that they are now perceived as indispensable sources of information and are fast gaining popularity amongst users, organizations and researchers in various disciplines. The rich and continuous mass of data made available by these platforms is being harnessed with the purpose of studying individual and group behavior as well as global patterns especially in regards to sentiment towards the stock market trend. Twitter is currently the 10th most popular website globally with over 300 million active monthly users and much data related to stock market trends can be harnessed from this platform. There is a predictive relationship between twitter mood and stock movements but only up to a certain point. It was found in many research activities that positive correlation exists between volume of positive tweets and close prices, but in a time lag of 48 h, negative tweets are better indicators.

So by making the use of twitter data and performing sentiment analysis on it, we can find the trend of the stock market. But as twitter is a huge platform, the data generated will be of large amount and here comes the big data tools in picture. There are many research works going on for understanding the trends via social media platforms but the first and foremost question for these researches is which tool is to be used? There are many big data tools present right now and choosing the correct tool is necessary. So here in our project we apply sentiment analysis for predicting the stock market and also we will be comparing the efficiency of some famous big data tools.

# **Objectives**

- 1.** Extraction of tweets related to stock market from twitter.
- 2.** Processing and cleansing of data.
- 3.** Hashtags extraction for future works.
- 4.** Sentiment analysis on the clean dataset.
- 5.** Visualization of the analysis result and prediction of trend.
- 6.** Comparing the efficiency of various big data tools like Hive, Apache Spark and Cassandra by doing analysis by each one of them on the twitter data.

# Software and Hardware requirements

- The software and the hardware requirements of our project are listed below
  - 1) Twitter developer account.
    - For generating secret keys for using the API to extract tweets.
  - 2) System with linux environment.
  - 3) Hadoop version 3 and above.
    - For storing the big data and processing it.
  - 4) Hive version 2 and above
    - For sentiment analysis and querying the data.
  - 5) Apache Spark
    - For querying the data.
  - 6) Cassandra
    - For querying and comparison.
  - 7) Public account on Tableau
    - For data visualization.
  - 8) Jupyter Notebook.
    - For executing pySpark scripts.

# Project Explanation

## Tweets extraction

To extract the tweets, we used different API's or library's in python. So we used to retrieve the tweets using searchtweets library. This code is importing the credentials from the other file and sending a query(here we are querying AAPL) between the dates 1st May 2021 to 21st May 2021.

```
initials.py
1 import yaml
2 config = dict(
3     search_tweets_api = dict(
4         account_name = 'premium',
5         endpoint = 'https://api.twitter.com/1.1/tweets/search/30day/development.json',
6         consumer_key = 'C1W2afGekemUQYFQnbLXL0AC',
7         consumer_secret = 'z18ko1VLovS1boW0TDG2M0LbYVKjaFGSrphaZpvpXHWjkzUJw'
8     )
9 )
10 with open('twitter_keys_fullarchive.yaml', 'w') as config_file:
11     yaml.dump(config, config_file, default_flow_style=False)
```

```
ApiCredentials.py
1 From searchtweets import gen_rule_payload
2 From searchtweets import load_credentials
3
4 premium_search_args = load_credentials("twitter_keys_fullarchive.yaml",
5                                         yaml_key="search_tweets_api",
6                                         env_overwrite=False)
7 print(premium_search_args)
8 query = "AAPL"
9 rule = gen_rule_payload(query, results_per_call=100, from_date="2021-05-01", to_date="2021-05-20")
10 from searchtweets import ResultStream
11
12 rs = ResultStream(rule_payload=rule,
13                     max_results=25000,
14                     **premium_search_args)
15 print(rs)
16 import json
17 with open('tweets1.json', 'a', encoding='utf-8') as f:
18     for tweet in rs.stream():
19         json.dump(tweet, f)
20         f.write('\n')
21 print('done')
```

```
RulePayload.py
```

## Storing in HDFS

```
divyanshi@divyanshi:~/jcomp/pyScripts$ hdfs dfs -put '/home/divyanshi/jcomp/pyScripts/tweets1.json' /LSDP_Jcomp
divyanshi@divyanshi:~/jcomp/pyScripts$ hdfs dfs -mkdir /LSDP_Jcomp/twitter
divyanshi@divyanshi:~/jcomp/pyScripts$ hdfs dfs -cp /LSDP_Jcomp/twitter1.json /LSDP_Jcomp/twitter
cp: '/LSDP_Jcomp/twitter1.json': No such file or directory
divyanshi@divyanshi:~/jcomp/pyScripts$ hdfs dfs -cp /LSDP_Jcomp/tweets1.json /LSDP_Jcomp/twitter
divyanshi@divyanshi:~/jcomp/pyScripts$
```

We created the table in hive tweets\_raw with the few required attributes from the tweet object which we extracted from twitter api.

```
divyanshi@divyanshi:~/opt/hive/apache-hive-3.1.2-bin/bin$ ./hive
Hive Session ID = a1ee8113-c5ec-465e-a373-3a0bd3b9816d
Logging initialized using configuration in jar:file:/opt/hive/apache-hive-3.1.2-bin/lib/hive-common-3.1.2.jar!/hive-log4j2.properties Async: true
Hive Session ID = 4e0bf358-95fa-4614-809a-4ff464ad3c2
Hive MR is deprecated in Hive 2 and may not be available in the future versions. Consider using a different execution engine (i.e. spark, tez) or using Hive 1.x releases.
hive> CREATE EXTERNAL TABLE tweets_raw (
    > id BIGINT,
    > created_at STRING,
    > source STRING,
    > favorited BOOLEAN,
    > retweet_count INT,
    > retweeted_status STRUCT<
        > text:STRING,
        > user:STRUCT<screen_name:STRING,name:STRING>,
        > entities STRUCT<
            > urls:ARRAY<STRUCT<expanded_url:STRING>>,
            > user_mentions:ARRAY<STRUCT<screen_name:STRING,name:STRING>>,
            > hashtags:ARRAY<STRUCT<text:STRING>>,
            > text:STRING,
            > user:STRUCT<screen_name:STRING,
            > name:STRING,
            > time_zone:STRING,
            > in_reply_to_screen_name STRING,
            > year int,
            > month int,
            > day int,
            > hour int
        > ;
    > ROW FORMAT SERDE 'org.openx.data.jsonserde.JsonSerDe'
    > LOCATION '/LSDP_Jcomp/twitter'
    > ;
)
OK
Time taken: 1.588 seconds
hive>
```

## Sentiment Analysis:

We separated the word using split function on text data in tweets \_raw table and saved into split\_words table using below command. To perform sentiment analysis, we are using the AFINN dictionary which contains more than 2.5K words and values of polarity which range from -5.0 to +5.0 (Negative to Positive).

```
hive> create table split_words as select id as id,split(text, ' ') as words from tweets_raw;
Query ID = divyanshi_20210522125933_d34f84ff-3f6c-4bdf-8d65-45ca1b63dc42
Total jobs = 3
Launching Job 1 out of 3
Number of reduce tasks is set to 0 since there's no reduce operator
Starting Job = job_1621661702117_0002, Tracking URL = http://divyanshi:8088/proxy/application_1621661702117_0002/
Kill Command = /opt/hadoop/hadoop-3.3.0/bin/mapred job -kill job_1621661702117_0002
Hadoop job information for Stage-1: number of mappers: 1; number of reducers: 0
2021-05-22 12:59:45,621 Stage-1 map = 0%, reduce = 0%
2021-05-22 13:00:29,870 Stage-1 map = 100%, reduce = 0%, Cumulative CPU 52.51 sec
MapReduce Total cumulative CPU time: 52 seconds 510 msec
MapReduce Job = job_1621661702117_0002
Ended Job = job_1621661702117_0002
Stage-4 is selected by condition resolver.
Stage-3 is filtered out by condition resolver.
Stage-5 is filtered out by condition resolver.
Moving data to directory hdfs://localhost:9000/user/hive/warehouse/hive_staging_hive_2021-05-22_12-59-33_572_8016279155289195424-1/-ext-10002
Moving data to directory hdfs://localhost:9000/user/hive/warehouse/split_words
MapReduce Jobs Launched:
Stage-Stage-1: Map: 1  Cumulative CPU: 52.51 sec  HDFS Read: 135973399 HDFS Write: 3710285 SUCCESS
Total MapReduce CPU Time Spent: 52 seconds 510 msec
OK
Time taken: 58.525 seconds
hive> select * from split_words limit 5;
OK
1395166990227189761      ["RT","@puppy_trades:","Update","on","Tesla,","Apple,","Bitcoin,","Microsfot,","Disney,","Amazon,","and","Travel","Stocks"]
NULL      NULL
NULL      NULL
NULL      NULL
1395166956026802182      ["RT","@_r_o_b_b:","I've","come","to","the","realization","that","trading","options","is","way","safer","than","trading","#crypto","","",".","","Probably","mentally","healthier","too","D."]
Time taken: 0.244 seconds, Fetched: 5 row(s)
hive>
```

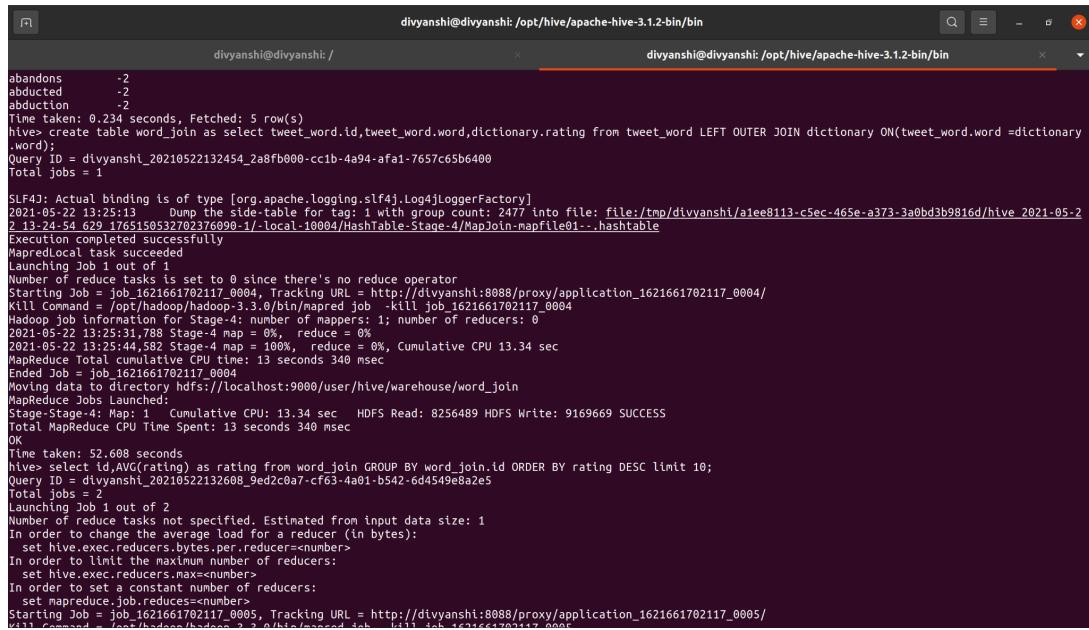
Now, we are splitting each word in the array as the new row. For this, we are using the explode function which extracts the element from array into new row using LATERAL VIEW. A table **tweet\_word** is created.

```
divyanshi@divyanshi:/opt/hive/apache-hive-3.1.2-bin/bin
divyanshi@divyanshi:/opt/hive/apache-hive-3.1.2-bin/bin
NULL      NULL
NULL      NULL
1395166956026802182      ["RT","@_r_o_b_b:","I've","come","to","the","realization","that","trading","options","is","way","safer","than","trading","#crypto","","",".","","Probably","mentally","healthier","too","D."]
Time taken: 0.244 seconds, Fetched: 5 row(s)
hive> create table tweet_word as select id,word from split_words LATERAL VIEW explode(words) w as word;
Query ID = divyanshi_20210522130221_707ae4cc-3683-437e-adf3-1ce12328d302
Total jobs = 3
Launching Job 1 out of 3
Number of reduce tasks is set to 0 since there's no reduce operator
Starting Job = job_1621661702117_0003, Tracking URL = http://divyanshi:8088/proxy/application_1621661702117_0003/
Kill Command = /opt/hadoop/hadoop-3.3.0/bin/mapred job -kill job_1621661702117_0003
Hadoop job information for Stage-1: number of mappers: 1; number of reducers: 0
2021-05-22 13:02:33,489 Stage-1 map = 0%, reduce = 0%
2021-05-22 13:02:41,803 Stage-1 map = 100%, reduce = 0%, Cumulative CPU 6.27 sec
MapReduce Total cumulative CPU time: 6 seconds 270 msec
Ended Job = job_1621661702117_0003
Stage-4 is selected by condition resolver.
Stage-5 is filtered out by condition resolver.
Moving data to directory hdfs://localhost:9000/user/hive/warehouse/hive_staging_hive_2021-05-22_13-02-21_979_4233365098245905421-1/-ext-10001
Moving data to directory hdfs://localhost:9000/user/hive/warehouse/tweet_word
MapReduce Jobs Launched:
Stage-Stage-1: Map: 1  Cumulative CPU: 6.27 sec  HDFS Read: 3715377 HDFS Write: 8248529 SUCCESS
Total MapReduce CPU Time Spent: 6 seconds 270 msec
OK
Time taken: 21.32 seconds
hive> describe tweet_word;
OK
id          bigint
word        string
Time taken: 0.115 seconds, Fetched: 2 row(s)
hive> select * from tweet_word limit 5;
OK
1395166990227189761      RT
1395166990227189761      @puppy_trades:
1395166990227189761      Update
1395166990227189761      on
1395166990227189761      Tesla,
Time taken: 0.179 seconds, Fetched: 5 row(s)
hive>
```

Now, we will create a new table named **dictionary** and will load the AFINN dictionary data into that table **dictionary**.

```
hive> LOAD DATA LOCAL INPATH '/home/divyanshi/jcomp/AFINN' into TABLE dictionary;
Loading data to table default.dictionary
OK
Time taken: 1.242 seconds
hive> select * from dictionary limit 5;
OK
abandon     -2
abandoned   -2
abandons    -2
abducted    -2
abduction   -2
Time taken: 0.234 seconds, Fetched: 5 row(s)
hive>
```

Create a new table word\_join by joining the dictionary table with tweet\_word table.

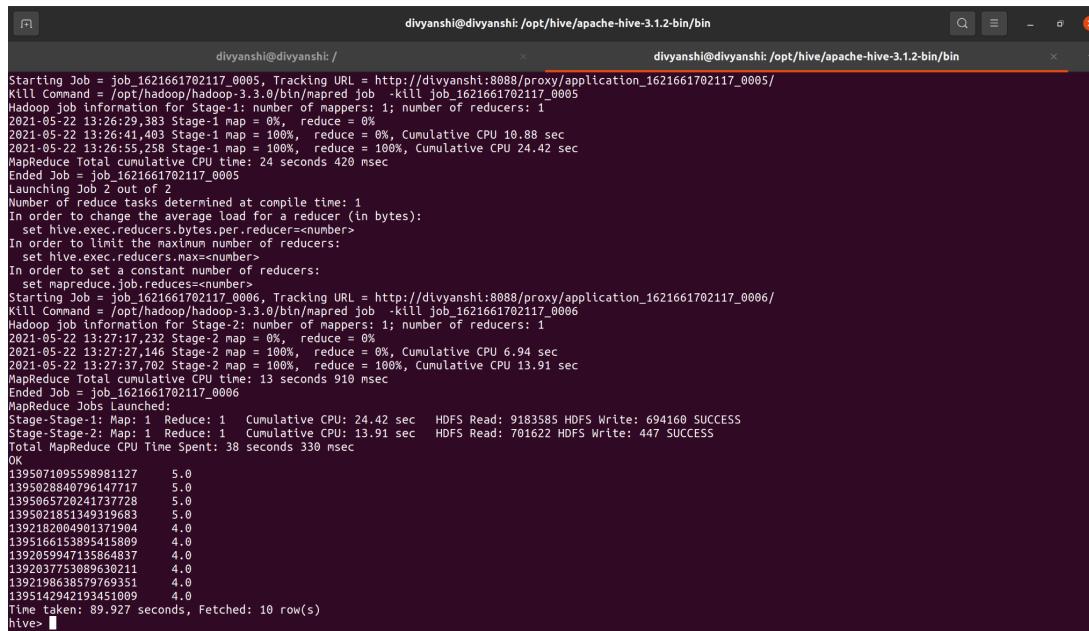


The screenshot shows two terminal windows side-by-side. Both windows have the title bar 'divyanshi@divyanshi: /opt/hive/apache-hive-3.1.2-bin/bin'. The left window shows the command 'hive> create table word\_join as select tweet\_word.id,tweet\_word.word,dictionary.rating from tweet\_word LEFT OUTER JOIN dictionary ON(tweet\_word.word =dictionary .word);' and its output, which includes a query ID and a note about actual binding. The right window shows the command 'hive> select id,AVG(rating) as rating from word join GROUP BY word join.id ORDER BY rating DESC limit 10;' and its output, which includes a query ID and the resulting data.

```
abandons      -2
abducted      -2
abduction     -2
Time taken: 0.234 seconds, Fetched: 5 row(s)
hive> create table word_join as select tweet_word.id,tweet_word.word,dictionary.rating from tweet_word LEFT OUTER JOIN dictionary ON(tweet_word.word =dictionary .word);
Query ID = divyanshi_20210522132454_2a8fb000-cc1b-4a94-afab-7657c65b6400
Total jobs = 1

SLF4J: Actual binding is of type [org.apache.logging.slf4j.Log4jLoggerFactory]
2021-05-22 13:25:13   Dump the side-table for tag: 1 with group count: 2477 into file: file:/tmp/divyanshi/a1ee8113-c5ec-465e-a373-3a0bd3b9816d/hive_2021-05-22_13-25-13_54_629_1765150532702376990-1-local-10004/HashTable-Stage-4/MapJoin-mapfile01...hashtable
Execution completed successfully
MapredLocal tasks succeeded
Launching Job 1 out of 1
Number of reduce tasks is set to 0 since there's no reduce operator
Starting Job = job_1621661702117_0004, Tracking URL = http://divyanshi:8088/proxy/application_1621661702117_0004/
Kill Command = /opt/hadoop/hadoop-3.3.0/bin/mapred job -kill job_1621661702117_0004
Hadoop job information for Stage-4: number of mappers: 1; number of reducers: 0
2021-05-22 13:25:31,785 Stage-4 map = 0%, reduce = 0%
2021-05-22 13:25:44,582 Stage-4 map = 100%, reduce = 0%, Cumulative CPU 13.34 sec
MapReduce Total cumulative CPU time: 13 seconds 340 msec
Ended Job = job_1621661702117_0004
Moving data to directory hdfs://localhost:9000/user/hive/warehouse/word_join
MapReduce Jobs Launched:
Stage-Stage-4: Map: 1  Cumulative CPU: 13.34 sec  HDFS Read: 8256489 HDFS Write: 9169669 SUCCESS
Total MapReduce CPU Time Spent: 13 seconds 340 msec
OK
Time taken: 52.608 seconds
hive> select id,AVG(rating) as rating from word join GROUP BY word join.id ORDER BY rating DESC limit 10;
Query ID = divyanshi_20210522132608_9ed2c0a7-cf63-4a01-b542-6d4549e8a2e5
Total jobs = 2
Launching Job 1 out of 2
Number of reduce tasks not specified. Estimated from input data size: 1
In order to change the average load for a reducer (in bytes):
  set hive.exec.reducers.bytes.per.reducer=<number>
In order to limit the maximum number of reducers:
  set hive.exec.reducers.max=<number>
In order to set a constant number of reducers:
  set mapreduce.job.reduces=<number>
Starting Job = job_1621661702117_0005, Tracking URL = http://divyanshi:8088/proxy/application_1621661702117_0005/
Kill Command = /opt/hadoop/hadoop-3.3.0/bin/mapred job -kill job_1621661702117_0005
```

After joining now, we got the polarity value for each tweet in the word\_join table.



The screenshot shows two terminal windows side-by-side. Both windows have the title bar 'divyanshi@divyanshi: /opt/hive/apache-hive-3.1.2-bin/bin'. The left window shows the command 'hive> select id,AVG(rating) as rating from word join GROUP BY word join.id ORDER BY rating DESC limit 10;' and its output, which includes a query ID and the resulting data. The right window shows the command 'hive> select id,AVG(rating) as rating from word join GROUP BY word join.id ORDER BY rating DESC limit 10;' and its output, which includes a query ID and the resulting data.

```
Starting Job = job_1621661702117_0005, Tracking URL = http://divyanshi:8088/proxy/application_1621661702117_0005/
Kill Command = /opt/hadoop/hadoop-3.3.0/bin/mapred job -kill job_1621661702117_0005
Hadoop job information for Stage-1: number of mappers: 1; number of reducers: 1
2021-05-22 13:26:29,383 Stage-1 map = 0%, reduce = 0%
2021-05-22 13:26:41,403 Stage-1 map = 100%, reduce = 0%, Cumulative CPU 10.88 sec
2021-05-22 13:26:55,250 Stage-1 map = 100%, reduce = 100%, Cumulative CPU 24.42 sec
MapReduce Total cumulative CPU time: 24 seconds 420 msec
Ended Job = job_1621661702117_0005
Launching Job 2 out of 2
Number of reduce tasks determined at compile time: 1
In order to change the average load for a reducer (in bytes):
  set hive.exec.reducers.bytes.per.reducer=<number>
In order to limit the maximum number of reducers:
  set hive.exec.reducers.max=<number>
In order to set a constant number of reducers:
  set mapreduce.job.reduces=<number>
Starting Job = job_1621661702117_0006, Tracking URL = http://divyanshi:8088/proxy/application_1621661702117_0006/
Kill Command = /opt/hadoop/hadoop-3.3.0/bin/mapred job -kill job_1621661702117_0006
Hadoop job information for Stage-2: number of mappers: 1; number of reducers: 1
2021-05-22 13:27:17,232 Stage-2 map = 0%, reduce = 0%
2021-05-22 13:27:18,146 Stage-2 map = 100%, reduce = 0%, Cumulative CPU 6.94 sec
2021-05-22 13:27:37,702 Stage-2 map = 100%, reduce = 100%, Cumulative CPU 13.91 sec
MapReduce Total cumulative CPU time: 13 seconds 910 msec
Ended Job = job_1621661702117_0006
MapReduce Jobs Launched:
Stage-Stage-1: Map: 1 Reduce: 1  Cumulative CPU: 24.42 sec  HDFS Read: 9183585 HDFS Write: 694160 SUCCESS
Stage-Stage-2: Map: 1 Reduce: 1  Cumulative CPU: 13.91 sec  HDFS Read: 701622 HDFS Write: 447 SUCCESS
Total MapReduce CPU Time Spent: 38 seconds 330 msec
OK
1395071095598981127  5.0
1395028840796147717  5.0
1395065720241737728  5.0
1395021851349319683  5.0
1392182004901371964  4.0
1395166153895415809  4.0
1392859947155864837  4.0
1392037753089636211  4.0
1392198638579769351  4.0
1395142942193451069  4.0
Time taken: 89.927 seconds, Fetched: 10 row(s)
hive>
```

Now, we are saving the same id based grouped data into the **testjoin** table and we are joining the tweets\_raw table and test join table based on id and creating new table **tweets**.

# How to use

- 1) Run the RulePayload.py file to extract all the tweets from a particular start date to end date and store them in tweets1.json file
- 2) Store this tweets1.json file in hdfs
- 3) Start hive
- 4) Add json serde dependency jar to hive. JSON SerDe is used to map the json object to the table schema.
- 5) Create tweets\_raw table in hive.
- 6) Now we perform sentiment analysis using the AFINN dictionary.
- 7) We created a split\_words table. We separated the words using split function on text data in tweets\_raw table and saved into split\_words table.
- 8) Now, we are splitting each word in the array as the new row. For this, we are using the explode function which extracts the element from the array into a new row. Since it has some limitations we are using a LATERAL VIEW. We are storing these words in a new table called tweet\_words
- 9) Now, we will create a new table named dictionary and will load the AFINN dictionary data into that table.
- 10) Create a new table word\_join by joining the ‘dictionary’ table and ‘tweet\_word’ table.
- 11) Find average rating or polarity value of each tweet using ID attribute.
- 12) Save the ID based grouped data in a new table called ‘testjoin’.
- 13) Now join the tweets\_raw table and testjoin table based on id and create a new table called ‘tweets’.
- 14) Now we can run sql queries in hive to analyse different sections of data.
- 15) Create another table know as ‘outtable’ by joining ‘tweets\_raw’ table and ‘testjoin’ table
- 16) Export the data stored in ‘outdata’ table in a ‘out\_data.csv’ file
- 17) Now we extract all the hashtags from the tweets by running th parse.py python script.
- 18) All hashtags will be stored in the hashtags.txt file.
- 19) Store the hashtags.txt file in hdfs.
- 20) Run the word count program to count each hashtag.
- 21) Store the result in a local json file called as ‘tweets1.json’.
- 22) Now we create data frames of out\_data.csv and tweets1.json in spark.
- 23) Then we run the pyspark script to perform data analysis on the two data frames in spark and compare the execution time of same queries running in hive.
- 24) Visualize the polarities in the ‘out\_data.csv’ file which tells us the percentage of positive, negative and neutral tweets.
- 25) Plot the graph of time vs avg. polarity.

# Dataset Description

## Dataset:

**tweets1.json:** it contains the data related to stock market extracted from the Twitter API

## Tables and HDFS files:

- tweets\_raw : with the few required attributes from the tweet object.
  - split\_word : separated the word using split function
  - tweet\_word : splitting each word in the array as the new row.
  - dictionary : load the AFINN dictionary data.
  - word\_join : by joining the dictionary table with tweet\_word table.
  - testjoin : the same id based grouped data.
  - tweets : joining the tweets\_raw table and test join table based on id.
- hashtags.txt : file containing the extracted hashtags.

Link to dataset and hdfs files:

[https://drive.google.com/drive/folders/1bGtk24KIJQzVMigj7qEFAj4mGK0\\_rIom?usp=sharing](https://drive.google.com/drive/folders/1bGtk24KIJQzVMigj7qEFAj4mGK0_rIom?usp=sharing)

# Results

## Tweets extraction

```
divyanshi@divyanshi:~/jcomp/py scripts$ python3 ApiCredentials.py
divyanshi@divyanshi:~/jcomp/py scripts$ python3 RulePayload.py
Grabbing bearer token from OAUTH
('bearer_token': 'AAAAAAAAAAAAAAJoxPgEAAAAAmVJ%2FoCGlh%2FfdfaxMFNGD05mj0kX3DiJHiTi000cKpCmRoael0VChjQchByMfpITDf45lvqCoL3WNGF', 'endpoint': 'https://api.twitter.com/1.1/tweets/search/30day/development.json', 'extra_headers_dict': None}
ResultStream:
{
  "username": null,
  "endpoint": "https://api.twitter.com/1.1/tweets/search/30day/development.json",
  "rule_payload": {
    "query": "AAPL",
    "maxResults": 100,
    "toDate": "202105200000",
    "fromDate": "202105010000"
  },
  "tweetify": true,
  "max_results": 25000
}
HTTP Error code: 429: {"error":{"message":"Exceeded rate limit","sent":"2021-05-22T05:15:44+00:00","transactionId":"f7c5bbf3d1403ddb"}}
Request payload: {'query': 'AAPL', 'maxResults': 100, 'toDate': '202105200000', 'fromDate': '202105010000', 'next': 'eyJtYXhJZCI6MTMSNDgzNTkwOTUyNzM2MzU4Nn0='}
Rate limit hit... Will retry...
Will retry in 4 seconds...
HTTP Error code: 429: {"error":{"message":"Exceeded rate limit","sent":"2021-05-22T05:15:48+00:00","transactionId":"27f20b84cd43380e"}}
Request payload: {'query': 'AAPL', 'maxResults': 100, 'toDate': '202105200000', 'fromDate': '202105010000', 'next': 'eyJtYXhJZCI6MTMSNDgzNTkwOTUyNzM2MzU4Nn0='}
Rate limit hit... Will retry...
Will retry in 10 seconds...
HTTP Error code: 429: {"error":{"message":"Exceeded rate limit","sent":"2021-05-22T05:16:05+00:00","transactionId":"152765aeba84c45"}}
Request payload: {'query': 'AAPL', 'maxResults': 100, 'toDate': '202105200000', 'fromDate': '202105010000', 'next': 'eyJtYXhJZCI6MTMSNDgzNTkwOTUyNzM2MzU4Nn0='}
Rate limit hit... Will retry...
Will retry in 36 seconds...
HTTP Error code: 429: {"error":{"message":"Exceeded rate limit","sent":"2021-05-22T05:17:10+00:00","transactionId":"97e722ad91863e2c"}}
Request payload: {'query': 'AAPL', 'maxResults': 100, 'toDate': '202105200000', 'fromDate': '202105010000', 'next': 'eyJtYXhJZCI6MTMSNDU0NTc40TI0NjkxODY1N30='}
Rate limit hit... Will retry...
Will retry in 4 seconds...
HTTP Error code: 429: {"error":{"message":"Exceeded rate limit","sent":"2021-05-22T05:17:15+00:00","transactionId":"54e94adfa5449613"}}
Request payload: {'query': 'AAPL', 'maxResults': 100, 'toDate': '202105200000', 'fromDate': '202105010000', 'next': 'eyJtYXhJZCI6MTMSNDU0NTc40TI0NjkxODY1N30='}
Rate limit hit... Will retry...
Will retry in 16 seconds...
HTTP Error code: 429: {"error":{"message":"Exceeded rate limit","sent":"2021-05-22T05:17:31+00:00","transactionId":"ff69f2542717b341"}}
Request payload: {'query': 'AAPL', 'maxResults': 100, 'toDate': '202105200000', 'fromDate': '202105010000', 'next': 'eyJtYXhJZCI6MTMSNDU0NTc40TI0NjkxODY1N30='}
Rate limit hit... Will retry...
Will retry in 36 seconds...
HTTP Error code: 429: {"error":{"message":"Exceeded rate limit","sent":"2021-05-22T05:18:37+00:00","transactionId":"d56c5850a22afa6e"}}
Request payload: {'query': 'AAPL', 'maxResults': 100, 'toDate': '202105200000', 'fromDate': '202105010000', 'next': 'eyJtYXhJZCI6MTMSNDU0NTc40TI0NjkxODY1N30='}
```

```
divyanshi@divyanshi:~/jcomp/py scripts$ wc -l tweets1.json
25000 tweets1.json
divyanshi@divyanshi:~/jcomp/py scripts$
```

### **Final table (tweets) after sentiment analysis (with polarity of every tweet)**

## Hive tables

Hadoop Overview Datanodes Datanode Volume Failures Snapshot Startup Progress Utilities ▾

### Browse Directory

/user/hive/warehouse

Show 25 entries Search:

	Permission	Owner	Group	Size	Last Modified	Replication	Block Size	Name	
	drwxr-xr-x	divyanshi	supergroup	0 B	May 22 13:23	0	0 B	dictionary	trash
	drwxr-xr-x	divyanshi	supergroup	0 B	May 22 18:29	0	0 B	outtable	trash
	drwxr-xr-x	divyanshi	supergroup	0 B	May 22 13:00	0	0 B	split_words	trash
	drwxr-xr-x	divyanshi	supergroup	0 B	May 22 13:29	0	0 B	testjoin	trash
	drwxr-xr-x	divyanshi	supergroup	0 B	May 22 13:02	0	0 B	tweet_word	trash
	drwxr-xr-x	divyanshi	supergroup	0 B	May 22 13:31	0	0 B	tweets	trash
	drwxr-xr-x	divyanshi	supergroup	0 B	May 22 13:25	0	0 B	word_join	trash

Showing 1 to 7 of 7 entries Previous 1 Next

Hadoop, 2020.

## Cassandra Analysis

```
divyanshi@divyanshi: ~
cqlsh:divyanshi> create table ctweets (created_at text,source text,retweet_count int,screenname text,username text,polarity double,review text,primary key(created_at));
cqlsh:divyanshi> copy ctweets (created_at,source,retweet_count,screenname,username,polarity,review) from '/home/divyanshi/jcomp/tweets_dataset.csv' with header=true and delimiter=';';
Using 7 child processes
Starting copy of divyanshi.ctweets with columns [created_at, source, retweet_count, screenname, username, polarity, review].
Processed: 25460 rows; Rate: 19482 rows/s; Avg. rate: 20355 rows/s
24998 rows imported from 1 files in 1.251 seconds (0 skipped).
```

```
divyanshi@divyanshi: ~
cqlsh:divyanshi> select count(*) from ctweets;
count
-----
17916
(1 rows)
Warnings :
Aggregation query used without partition key
cqlsh:divyanshi> select count(*) as positive from ctweets where polarity > 3 allow filtering;
positive
-----
77
(1 rows)
Warnings :
Aggregation query used without partition key
cqlsh:divyanshi> select count(*) as positive from ctweets where polarity > 0 allow filtering;
positive
-----
2913
(1 rows)
Warnings :
Aggregation query used without partition key
cqlsh:divyanshi> select count(*) as negative from ctweets where polarity < 0 allow filtering;
negative
-----
982
(1 rows)
Warnings :
Aggregation query used without partition key
```

## Tools Comparison (Hive vs SparkSQL)

- **Query: select count(\*) from tweets;**

### Hive



```
hive> select count (*) from tweets;
OK
25000
Time taken: 0.108 seconds, Fetched: 1 row(s)
```

**Time Taken: 0.108 seconds**

### SparkSQL

```
In [104]: start = time.time()
q9 = spark.sql("SELECT count(*) FROM out_data")
q9.show()
end = time.time()
print("Time taken = " + str(end - start))

+-----+
|count(1)|
+-----+
| 25461|
+-----+

Time taken = 0.08100366592407227
```

**Time Taken: 0.08100366592407227 seconds**

- **Query :select distinct `user`.screen\_name as name, `user`.followers\_count as count from tweets where size(entities.hashtags) > 0 order by count desc limit 5;**

### Hive

```
hive> select distinct `user`.screen_name as name, `user`.followers_count as count
   > from tweets
   > where size(entities.hashtags) > 0
   > order by count desc
   > limit 5;
Query ID = divyanshi_20210522181440_8f5b0a8f-c1f3-48c7-a075-e06928aa5dc4
Total jobs = 2
Launching Job 1 out of 2
Number of reduce tasks not specified. Estimated from input data size: 1
In order to change the average load for a reducer (in bytes):
  set hive.exec.reducers.bytes.per.reducer=<number>
In order to limit the maximum number of reducers:
  set hive.exec.reducers.max=<number>
In order to set a constant number of reducers:
  set mapreduce.job.reduces=<number>
Starting Job = job_1621661702117_0012, Tracking URL = http://divyanshi:8088/proxy/application_1621661702117_0012/
Kill Command = /opt/hadoop/hadoop-3.3.0/bin/mapred job -kill job_1621661702117_0012
Hadoop job information for Stage-1: number of mappers: 1; number of reducers: 1
2021-05-22 18:14:46,800 Stage-1 map = 0%, reduce = 0%
2021-05-22 18:14:51,940 Stage-1 map = 100%, reduce = 0%, Cumulative CPU 5.15 sec
2021-05-22 18:14:57,060 Stage-1 map = 100%, reduce = 100%, Cumulative CPU 7.58 sec
MapReduce Total cumulative CPU time: 7 seconds 580 msec
Ended Job = job_1621661702117_0012
Launching Job 2 out of 2
Number of reduce tasks determined at compile time: 1
In order to change the average load for a reducer (in bytes):
  set hive.exec.reducers.bytes.per.reducer=<number>
In order to limit the maximum number of reducers:
  set hive.exec.reducers.max=<number>
In order to set a constant number of reducers:
  set mapreduce.job.reduces=<number>
Starting Job = job_1621661702117_0013, Tracking URL = http://divyanshi:8088/proxy/application_1621661702117_0013/
Kill Command = /opt/hadoop/hadoop-3.3.0/bin/mapred job -kill job_1621661702117_0013
Hadoop job information for Stage-2: number of mappers: 1; number of reducers: 1
2021-05-22 18:15:09,282 Stage-2 map = 0%, reduce = 0%
2021-05-22 18:15:13,383 Stage-2 map = 100%, reduce = 0%, Cumulative CPU 1.33 sec
2021-05-22 18:15:18,504 Stage-2 map = 100%, reduce = 100%, Cumulative CPU 3.65 sec
```

```

divyanshi@divyanshi: /opt/hive/apache-hive-3.1.2-bin/bin
divyanshi@divyanshi: /opt/hive/apache-hive-3.1.2-bin/bin

set hive.exec.reducer.bytes_per.reducer=<number>
In order to limit the maximum number of reducers:
set hive.exec.reducer.max=<number>
In order to set a constant number of reducers:
set mapreduce.job.reduces=<number>
Starting Job : job_1621661702117_0012, Tracking URL = http://divyanshi:8088/proxy/application_1621661702117_0012/
Kill Command = /opt/hadoop/hadoop-3.3.0/bin/mapred job -kill job_1621661702117_0012
Hadoop job information for Stage-1: number of mappers: 1; number of reducers: 1
2021-05-22 18:14:46,880 Stage-1 map = 0%, reduce = 0%
2021-05-22 18:14:51,940 Stage-1 map = 100%, reduce = 0%, Cumulative CPU 5.15 sec
2021-05-22 18:14:57,060 Stage-1 map = 100%, reduce = 100%, Cumulative CPU 7.58 sec
MapReduce Total cumulative CPU time: 7 seconds 580 msec
Ended Job : job_1621661702117_0012
Launching Job 2 out of 2
Number of reduce tasks determined at compile time: 1
In order to change the average load for a reducer (in bytes):
set hive.exec.reducer.bytes.per.reducer=<number>
In order to limit the maximum number of reducers:
set hive.exec.reducer.max=<number>
In order to set a constant number of reducers:
set mapreduce.job.reduces=<number>
Starting Job : job_1621661702117_0013, Tracking URL = http://divyanshi:8088/proxy/application_1621661702117_0013/
Kill Command = /opt/hadoop/hadoop-3.3.0/bin/mapred job -kill job_1621661702117_0013
Hadoop job information for Stage-2: number of mappers: 1; number of reducers: 1
2021-05-22 18:15:09,282 Stage-2 map = 0%, reduce = 0%
2021-05-22 18:15:13,383 Stage-2 map = 100%, reduce = 0%, Cumulative CPU 1.33 sec
2021-05-22 18:15:18,504 Stage-2 map = 100%, reduce = 100%, Cumulative CPU 3.65 sec
MapReduce Total cumulative CPU time: 3 seconds 650 msec
Ended Job : job_1621661702117_0013
MapReduce Jobs Launched:
Stage-Stage-1: Map: 1 Reduce: 1 Cumulative CPU: 7.58 sec HDFS Read: 11100165 HDFS Write: 231 SUCCESS
Stage-Stage-2: Map: 1 Reduce: 1 Cumulative CPU: 3.65 sec HDFS Read: 7677 HDFS Write: 217 SUCCESS
Total MapReduce CPU Time Spent: 11 seconds 230 msec
OK
CNBC    4150506
CNBC    4150493
CNBC    4150469
MarketWatch  4021394
BW     1848295
Time taken: 40.167 seconds, Fetched: 5 row(s)
hive> 

```

**Time taken: 40.167 seconds**

## SparkSQL

```

In [21]: start = time.time()
q3 = spark.sql("SELECT user.name, max(user.followers_count) as followers_count FROM tweets1 WHERE text like '%AAPL%'")
q3.show()
end = time.time()
print("Time taken = " + str(end - start))

+-----+-----+
|      name|followers_count|
+-----+-----+
| CNBC|        4150506|
| MarketWatch|        4021394|
| TheStreet|        782944|
| Charles V Payne|        556900|
| Steve Burns|        363566|
| Carl Quintanilla|        300922|
| Hank Lockwood|        274053|
| Jon Najarian|        262431|
| Investors.com|        252172|
| Amit Bhawani|        249170|
| Cheddar News|        249027|
| Ramp Capital|        243000|
| FXStreet News|        204190|
| LiveSquawk|        194892|
| Morning Brew =|        187974|
+-----+-----+
Time taken = 1.4111835956573486

```

**Time Taken: 1.4111835956573486 seconds**

- **Query: select `user`.name,count(\*) as count from tweets\_raw group by `user`.name order by count desc limit 10;**

## Hive

```
divyanshi@divyanshi: /opt/hive/apache-hive-3.1.2-bin/bin
divyanshi@divyanshi: /opt/hive/apache-hive-3.1.2-bin/bin

hive> select `user`.name,count(*) as count from tweets_raw group by `user`.name order by count desc limit 10;
Query ID = divyanshi_20210522183743_5051741c-09b6-4aed-831a-53e38e9eedaa
Total jobs = 2
Launching Job 1 out of 2
Number of reduce tasks not specified. Estimated from input data size: 1
In order to change the average load for a reducer (in bytes):
  set hive.exec.reducers.bytes.per.reducer=<number>
In order to limit the maximum number of reducers:
  set hive.exec.reducers.max=<number>
In order to set a constant number of reducers:
  set mapreduce.job.reduces=<number>
Starting Job = job_1621661702117_0021, Tracking URL = http://divyanshi:8088/proxy/application_1621661702117_0021/
Kill Command = /opt/hadoop/hadoop-3.3.0/bin/mapred job -kill job_1621661702117_0021
Hadoop job information for Stage-1: number of mappers: 1; number of reducers: 1
2021-05-22 18:37:51.037 Stage-1 map = 0%, reduce = 0%
2021-05-22 18:37:58.192 Stage-1 map = 100%, reduce = 0%, Cumulative CPU 6.78 sec
2021-05-22 18:38:03.298 Stage-1 map = 100%, reduce = 100%, Cumulative CPU 9.81 sec
MapReduce Total cumulative CPU time: 9 seconds 810 msec
Ended Job = job_1621661702117_0021
Launching Job 2 out of 2
Number of reduce tasks determined at compile time: 1
In order to change the average load for a reducer (in bytes):
  set hive.exec.reducers.bytes.per.reducer=<number>
In order to limit the maximum number of reducers:
  set hive.exec.reducers.max=<number>
In order to set a constant number of reducers:
  set mapreduce.job.reduces=<number>
Starting Job = job_1621661702117_0022, Tracking URL = http://divyanshi:8088/proxy/application_1621661702117_0022/
Kill Command = /opt/hadoop/hadoop-3.3.0/bin/mapred job -kill job_1621661702117_0022
Hadoop job information for Stage-2: number of mappers: 1; number of reducers: 1
2021-05-22 18:38:15.465 Stage-2 map = 0%, reduce = 0%
2021-05-22 18:38:19.577 Stage-2 map = 100%, reduce = 0%, Cumulative CPU 2.01 sec
2021-05-22 18:38:24.689 Stage-2 map = 100%, reduce = 100%, Cumulative CPU 4.37 sec
MapReduce Total cumulative CPU time: 4 seconds 370 msec
Ended Job = job_1621661702117_0022
MapReduce Jobs Launched:
Stage-Stage-1: Map: 1 Reduce: 1 Cumulative CPU: 9.81 sec HDFS Read: 135982491 HDFS Write: 301892 SUCCESS
Stage-Stage-2: Map: 1 Reduce: 1 Cumulative CPU: 4.37 sec HDFS Read: 309499 HDFS Write: 349 SUCCESS
Total MapReduce CPU Time Spent: 14 seconds 180 msec
OK
TSLA bot      761
```

```
divyanshi@divyanshi: /opt/hive/apache-hive-3.1.2-bin/bin
divyanshi@divyanshi: /opt/hive/apache-hive-3.1.2-bin/bin

hive> select `user`.name,count(*) as count from tweets_raw group by `user`.name order by count desc limit 10;
Starting Job = job_1621661702117_0021, Tracking URL = http://divyanshi:8088/proxy/application_1621661702117_0021/
Kill Command = /opt/hadoop/hadoop-3.3.0/bin/mapred job -kill job_1621661702117_0021
Hadoop job information for Stage-1: number of mappers: 1; number of reducers: 1
2021-05-22 18:37:51.037 Stage-1 map = 0%, reduce = 0%
2021-05-22 18:37:58.192 Stage-1 map = 100%, reduce = 0%, Cumulative CPU 6.78 sec
2021-05-22 18:38:03.298 Stage-1 map = 100%, reduce = 100%, Cumulative CPU 9.81 sec
MapReduce Total cumulative CPU time: 9 seconds 810 msec
Ended Job = job_1621661702117_0021
Launching Job 2 out of 2
Number of reduce tasks determined at compile time: 1
In order to change the average load for a reducer (in bytes):
  set hive.exec.reducers.bytes.per.reducer=<number>
In order to limit the maximum number of reducers:
  set hive.exec.reducers.max=<number>
In order to set a constant number of reducers:
  set mapreduce.job.reduces=<number>
Starting Job = job_1621661702117_0022, Tracking URL = http://divyanshi:8088/proxy/application_1621661702117_0022/
Kill Command = /opt/hadoop/hadoop-3.3.0/bin/mapred job -kill job_1621661702117_0022
Hadoop job information for Stage-2: number of mappers: 1; number of reducers: 1
2021-05-22 18:38:15.465 Stage-2 map = 0%, reduce = 0%
2021-05-22 18:38:19.577 Stage-2 map = 100%, reduce = 0%, Cumulative CPU 2.01 sec
2021-05-22 18:38:24.689 Stage-2 map = 100%, reduce = 100%, Cumulative CPU 4.37 sec
MapReduce Total cumulative CPU time: 4 seconds 370 msec
Ended Job = job_1621661702117_0022
MapReduce Jobs Launched:
Stage-Stage-1: Map: 1 Reduce: 1 Cumulative CPU: 9.81 sec HDFS Read: 135982491 HDFS Write: 301892 SUCCESS
Stage-Stage-2: Map: 1 Reduce: 1 Cumulative CPU: 4.37 sec HDFS Read: 309499 HDFS Write: 349 SUCCESS
Total MapReduce CPU Time Spent: 14 seconds 180 msec
OK
TSLA bot      761
stock trader  454
Joe Gambiste  452
RAJU RAY      417
FINAPSE 415
tarik        413
Maaxxo069    406
MacHash 341
Right Stock Alerts 265
Jecika 199
Time taken: 41.965 seconds, Fetched: 10 row(s)
hive>
```

**Time Taken: 41.965 seconds**

## SparkSQL

```
In [22]: start = time.time()
q5 = spark.sql("SELECT count(*) as count, user.name from tweets1 where user.name is not null group by user.name order by count desc")
q5.show()
end = time.time()
print("Time taken = " + str(end - start))

+---+-----+
|count|      name|
+---+-----+
| 761|      TSLA bot|
| 454| stock trader|
| 452| Joe Gambiste|
| 417| RAJU RAY|
| 415| FINAPSE|
| 413| tarik|
| 406| Maaxxo069|
| 341| MacHash|
| 265| Right Stock Alerts|
| 190| Jecika|
+---+-----+
Time taken = 1.505272388458252
```

Time Taken: 1.505272388458252

- **Query:** select `user`.name,retweet\_count from tweets\_raw where retweet\_count in (select max(retweet\_count) from tweets);

## Hive

```
divyanshi@divyanshi: /opt/hive/apache-hive-3.1.2-bin/bin
divyanshi@divyanshi: /opt/hive/apache-hive-3.1.2-bin/bin

Total MapReduce CPU Time Spent: 3 seconds 80 msec
OK
Time taken: 13.037 seconds
hive> select `user`.name,retweet_count from tweets_raw where retweet_count in (select max(retweet_count) from tweets);
Query ID = divyanshi_20210522183413_6ee6d0e5-8d68-4bb5-816e-3ab60836f977
Total jobs = 4
Launching Job 1 out of 4
Number of reduce tasks determined at compile time: 1
In order to change the average load for a reducer (in bytes):
  set hive.exec.reducers.bytes.per.reducer=<number>
In order to limit the maximum number of reducers:
  set hive.exec.reducers.max=<number>
In order to set a constant number of reducers:
  set mapreduce.job.reduces=<number>
Starting Job = job_1621661702117_0019, Tracking URL = http://divyanshi:8088/proxy/application_1621661702117_0019/
Kill Command = /opt/hadoop/hadoop-3.3.0/bin/mapred job -kill job_1621661702117_0019
Hadoop job information for Stage-2: number of mappers: 1; number of reducers: 1
2021-05-22 18:34:22,390 Stage-2 map = 0%, reduce = 0%
2021-05-22 18:34:26,515 Stage-2 map = 100%, reduce = 0%, Cumulative CPU 2.27 sec
2021-05-22 18:34:31,619 Stage-2 map = 100%, reduce = 100%, Cumulative CPU 5.37 sec
MapReduce Total cumulative CPU time: 5 seconds 370 msec
Ended Job = job_1621661702117_0019
Stage-6 is selected by condition resolver.
Stage-7 is filtered out by condition resolver.
Stage-1 is filtered out by condition resolver.

SLF4J: See http://www.slf4j.org/codes.html#multiple_bindings for an explanation.
SLF4J: Actual binding is of type [org.apache.logging.slf4j.Log4jLoggerFactory]
2021-05-22 18:34:37    Uploaded 1 File to: /tmp/divyanshi/_ai0e8113-c5ec-465e-a373-3a0bd3b9816d/hive_2021-05-22_18-34-13_519_2948164305926656225-1-local-1
0005/HashTable-Stage-3/MapJoin-mapfile31--hashtable (280 bytes)
2021-05-22 18:34:37    End of local task; Time Taken: 0.987 sec.
Execution completed successfully
MapredLocal task succeeded
Launching Job 3 out of 4
Number of reduce tasks is set to 0 since there's no reduce operator
Starting Job = job_1621661702117_0020, Tracking URL = http://divyanshi:8088/proxy/application_1621661702117_0020/
Kill Command = /opt/hadoop/hadoop-3.3.0/bin/mapred job -kill job_1621661702117_0020
Hadoop job information for Stage-3: number of mappers: 1; number of reducers: 0
2021-05-22 18:34:46,753 Stage-3 map = 0%, reduce = 0%
2021-05-22 18:34:53,919 Stage-3 map = 100%, reduce = 0%, Cumulative CPU 7.47 sec
MapReduce Total cumulative CPU time: 7 seconds 470 msec
Ended Job = job_1621661702117_0020
MapReduce Jobs Launched:
Stage-Stage-2: Map: 1 Reduce: 1 Cumulative CPU: 5.37 sec HDFS Read: 11098591 HDFS Write: 116 SUCCESS
Stage-Stage-3: Map: 1 Cumulative CPU: 7.47 sec HDFS Read: 135978774 HDFS Write: 113 SUCCESS
Total MapReduce CPU Time Spent: 12 seconds 840 msec
OK
```

```
Raoul Pal      623
Time taken: 41.486 seconds, Fetched: 1 row(s)
hive> []
```

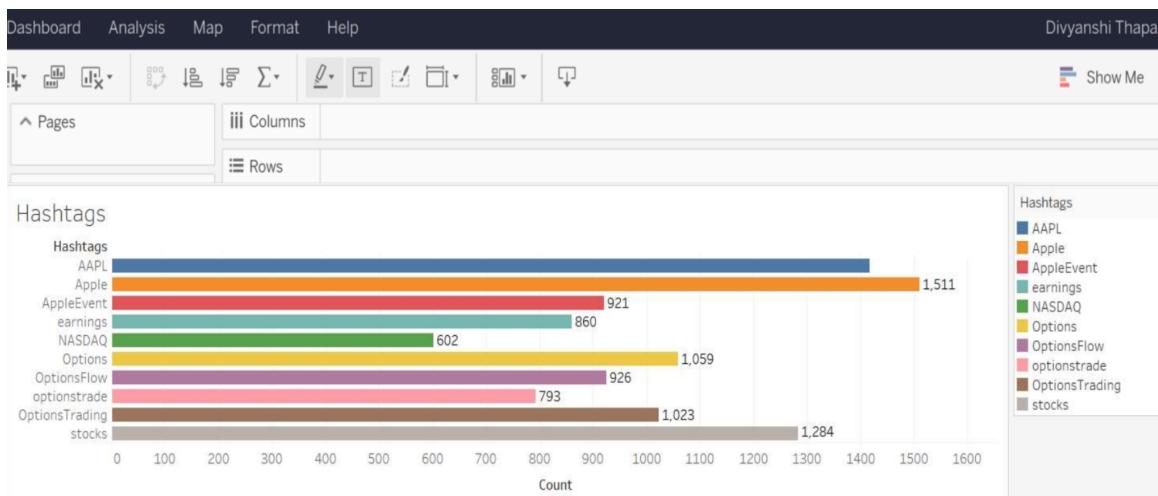
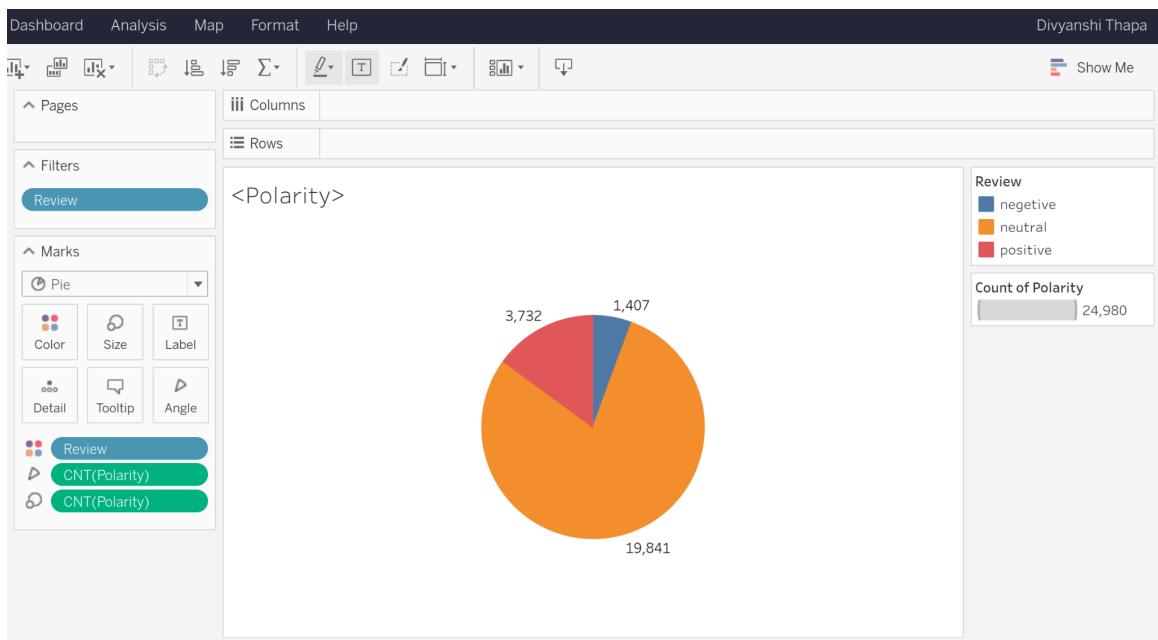
Time Taken: 41.486 seconds

## SparkSQL

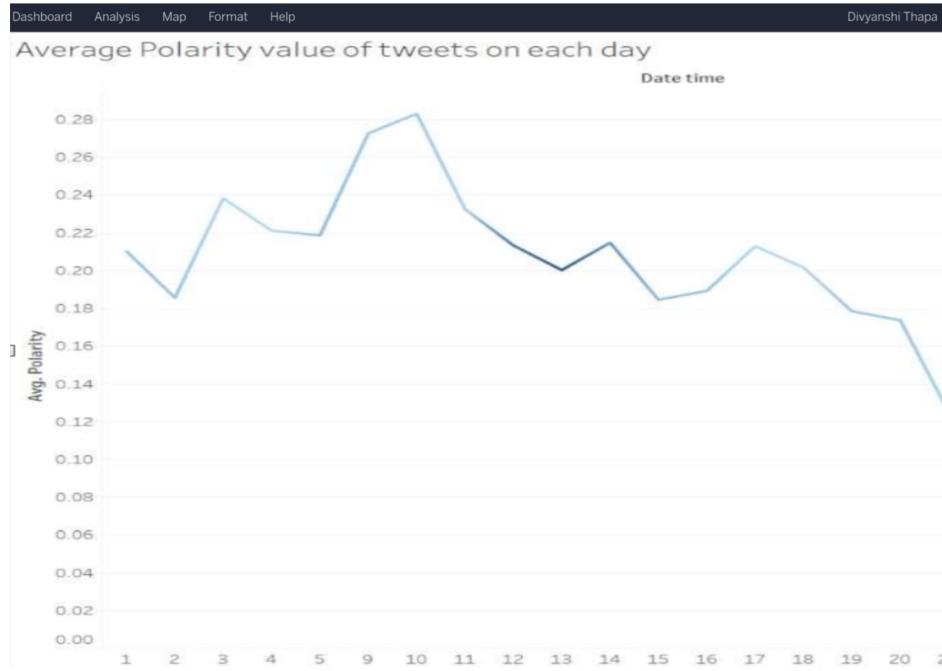
```
In [24]: time.time()
k.sql("select user.name, retweet_count from tweets1 where retweet_count in (select max(retweet_count) from tweets1);"
e.time()
me taken = " + str(end - start))
+-----+-----+
| name|retweet_count|
+-----+-----+
|[Raoul Pal] |623|
+-----+-----+
Time taken = 1.3617329597473145
```

Time Taken: 1.3617329597473145 seconds

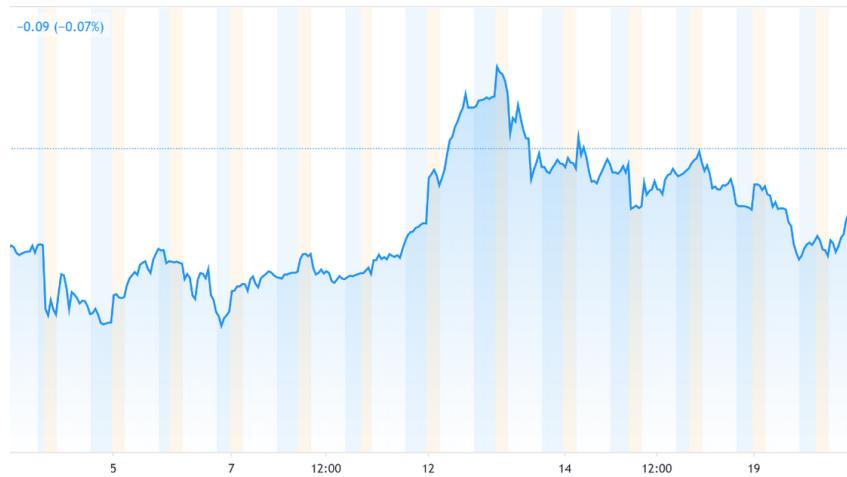
## Visualization



## Our predicted trend



## Original trend



# **Conclusion**

In general, the task of stock market prediction is quite challenging, as the data keeps on changing. Data Analysis and Big Data are on the cusp of completely revolutionising how the stock markets will function. Using this system, we successfully analyzed the data extracted from the twitter social opinions based on stock ticker with the help of hadoop ecosystem. In the above results, we are able to get the trend after analysis which can be used for predictions of movement of stock. This trend might help for the investors in future investments. Along with prediction we also compared the efficiency of different big data tools and deduced the best among them. In our case, Apache SparkSQL performed better than Hive when it came to querying big data. This comparative analysis can be of use for other big data prediction models also.

# References

- 1) <https://www.ijert.org/big-data-analysis-in-stock-market-prediction#:~:text=Abstract%20Big%20data%20analytics%20can,data%2C%20which%20is%20hidden%20otherwise.>
- 2) <https://www.smartdatacollective.com/big-data-analytics-has-potential-to-massively-disrupt-stock-market/>
- 3) <https://ieeexplore.ieee.org/document/8669649>
- 4) <https://github.com/topics/twitter-sentiment-analysis>
- 5) <https://www.intechopen.com/books/e-business-higher-education-and-intelligence-applications/recent-advances-in-stock-market-prediction-using-text-mining-a-survey>
- 6) <https://link.springer.com/article/10.1007/s12652-020-01839-w>