

TWESPA: A MODERN APPROACH TO TWITTER SENTIMENT ANALYSIS

Isha Rane, Sagarika Sardesai , Rucha Tallur

Department of Computer Engineering , St. Xavier Institute of Engineering, Mumbai University, Mumbai, India

Received 21 April 2022 Received in revised form 23 April 2022 Accepted 24 April 2022

Available online 03 May 2022

ABSTRACT

Sentiment analysis has become a stimulating field of research. The term sentiment refers to the emotions or thought of certain issues. Further, it's also considered as an immediate application for opinion mining. The large amount of tweets which appear daily makes Twitter an upscale source of written data and one among the foremost vital data volumes. This data has different purposes , which includes business, industrial or social aims which are consistent with the info requirement and it need to be processed. Actually, the quantity of knowledge, which is very large , grows rapidly per second and this is often called big data which needs special processing techniques and high computational power so as to perform the specified mining tasks. During this work, we perform a sentiment analysis with the assistance of PySpark framework, an interface for Apache Spark in Python - an open-source distributed processing platform which utilizes distributed memory abstraction. The purpose of using PySpark is that one can run applications parallelly on the distributed cluster (multiple nodes). The effectiveness of our proposed approach is proved against other approaches achieving better classification results when using Naïve Bayes, Logistic Regression and Decision trees classification algorithms. Finally, our solution estimates the performance of Apache Spark concerning its scalability.

Keywords— Python, Pyspark, Machine learning, tweepy, tokenization.

1. INTRODUCTION

Now social networks sites make the entire world a small village, where users can share their views, feelings, experiences, advice through those sites so that others can get help from these. Since many of us use social media daily, a huge quantity of comments, opinion, article have been created. An automatic method for studying and categorizing users' views and opinion in social networks like Twitter is highly essential.

This considered as the main tool for gathering direct information from users. Sentiment analysis is the method of classifying texts or documents for keeping their(Social network's) polarity[1] . Sentiment analysis can be described as a major branch of Natural Language Processing (NLP), its aim to identify the meaning from a document in order to discover the polarity of the text[2,3]. For the sentiment analysis, we focus our attention in the direction of the Twitter, where users can communicate with each other or share their opinions in short blogs.

A Large number of variety of text posts are shared on twitter and it increases every day. The fast vast data growth make the already existing databases difficult to handle an extensive amount of data in a short time.

Moreover, these databases type designed to process structured data but there is a limitation on it when dealing with huge data. So the conventional solutions are not helpful for organizations to manage and process unstructured or large data. Frameworks, such as Hadoop, Apache Spark, Apache flume and distributed data storages like Hadoop Distributed File System (HDFS), Cassandra and HBase are being very widespread, as they are designed in a manner which facilitates the process of huge amounts of big data and makes it almost effortless[4-6]

II. MATERIALS AND METHODS

A. Authentication: For fetching tweets through Twitter API, one needs to register an App through their twitter account.

1.Open the link and click the button: 'Create New App'.

2.Fill the application details. You can leave the call back url field empty.

3.Once the app is created, you will be redirected to the app page

4.Open the 'Keys and Access Tokens' tab.

5.Copy 'Consumer Key', 'Consumer Secret', 'Access token' and 'Access Token Secret'

B. Building the Twitter HTTP Client.

In this step, we get the tweets from Twitter API using Python and pass them to the Spark Streaming instance. First, we created a file called `Twitter_Data_Streaming.py` and then we'll import libraries i.e `import socket`, `import sys`, `import requests`, `import requests_oauthlib`, `import json` and most importantly `import tweepy`. Tweepy is publicly released, facilitated on GitHub and enables Python to talk with Twitter stage and utilize its API. Tweepy supports accessing Twitter via Basic Authentication and therefore the newer methodology, OAuth. Twitter has stopped accepted Basic Authentication thus OAuth is currently the sole thing to use the Twitter API

C. Streaming Tweets

Streaming the data from the Twitter API requires creating a listening TCP socket in the local machine (server) on a predefined local IP address and port. A socket consists of a server-side, which is our local machine, and a client-side, which is Twitter API. The open socket from the server-side listens for connections from the client. When the client-side is up and running, the socket will receive the data from the Twitter API based on the topic or keywords defined in the Stream Listener instance.

Now, we'll make the app host socket connections that spark will connect with. We configured the IP here to be localhost as all will run on the same machine and the port 5559. Then we'll getting the tweets from Twitter and pass its response along with the socket connection to `send_tweets_to_spark` for sending the tweets to Spark

```

In [4]: import tweepy
        from tweepy import OAuthHandler
        from IPython.display import display
        from tweepy import Stream
        import emoji
        import json
        from bs4 import BeautifulSoup

        import socket

        consumer_key='6JE4WgZWQcRJddXjv3NZ2c8Ub'
        consumer_secret='g9fzIjDZV9eDwfbEe4iLOFv8tmwcE6zx9i3tB4gBhrPDx8yJ8k'
        access_token = '1644365797-zK4tqxUWYxg7DercDDCdDN13z2BZJCE9sID1hg8'
        access_secret='wQsopDhueX38pS0hYV0ZzNidLnLnia1VpWTgyweE6ccEf'

        def twitter_Access():
            auth = OAuthHandler(consumer_key, consumer_secret)
            auth.set_access_token(access_token, access_secret)
            api = tweepy.API(auth, wait_on_rate_limit=True)
            return api

        def send_tweets_to_spark(http_resp, tcp_connection):
            try:
                full_tweet = json.loads(http_resp)
                tweet_text = full_tweet['text']
                print("Tweet Text: " + tweet_text)
                print("-----")
                tweet_screen_name = "SN:"+full_tweet['user']['screen_name']
                print("SCREEN NAME IS : " + tweet_screen_name)
                print("-----")
                # source = full_tweet['source']
    
```

```

#Streaming tweets
class listener(Stream):
    def on_data(self, data):
        y = json.loads(data)
        # print(y['text'])
        send_tweets_to_spark(data, conn)
        return True
    def on_error(self, status):
        print(status)

TCP_IP = 'localhost'
TCP_PORT = 5559
conn = None
s = socket.socket(socket.AF_INET, socket.SOCK_STREAM)
s.bind((TCP_IP, TCP_PORT))
s.listen(1)
print("Waiting for TCP connection...")
conn, addr = s.accept()
print("Connected... Starting getting tweets.")

twitter_stream = listener(
    consumer_key, consumer_secret,
    access_token, access_secret
)
twitter_stream.filter(track=["#"]) #specify what to search as track parameter
    
```

D. Setting Up Our PySpark Streaming Application

In this part we make a different read_tweets.py file which will do real-time processing for the incoming tweets and extract the hashtags from them. Here we import some libraries like SparkConf and Spark Context from PySpark. We have created an instance of Spark Context sc, then we created the Streaming Context ssc from sc with a batch interval two seconds that will do the transformation on all streams received every two seconds. We also set the log level to ERROR in order to disable most of the logs that Spark writes.

We defined a checkpoint here in order to allow periodic RDD checkpointing. Then we define our main DStream. Data Stream that will connect to the socket server we created before on port 5559 and read the tweets from that port. Each record in the DStream will be a tweet. We split all the tweets into words and put them in words RDD. Then we filtered only hashtags from all words and mapped them to pair of (hashtag, 1) and put them in hashtags RDD. Then we calculate how many times the hashtag has been mentioned. We need to calculate the counts across all the batches, so we used another function

called update State ByKey, as this function allows us to maintain the state of RDD while updating it with new data. This way is called Stateful Transformation.

The update State By Key takes a function as a parameter called the update function. It runs on each item in RDD and does the desired logic. In our case, we've created an update function called aggregate_tags_count that will sum all the new values for each hashtag and add them to the total_sum that is the sum across all the batches and save the data into tags_totals RDD. Then we do processing on tags_totals RDD in every batch in order to convert it to temp table using Spark SQL Context and then perform a select statement in order to retrieve the top ten hashtags with their counts and put them into hashtag_counts_df data frame. Then next step in our Spark application is to send the hashtag_counts_df data frame to the dashboard application. So we'll convert the data frame into two arrays, one for the hashtags and the other for their counts. Then we'll send them to the dashboard application through the REST API.

jupyter read_tweets Last Checkpoint: 5 hours ago (autosaved) Python 3 (ipykernel) Logout

File Edit View Insert Cell Kernel Widgets Help Trusted

```

pass

In [5]: # create spark configuration
conf = SparkConf()
conf.setAppName("TwitterStreamApp")

Out[5]: <pyspark.conf.SparkConf at 0x205b2d8bf40>

In [6]: # create spark context with the above configuration
sc = SparkContext(conf=conf)
sc.setLogLevel("ERROR")

In [7]: # create the Streaming Context from the above spark context with interval size 2 seconds
ssc = StreamingContext(sc, 2)

# setting a checkpoint to allow RDD recovery
ssc.checkpoint("checkpoint_TwitterApp")

# read data from port 5556
dataStream = ssc.socketTextStream("localhost",5559)

In [8]: # split each tweet into words
words = dataStream.flatMap(lambda line: line.split(" "))

# filter the words to get only hashtags, then map each hashtag to be a pair of (hashtag,1)
hashtags = words.filter(lambda x: "#" in x and len(x)>1).map(lambda x: (x, 1))
# hashtags = words.filter(lambda w: '#' in w).map(lambda x: (x, 1))

# adding the count of each hashtag to its last count
tags_totals = hashtags.updateStateByKey(aggregate_tags_count)

```

jupyter read_tweets Last Checkpoint: 5 hours ago (autosaved) Python 3 (ipykernel) Logout

File Edit View Insert Cell Kernel Widgets Help Trusted

```

pass

In [5]: # create spark configuration
conf = SparkConf()
conf.setAppName("TwitterStreamApp")

Out[5]: <pyspark.conf.SparkConf at 0x205b2d8bf40>

In [6]: # create spark context with the above configuration
sc = SparkContext(conf=conf)
sc.setLogLevel("ERROR")

In [7]: # create the Streaming Context from the above spark context with interval size 2 seconds
ssc = StreamingContext(sc, 2)

# setting a checkpoint to allow RDD recovery
ssc.checkpoint("checkpoint_TwitterApp")

# read data from port 5556
dataStream = ssc.socketTextStream("localhost",5559)

In [8]: # split each tweet into words
words = dataStream.flatMap(lambda line: line.split(" "))

# filter the words to get only hashtags, then map each hashtag to be a pair of (hashtag,1)
hashtags = words.filter(lambda x: "#" in x and len(x)>1).map(lambda x: (x, 1))
# hashtags = words.filter(lambda w: '#' in w).map(lambda x: (x, 1))

# adding the count of each hashtag to its last count
tags_totals = hashtags.updateStateByKey(aggregate_tags_count)

```

E. Sentiment Analysis of tweets:

Here a function is written to firstly fetch the tweets from the API which are related to the word entered by the user and then the function analyses every tweet. Every word is passed through a function which contains a list of words which are mentioned in a file as bag of words which denote positive ,negative or neutral sentiments. The keyword in the tweet are analysed from this list of words and then it is decided if its impact is positive ,negative or

neutral. Then this returns 0 if neutral,-1 is its negative and +1 if its positive and then it calculates the maximum number of these tweets to conclude if the tweet was positive or negative. The percentage of positive, negative and neutral is also calculated which later on helps to give us the final grade of the tweets of that particular keyword. Lastly a pie chart is displayed which represents the sentiment of the data.

III. RESULTS



Top Trending Twitter Hashtags

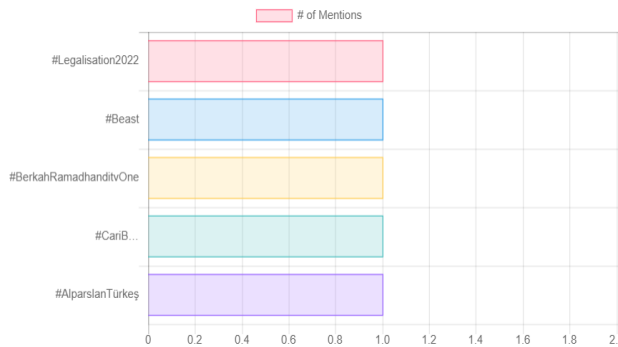


Fig. 3 Front page of the system

The user on opening our webapp is shown a homepage which displays the top trending tweets. Here ,the data analysis of the tweets takes place and it gives us a easy user friendly output in the form of Bar graph. This is a real time analysis which keeps updating itself at constant interval of 2 seconds and displays the tweet

mentions .For better understanding and visual effect we have given different colours to the bar graph representation. When a user hovers over the bar graph he/she can see the number of mentions of that particular hashtag. It keeps updating itself and the highest mentioned hashtag will be at the top most position.

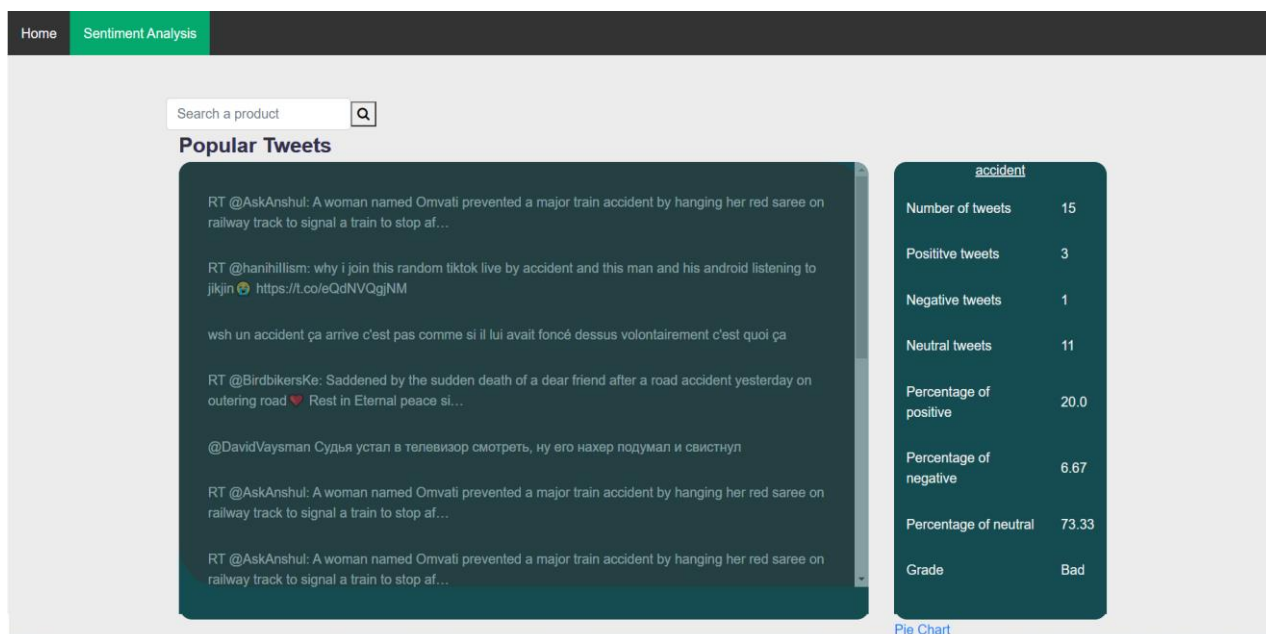


Fig. 4 Twitter Sentiment Analysis of tweets

The second section on our webapp is Sentiment analysis. Once the user clicks on the icon he/she is directed to a new page where any desired word can be given as input in the search box and it will display the results related to that word. All the top recent tweets related to the input word will be shown in the dialog box in the form of tweets. Along with this the sentiment of

all the tweets is shown which consists of the total number of tweets displayed, positive tweets, negative tweets and neutral tweets. The percentage of positive, negative and neutral tweets is also displayed for better understanding. Depending on the polarity of the maximum tweets the final sentiment analysis is done and the grade is given for that particular word.

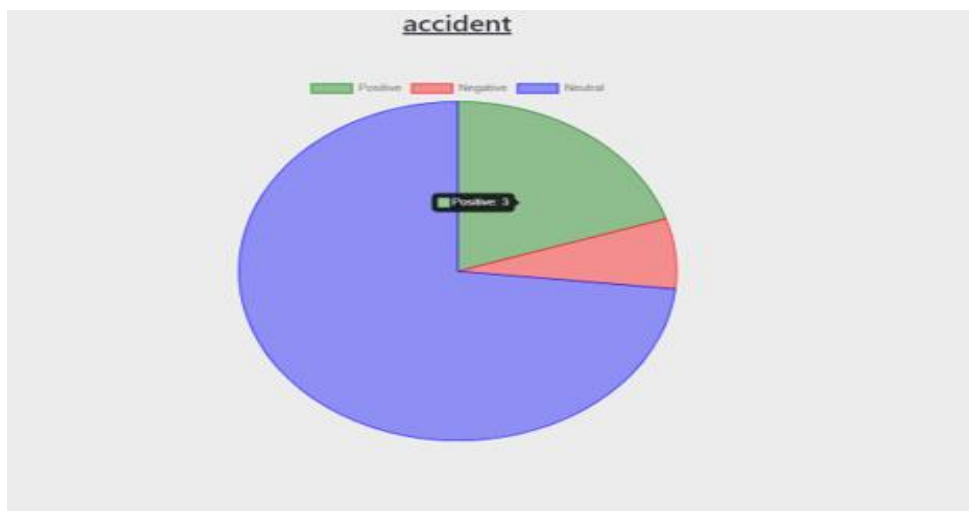


Fig. 5 Pie Chart of Tweets analysed according to sentiments

IV. CONCLUSIONS

Twitter Data in the form of reviews, thoughts, opinion, comments, feedback, and grievance are treated as big data and it cannot be interpreted directly; it should be preprocessed in order to be suitable for mining tasks. In this research, we propose an efficient sentiment analysis technique, utilizing PySpark. The results indicate a significant enhancement in the accuracy of the polarity of the tweets and the additional feature of searching for any desired word is achieved. Along with this we have also did data analysis on the large Twitter datasets and represented the real time data of top trending tweets. From the former outcomes, our system can be described as effective and scalable.

ACKNOWLEDGMENT

We would like to thank our guide "Prof. Kunal Meher." who gave us his valuable suggestions and ideas when we were in need of them. He encouraged us to work on this project. We are also grateful to our college for giving us the opportunity to work with them and providing us the necessary resources for the project.

REFERENCES

- [1]. Vishal.A.Kharde, Prof. Sheetal.Sonawane, 'Sentiment Analysis of Twitter Data: A Survey of Techniques', International Journal of Computer Applications 139(11):2016 5-15,
- [2]. R. Plutchick. "Emotions: A general psycho evolutionary theory." In K.R. Scherer & P. Ekman (Eds) Approaches to emotion. Hillsdale, NJ; Lawrence Erlbaum Associates, 1984
- [3]. P. Basile, V. Basile, M. Nissim, N. Novielli, V. Patti. "Sentiment Analysis of Microblogging Data". Book: Encyclopedia of Social Network Analysis and Mining, Springer 2017 pp-1-17.
- [4]. H. Elzayady, K. M. Badran and G. I. Salama, "Sentiment Analysis on Twitter Data using Apache Spark Framework," 13th International Conference on Computer Engineering and Systems (ICCES), 2018, pp. 171-176, doi: 10.1109/ICCES.2018.8639195.

[5] J. Ranganathan, A. S. Irudayaraj and A. A. Tzacheva, "Action Rules for Sentiment Analysis on Twitter Data Using Spark," IEEE International Conference on Data Mining Workshops (ICDMW), 2017, pp. 51-60, doi: 10.1109/ICDMW.2017.14.

[6]. Baltas, Alexandros, Andreas Kanavos, and Athanasios K. Tsakalidis. "An apache spark implementation for sentiment analysis on twitter data." International Workshop of Algorithmic Aspects of Cloud Computing. Springer, Cham, 2016.