# BIG DATA PROGRAMMING

## PROJECT_BASED EXAM-2

## Spark Streaming Task

**OBJECTIVE:**

Perform Word-Count on Twitter Streaming Data.

**INTRODUCTION:**

- Spark streaming helps to access and process Real-time data with the help of different algorithms like map, Reduce, join etc.
- In this project we use MapReduce algorithm to perform word count on streamed **data.**

**QUESTION 3: (1)**

**IDEA OF THE PROJECT:**

Idea of the project is to apply different concepts learnt in Big-Data-Programming so far. Here we use MapReduce method to calculate word frequency.

**QUESTION 3: (2)**

**USAGE OF PROJECT IN REAL-TIME:**

Spark Streaming Context is used for processing the real-time data streams. In real-time, this idea of streaming helps in prediction, analyzing and data processing workloads etc**.**

**QUESTION 3: (3)**

**IMPLEMENTATION:**

- Initiate a socket object with local machine's IP address and a service specific port number.
- Bind the host and port.
- Make client connection.

```
if __name__ == "__main__":
    new_skt = socket.socket()  # initiate a socket object
    host = "127.0.0.1"  # local machine address
    port = 8085  # specific port for your service.
    new_skt.bind((host, port))  # Binding host and port

    print("Now listening on port: %s" % str(port))

    new_skt.listen(5)  # waiting for client connection.
    c, addr = new_skt.accept()  # Establish connection with client. it returns first a socket object,c, and the address bound to the socket
```

- Create a twitter developer account to get access.
- Get authorization and collect tweets under desired topic.
- Here I have extracted tweets under topic football.

```
def send_tweets(c_socket):
    auth = OAuthHandler(consumer_key, consumer_secret)
    auth.set_access_token(access_token, access_secret)

    twitter_stream = Stream(auth, TweetsListener(c_socket))
    twitter_stream.filter(track=['football'])  # this is the topic we are interested in
```

- Create a class and a model where it extracts only text content from the entire data collected.
- Collect the data and read it to self.
- Run the program and the status gets displayed.

```
# we create this class that inherits from the StreamListener in tweepy StreamListener
class TweetsListener(StreamListener):

    def __init__(self, csocket):
        self.client_socket = csocket

    # we override the on_data() function in StreamListener
    def on_data(self, data):
        # try:
            message = json.loads(data)
            # print(message['text'].encode('utf-8'))
            print(message['text'].encode('utf-8'))
            # c = message['text'] + "\n"
            self.client_socket.send(message['text'].encode('utf-8'))
            # time.sleep(30)
            return True
        # except BaseException as e:
        #     print("Error on_data: %s" % str(e))
        # return True

    def if_error(self, status):
        print(status)
        return True


def send_tweets(c_socket):
```

OUTPUT:



- Later initialize the streaming part.
- Give the same port number as given earlier with which client receives the data from server.
- Use flatMap, where it helps in splitting up a string or a sentence Separated and terminated by a delimiter space as shown in the code below.
- Using Map and reduce methods, split up the words and count their repetition.
- Store the resultant word frequency and print them. Thus, output is obtained.

```python
sc = SparkContext.getOrCreate()
ssc = StreamingContext(sc, 10)
lines = ssc.socketTextStream("localhost", 8085)
Eachwords = lines.flatMap(lambda line: line.split(" "))
KeyPair = Eachwords.map(lambda word: (word, 1))
ResultCount = KeyPair.reduceByKey(lambda x, y: x + y)
ResultCount.pprint()
```

OUTPUT:

```
-------------------------------------------
Time: 2020-07-27 14:31:40
-------------------------------------------
('https://t.co/MO8NfIkyid@dougiecsf', 1)
('@claireeadam', 1)
('of', 13)
('21st', 1)
('are', 2)
('sfa', 1)
('@Matthew_4_Trump:', 1)
('Imagine', 1)
('Makasi', 1)
('Its', 1)
...

[Stage 0:>                (0 + 1) / 1][Stage 67:>              (0 + 5) / 5]20/07/27 14:33:07 WARN RandomBlockReplicationPolicy: Expecting 1 replicas with only 0 peer/s.
20/07/27 14:33:14 WARN BlockManager: Block input-0-1595878394600 replicated to only 0 peer(s) instead of 1 peers
-------------------------------------------
Time: 2020-07-27 14:31:50
-------------------------------------------
('good', 1)
('Manager', 1)
('of', 6)
('treating', 1)
('coupe', 1)
('du', 2)
('watches', 1)
('@guardian_sport:', 2)
('Havertz', 2)
('step', 2)
...

[Stage 0:>                (0 + 1) / 1][Stage 69:>              (0 + 7) / 7]20/07/27 14:33:16 WARN RandomBlockReplicationPolicy: Expecting 1 replicas with only 0 peer/s.
```

## QUESTION 3: (4)

## CHALLENGES FACED:

Had minor issues while writing tweets into socket as it consumes little more time comparatively. And also getting access from twitter also takes time. Irrespective of these two, entire execution procedure went very effectively without any obstacles.

## QUESTION 3: (5)

## MILESTONES AND INTEGRATION OF THE PROJECT:

Work went very effective and smooth without any obstacles. Team split-up made it easy to perform tasks. As there is no dependency for one task to other, individuals performed one task each and accomplished the tasks given.

**TEAM MEMBERS AND CONTRIBUTION:**

Roshini varada -- Hadoop MapReduce Algorithm

 Link -- https://github.com/RoshiniVarada/BDP_Project2/wiki/CASE1

Sarika Reddy Kota -- Spark Data Frames

Link -- https://github.com/RoshiniVarada/BDP_Project2/wiki/CASE2

Pallavi Arikatla -- Spark streaming

 Link -- https://github.com/RoshiniVarada/BDP_Project2/wiki/CASE3

Zakari, Abdulmuhaymin -- Spark Graphx

Link -- https://github.com/RoshiniVarada/BDP_Project2/wiki/CASE4

**Video Link:**

https://youtu.be/UtsiVZaijyg

**PALLAVI ARIKATLA**

**16301244**