# EECS E6893 Big Data Analytic HW3

Chong Hu ch3467

October 31, 2019

**Problem 1.** `Twitter data analysis with Spark Streaming`

1. Hashtag result



Figure 1: Hashtag result preview in Big Query

2. Word Count result



Figure 2: WordCount result preview in Big Query



Figure 3: Finish successfully in terminal

Code: `sparkStreaming.py`

```python
#!/usr/bin/env python
# -*- coding: utf-8 -*-
# Columbia EECS E6893 Big Data Analytics
"""
```

```
5    This module is the spark streaming analysis process.

6


7

8    Usage:

9        If used with dataproc:

10           gcloud dataproc jobs submit pyspark --cluster <Cluster Name> twitterHTTPClient.py

11

12       Create a dataset in BigQurey first using

13           bq mk bigdata_sparkStreaming

14

15       Remeber to replace the bucket with your own bucket name

16


17

18   Todo:

19       1. hashtagCount: calculate accumulated hashtags count

20       2. wordCount: calculate word count every 60 seconds

21           the word you should track is listed below.

22       3. save the result to google BigQuery

23

24   """

25

26   import subprocess

27   import time

28

29   from pyspark import SparkConf, SparkContext

30   from pyspark.sql import SQLContext

31   from pyspark.streaming import StreamingContext

32

33   # global variables

34   bucket = "big_data_storage"

35   output_directory_hashtags = 'gs://{}/hadoop/tmp/bigquery/pyspark_output/hashtagsCount'.format(

36       bucket)

37   output_directory_wordcount = 'gs://{}/hadoop/tmp/bigquery/pyspark_output/wordcount'.format(

38       bucket)

39

40   # output table and columns name

41   output_dataset = 'twitter_analysis'  # the name of your dataset in BigQuery

42   output_table_hashtags = 'hashtags'

43   columns_name_hashtags = ['hashtags', 'count']

44   output_table_wordcount = 'wordcount'

45   columns_name_wordcount = ['word', 'count', 'time']

46

47   # parameter

48   IP = 'localhost'  # ip port

49   PORT = 9001  # port
```

```python
50
51   STREAMTIME = 600   # time that the streaming process runs
52   # STREAMTIME = 20   # for test
53
54   WORD = ['data', 'spark', 'ai', 'movie',
55           'good']   # the words you should filter and do word count
56
57
58   # Helper functions
59   def saveToStorage(rdd, output_directory, columns_name, mode):
60       """
61       Save each RDD in this DStream to google storage
62       Args:
63           rdd: input rdd
64           output_directory: output directory in google storage
65           columns_name: columns name of dataframe
66           mode: mode = "overwirte", overwirte the file
67                 mode = "append", append data to the end of file
68       """
69       if not rdd.isEmpty():
70           (rdd.toDF(columns_name)
71            .write.save(output_directory, format="json", mode=mode))
72
73
74   def saveToBigQuery(sc, output_dataset, output_table, directory):
75       """
76       Put temp streaming json files in google storage to google BigQuery
77       and clean the output files in google storage
78       """
79       files = directory + '/part-*'
80       subprocess.check_call(
81           'bq load --source_format NEWLINE_DELIMITED_JSON '
82           '--replace '
83           '--autodetect '
84           '{dataset}.{table} {files}'.format(
85               dataset=output_dataset, table=output_table, files=files
86           ).split())
87       output_path = sc._jvm.org.apache.hadoop.fs.Path(directory)
88       output_path.getFileSystem(sc._jsc.hadoopConfiguration()).delete(
89           output_path, True)
90
91
92   def hashtagCount(words):
93       """
94       Calculate the accumulated hashtags count sum from the beginning of the stream
```

4

```python
95          and sort it by descending order of the count.
96          Ignore case sensitivity when counting the hashtags:
97              "#Ab" and "#ab" is considered to be a same hashtag
98          You have to:
99          1. Filter out the word that is hashtags.
100             Hashtag usually start with "#" and followed by a serious of alphanumeric
101         2. map (hashtag) to (hashtag, 1)
102         3. sum the count of current DStream state and previous state
103         4. transform unordered DStream to a ordered Dstream
104         Hints:
105             you may use regular expression to filter the words
106             You can take a look at updateStateByKey and transform transformations
107         Args:
108             dstream(DStream): stream of real time tweets
109         Returns:
110             DStream Object with inner structure (hashtag, count)
111         """
112
113         def updateFunc(new_values, last_sum):
114             return sum(new_values) + (last_sum or 0)
115
116         hashtag = words.map(lambda x: x.lower()).filter(
117             lambda x: len(x) > 2 and x[0] == "#").map(
118             lambda x: (x, 1))
119         hashtag_cnt = hashtag.reduceByKey(lambda cnt1, cnt2: cnt1 + cnt2)
120         hashtag_cnt_total = hashtag_cnt.updateStateByKey(updateFunc)
121         return hashtag_cnt_total
122
123
124 def wordCount(words):
125     """
126     Calculte the count of 5 sepcial words for every 60 seconds (window no overlap)
127     You can choose your own words.
128     Your should:
129     1. filter the words
130     2. count the word during a special window size
131     3. add a time related mark to the output of each window, ex: a datetime type
132     Hints:
133         You can take a look at reduceByKeyAndWindow transformation
134         Dstream is a serious of rdd, each RDD in a DStream contains data from a certain interval
135         You may want to take a look of transform transformation of DStream when trying to add a time
136     Args:
137         dstream(DStream): stream of real time tweets
138     Returns:
139         DStream Object with inner structure (word, count, time)
```

5

```python
140        """
141        word_cnt = words.map(lambda x: x.lower()).filter(lambda x: x in WORD).map(
142            lambda x: (x, 1)).reduceByKeyAndWindow(lambda x, y: x + y,
143                                                    lambda x, y: x - y, 60, 60)
144        word_cnt_total = word_cnt.transform(
145            lambda time, rdd: rdd.map(
146                lambda x: (x[0], x[1], time.strftime("%Y-%m-%d %H:%M:%S"))))
147        return word_cnt_total
148
149
150    if __name__ == '__main__':
151        # Spark settings
152        conf = SparkConf()
153        conf.setMaster('local[2]')
154        conf.setAppName("TwitterStreamApp")
155
156        # create spark context with the above configuration
157        sc = SparkContext(conf=conf)
158        sc.setLogLevel("ERROR")
159
160        # create sql context, used for saving rdd
161        sql_context = SQLContext(sc)
162
163        # create the Streaming Context from the above spark context with batch interval size 60 seconds
164        ssc = StreamingContext(sc, 60)
165        # setting a checkpoint to allow RDD recovery
166        ssc.checkpoint("~/checkpoint_TwitterApp")
167
168        # read data from port 9001
169        dataStream = ssc.socketTextStream(IP, PORT)
170        dataStream.pprint()
171
172        words = dataStream.flatMap(lambda line: line.split(" "))
173
174        # calculate the accumulated hashtags count sum from the beginning of the stream
175        topTags = hashtagCount(words)
176        topTags.pprint()
177
178        # Calculte the word count during each time period 6s
179        wordCount = wordCount(words)
180        wordCount.pprint()
181
182        # save hashtags count and word count to google storage
183        # used to save to google BigQuery
184        # You should:
```

```python
185     #   1. topTags: only save the lastest rdd in DStream
186     #   2. wordCount: save each rdd in DStream
187     # Hints:
188     #   1. You can take a look at foreachRDD transformation
189     #   2. You may want to use helper function saveToStorage
190     #   3. You should use save output to output_directory_hashtags, output_directory_wordcount,
191     #      and have output columns name columns_name_hashtags and columns_name_wordcount.
192
193     topTags.foreachRDD(lambda rdd: saveToStorage(rdd, output_directory_hashtags,
194                                                  columns_name_hashtags,
195                                                  mode="overwrite"))
196     wordCount.foreachRDD(
197         lambda rdd: saveToStorage(rdd, output_directory_wordcount,
198                                   columns_name_wordcount, mode="append"))
199     # start streaming process, wait for 600s and then stop.
200     ssc.start()
201     time.sleep(STREAMTIME)
202     ssc.stop(stopSparkContext=False, stopGraceFully=True)
203     # put the temp result in google storage to google BigQuery
204     saveToBigQuery(sc, output_dataset, output_table_hashtags,
205                    output_directory_hashtags)
206     saveToBigQuery(sc, output_dataset, output_table_wordcount,
207                    output_directory_wordcount)
```

twitterHTTPClient.py

```python
1   #!/usr/bin/env python
2   # -*- coding: utf-8 -*-
3   # Columbia EECS E6893 Big Data Analytics
4   """
5   This module is used to pull data from twitter API and send data to
6   Spark Streaming process using socket. It acts like a client of
7   twitter API and a server of spark streaming. It open a listening TCP
8   server socket, and listen to any connection from TCP client. After
9   a connection established, it send streaming data to it.
10
11
12  Usage:
13    If used with dataproc:
14      gcloud dataproc jobs submit pyspark --cluster <Cluster Name> twitterHTTPClient.py
15
16    Make sure that you run this module before you run spark streaming process.
17    Please remember stop the job on dataproc if you no longer want to stream data.
18
```

```python
Todo:
    1. change the credentials to your own

"""

import json
import socket

from tweepy import OAuthHandler
from tweepy import Stream
from tweepy.streaming import StreamListener

# credentials
# replace with your own credentials
ACCESS_TOKEN = '1186641375900647425-Yk6Mr116gbZixYb0pt2bXGi0DWZiZa'  # your access token
ACCESS_SECRET = '8xYHZbwKlF4zMdqHEWuwgO1v0ocHuictsQnCn4ElPT1Ny'  # your access token secret
CONSUMER_KEY = 'Wg4xVrzE0DU8Ey4sYPQv0ntcv'  # your API key
CONSUMER_SECRET = 'FBXUfei2Sp7W7dfuLtK8wbY5BF20KB4Ck5T5IGiE7AoEdWtQZW'  # your API secret key

# the tags to track
tags = ['#', 'bigdata', 'spark', 'ai', 'movie']


class TweetsListener(StreamListener):
    """
    tweets listener object
    """

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

    def on_data(self, data):
        try:
            msg = json.loads(data)
            print('TEXT:{}\n'.format(msg['text']))
            self.client_socket.send(msg['text'].encode('utf-8'))
            return True
        except BaseException as e:
            print("Error on_data: %s" % str(e))
            return False
        # return True

    def on_error(self, status):
        print(status)
```

```python
64          return False
65
66
67  def sendData(c_socket, tags):
68      """
69      send data to socket
70      """
71      auth = OAuthHandler(CONSUMER_KEY, CONSUMER_SECRET)
72      auth.set_access_token(ACCESS_TOKEN, ACCESS_SECRET)
73      twitter_stream = Stream(auth, TweetsListener(c_socket))
74      twitter_stream.filter(track=tags, languages=['en'])
75
76
77  class twitter_client:
78      def __init__(self, TCP_IP, TCP_PORT):
79          self.s = s = socket.socket(socket.AF_INET, socket.SOCK_STREAM)
80          self.s.bind((TCP_IP, TCP_PORT))
81
82      def run_client(self, tags):
83          try:
84              self.s.listen(1)
85              while True:
86                  print("Waiting for TCP connection...")
87                  conn, addr = self.s.accept()
88                  print("Connected... Starting getting tweets.")
89                  sendData(conn, tags)
90                  conn.close()
91          except KeyboardInterrupt:
92              exit
93
94
95  if __name__ == '__main__':
96      client = twitter_client("localhost", 9001)
97      client.run_client(tags)
```