

EECS E6893 Big Data Analytic HW3

Chong Hu ch3467

November 1, 2019

Problem 1. Twitter data analysis with Spark Streaming

1. Hashtag result

hashtags

| Schema | Details | Preview |
|--------|---------|-------------------------|
| 1024 | 29 | #iot |
| 1025 | 30 | #petta |
| 1026 | 34 | #bigdata |
| 1027 | 36 | #viswasam |
| 1028 | 111 | #bigil |
| 1029 | 114 | #ai |
| 1030 | 116 | #sanki. |
| 1031 | 117 | #shahrukhkhan |
| 1032 | 121 | #atlee |
| 1033 | 126 | #sanki |
| 1034 | 235 | #srk |

Figure 1: Hashtag result preview in Big Query

```
-----  
( '#valimai', 2)  
( '#ai', 2)  
( '#mvsales', 1)  
( '#undafoqo', 1)  
( '#ml', 1)  
( '#valimaiinhindi@travizm', 1)  
( '#thalapathy', 1)  
( '#newsfeed', 1)  
( '#intoainews@resistjockey', 1)  
( '#sittizone.', 1)  
...
```

Figure 2: Partial Hashtag result in order in screen output

2. Word Count result

wordcount

| Schema | Details | Preview |
|--------|-------------------------|-----------|
| 25 | 2019-10-31 10:19:00 UTC | 3 good |
| 26 | 2019-10-31 10:15:00 UTC | 4 good |
| 27 | 2019-10-31 10:17:00 UTC | 5 good |
| 28 | 2019-10-31 10:18:00 UTC | 9 good |
| 29 | 2019-10-31 10:10:00 UTC | 10 movie |
| 30 | 2019-10-31 10:18:00 UTC | 118 movie |
| 31 | 2019-10-31 10:17:00 UTC | 121 movie |
| 32 | 2019-10-31 10:13:00 UTC | 124 movie |
| 33 | 2019-10-31 10:14:00 UTC | 126 movie |
| 34 | 2019-10-31 10:12:00 UTC | 128 movie |
| 35 | 2019-10-31 10:19:00 UTC | 131 movie |
| 36 | 2019-10-31 10:15:00 UTC | 135 movie |
| 37 | 2019-10-31 10:16:00 UTC | 139 movie |
| 38 | 2019-10-31 10:11:00 UTC | 142 movie |
| 39 | 2019-10-31 10:17:00 UTC | 1 spark |
| 40 | 2019-10-31 10:18:00 UTC | 2 spark |

Figure 3: WordCount result preview in Big Query

```

.....
Time: 2019-10-31 10:31:00
.....
Waiting on bqjob_r7ecef6063fd029c6_0000016e215e688c_1 ... (0s) Current status: R
ase088c_1 ... (1s) Current status: R
Waiting on bqjob_r25ad3a09c436b116_0000016e215e7aaa_1 ... (0s) Current status: R
ase7aaa_1 ... (1s) Current status: R
Waiting on bqjob_r25ad3a09c436b116_0000016e215e7aaa_1 ...
r25ad3a09c436b116_0000016e215e7aaa_1 ... (5s) Current status: R
5) Current status: R
Job [b35241b7bb874f3680cc3c9d223d50a6] finished successfully.
DriverControlFilesUri: gs://big_data_storage/google-cloud-dataproc-metaInfo/bcc7020e-1f9b-4641-96d4-6c264c9d226b
DriverOutputResourceUri: gs://big_data_storage/google-cloud-dataproc-metaInfo/bcc7020e-1f9b-4641-96d4-6c264c9d226b
JobUuid: aa5b5d7e-ef7f-3817-8336-2e791ca6d092
Placement:
  clusterName: hw3
  clusterUuid: bcc7020e-1f9b-4641-96d4-6c264c9d226b
  pysparkJob:
    mainPythonFileUri: gs://big_data_storage/google-cloud-dataproc-metaInfo/bcc7020e-1f9b-4641-96d4-6c264c9d226b
  reference:
    jobId: b35241b7bb874f3680cc3c9d223d50a6
    projectId: hardy-symbol-252200
  status:
    state: DONE
    stateStartTime: '2019-10-31T10:31:20.509Z'
  statusHistory:
    state: PENDING
    stateStartTime: '2019-10-31T10:09:46.590Z'
    state: SETUP_DONE
    stateStartTime: '2019-10-31T10:09:46.623Z'
    details: Agent reported job success
    state: RUNNING
    stateStartTime: '2019-10-31T10:09:47.053Z'

```

Figure 4: Finish successfully in terminal

Code: sparkStreaming.py

```

1  #!/usr/bin/env python
2  # -*- coding: utf-8 -*-
3  # Columbia EECS E6893 Big Data Analytics
4  """

```

```

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 BigQuery 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/hashtagCount'.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

```

```

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 = "overwrite", overwrite 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

```

```

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 series 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).transform(
121         lambda rdd: rdd.sortBy(lambda x: x[1], ascending=False))
122     return hashtag_cnt_total
123
124
125 def wordCount(words):
126     """
127     Calculate the count of 5 special words for every 60 seconds (window no overlap)
128     You can choose your own words.
129     Your should:
130     1. filter the words
131     2. count the word during a special window size
132     3. add a time related mark to the output of each window, ex: a datetime type
133     Hints:
134         You can take a look at reduceByKeyAndWindow transformation
135         Dstream is a series of rdd, each RDD in a DStream contains data from a certain interval
136         You may want to take a look of transform transformation of DStream when trying to add a time
137     Args:
138         dstream(DStream): stream of real time tweets
139     Returns:

```

```

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

```

```

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

```

twitterHTTPClient.py

```

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

```

```

18
19 Todo:
20     1. change the credentials to your own
21
22 """
23
24 import json
25 import socket
26
27 from tweepy import OAuthHandler
28 from tweepy import Stream
29 from tweepy.streaming import StreamListener
30
31 # credentials
32 # replace with your own credentials
33 ACCESS_TOKEN = '1186641375900647425-Yk6Mr116gbZixYb0pt2bXGi0DWZiZa' # your access token
34 ACCESS_SECRET = '8xYHZbwKlF4zMdqHEWuwg01v0ocHuictsQnCn4ElPT1Ny' # your access token secret
35 CONSUMER_KEY = 'Wg4xVrzEODU8Ey4sYPQv0ntcv' # your API key
36 CONSUMER_SECRET = 'FBXUfei2Sp7W7dfuLtK8wbY5BF20KB4Ck5T5IGiE7AoEdWtQZW' # your API secret key
37
38 # the tags to track
39 tags = ['#', 'bigdata', 'spark', 'ai', 'movie']
40
41
42 class TweetsListener(StreamListener):
43     """
44     tweets listener object
45     """
46
47     def __init__(self, csocket):
48         super(TweetsListener, self).__init__()
49         self.client_socket = csocket
50
51     def on_data(self, data):
52         try:
53             msg = json.loads(data)
54             print('TEXT:{}'.format(msg['text']))
55             self.client_socket.send(msg['text'].encode('utf-8'))
56             return True
57         except BaseException as e:
58             print("Error on_data: %s" % str(e))
59             return False
60         # return True
61
62     def on_error(self, status):

```



```

63         print(status)
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)

```
