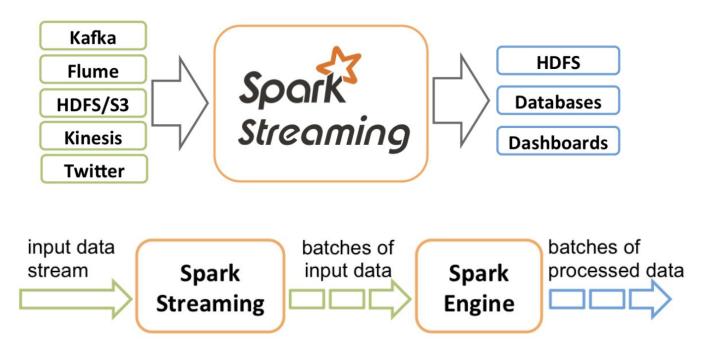


EECS E6893 Big Data Analytics HW3: Twitter data analysis with Spark Streaming

Hritik Jain, hj2533@columbia.edu

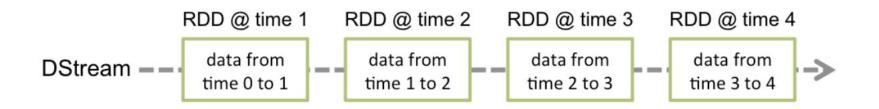
11/20/2020

Spark Streaming

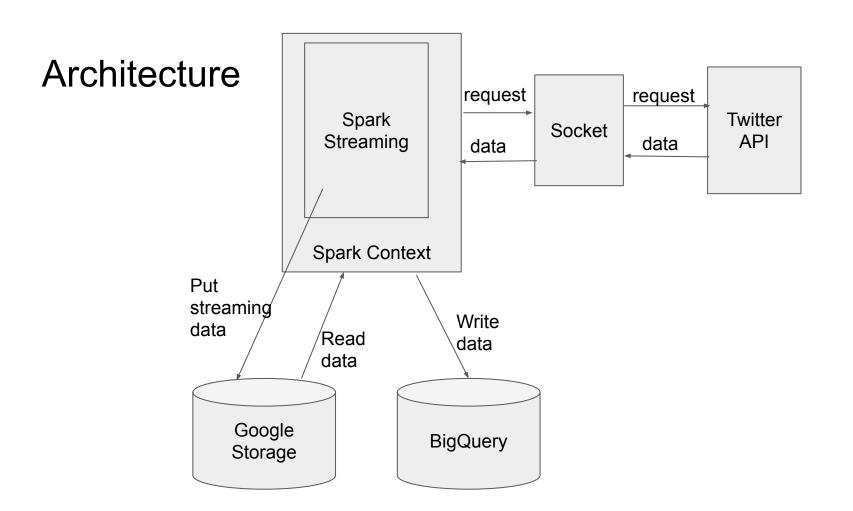


https://spark.apache.org/docs/latest/streaming-programming-guide.html

Dstream

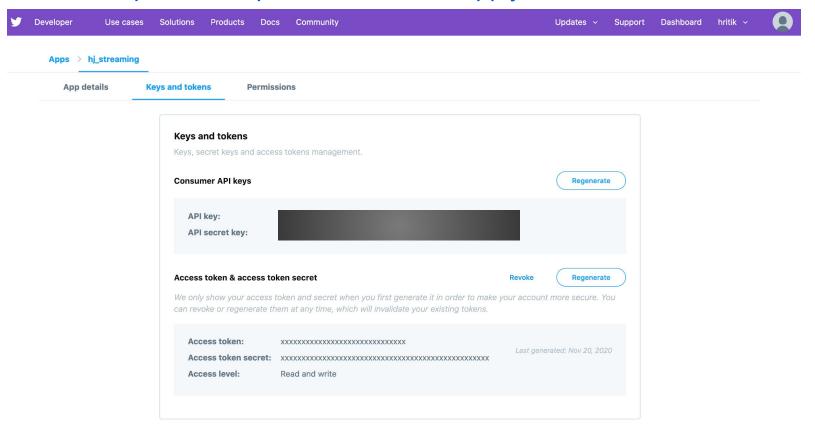


- A basic abstraction provided by Spark Streaming
- Represents a continuous stream of data
- Contains a continuous series of RDDs at different time



Register on Twitter Apps (Do this ASAP)

https://developer.twitter.com/en/apply-for-access.html



Register on Twitter Apps (Do this ASAP)



https://developer.twitter.com/en/apply-for-access.html

```
# credentials
# TODO: replace with your own credentials
ACCESS_TOKEN = ''  # your access token
ACCESS_SECRET = ''  # your access token secret
CONSUMER_KEY = ''  # your API key
CONSUMER_SECRET = ''  # your API secret key
```

Socket

Using TCP source -- need to provide IP and Port for the client to connect

```
class twitter_client:
   def __init__(self, TCP_IP, TCP_PORT):
      self.s = s = socket.socket(socket.AF_INET, socket.SOCK_STREAM)
      self.s.bind((TCP_IP, TCP_PORT))
    def run_client(self, tags):
      try:
        self.s.listen(1)
        while True:
          print("Waiting for TCP connection...")
          conn, addr = self.s.accept()
          print("Connected... Starting getting tweets.")
          sendData(conn,tags)
          conn.close()
      except KeyboardInterrupt:
        exit
```

Spark Streaming

```
if __name__ == '__main__':
                                            Create a local StreamingContext with two
    conf = SparkConf()
                                            working threads and batch interval of 5
    conf.setMaster('local[2]')
    conf.setAppName("TwitterStreamApp")
                                            second.
    sc = SparkContext(conf=conf)
    sc.setLogLevel("ERROR")
    sql_context = SQLContext(sc)
    ssc = StreamingContext(sc, 5)
                                                 Create stream from TCP socket IP localhost
    ssc.checkpoint("~/checkpoint_TwitterApp")
                                                 and Port 9001
    dataStream = ssc.socketTextStream("localhost",9001)
    dataStream.pprint()
```

Spark Streaming

```
ssc.start()
time.sleep(STREAMTIME)
ssc.stop(stopSparkContext=False, stopGraceFully=True)
```

STREAMTIME = 600 # time that the streaming process runs

Start streaming context

Stop after 600 seconds (You can set STREAMTIME to a smaller value at first)

Save results to BigQuery

Start streaming

- 1. Run twitterHTTPClient.py
- 2. Run sparkStreaming.py

You can test sparkStreaming.py multiple times and leave twitterHTTPClient.py running

Stop twitterHTTPClient.py (on job page of the cluster or use gcloud command)

Task1: hashtagCount

```
def hashtagCount(words):
    pass
```

Task2: wordCount

```
WORD = ['data', 'spark', 'ai', 'movie', 'good'] #the words you should filter and do word count
# Helper functions
```

```
def wordCount(words):
    pass
```

Task3: Save results

Create a dataset:

bq mk <Dataset name>

Replace with your own bucket and dataset name:

```
# global variables
bucket = ""  # TODO : replace with your own bucket name
output_directory_hashtags = 'gs://{}/hadoop/tmp/bigquery/pyspark_output/hashtagsCount'.format(bucket)
output_directory_wordcount = 'gs://{}/hadoop/tmp/bigquery/pyspark_output/wordcount'.format(bucket)

# output table and columns name
output_dataset = ''  #the name of your dataset in BigQuery
output_table_hashtags = 'hashtags'
columns_name_hashtags = ['hashtags', 'count']
output_table_wordcount = 'wordcount'
columns_name_wordcount = ['word', 'count', 'timestamp']
```

Task3: Save results

```
# save hashtags count and word count to google storage
# used to save to google BigQuery
# You should:
# 1. topTags: only save the last DStream
# 2. wordCount: save each DStream
# Hints:
# 1. You can take a look at foreachRDD transformation
# 2. You may want to use helper function saveToStorage
# TODO: insert your code here
```

Sample Results

hashtags

| Schema Det | | ails Preview |
|------------|-------|-----------------------|
| Row | count | hashtags |
| 1 | 25 | #valimai |
| 2 | 14 | #ai |
| 3 | 11 | #isapafeelgoodweekend |
| 4 | 8 | #isapawithfeelings |
| 5 | 7 | #วาร์ปสิวะ |
| 6 | 7 | #thalaajith |
| 7 | 6 | #warpsiwax |
| 8 | 5 | #mainemendoza |
| 9 | 5 | #jojorabbit |
| 10 | 4 | #bigdata |
| 11 | 4 | #uncen |

wordcount

Details

Schema

| Scrienta Details Fleview | | | | | |
|--------------------------|--------------------------|-------|------|--|--|
| Row | time | count | word | | |
| 1 | 2019-10-18 23:30:10 UTC | 2 | ai | | |
| 2 | 2019-10-18 23:30:50 UTC | 2 | ai | | |
| 3 | 2019-10-18 23:31:30 UTC | 3 | ai | | |
| 4 | 2019-10-18 23:31:50 UTC | 5 | ai | | |
| 5 | 2019-10-18 23:31:10 UTC | 8 | ai | | |
| 6 | 2019-10-18 23:30:30 UTC | 10 | ai | | |
| 7 | 2019-10-18 23:30:10 UTC | 1 | data | | |
| 8 | 2019-10-18 23:31:50 UTC | 1 | data | | |
| 9 | 2019-10-18 23:30:50 UTC | 1 | data | | |
| 10 | 2019-10-18 23:31:10 UTC | 2 | data | | |
| 11 | 2010-10-18 23-31-30 LIT€ | 2 | annd | | |

Preview

Reference

https://spark.apache.org/docs/latest/streaming-programming-guide.html