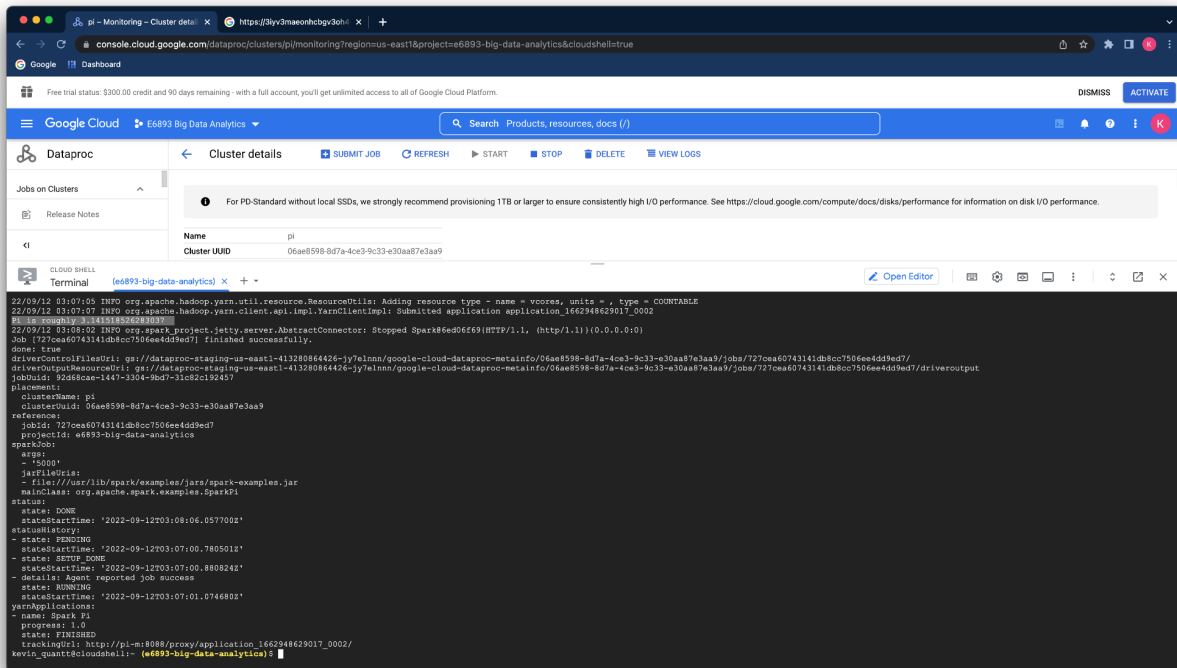EECS E6893 Big Data Analytics

# Homework #0 Report

Name: Yi Zhang

Uni: yz4140

September 2022

# Question 1

(1) As shown below, the pi is successfully estimated via the cluster called pi and the value is correct.



(2) According to the source code provided by the Spark below

Python    Scala    Java

```python
def inside(p):
    x, y = random.random(), random.random()
    return x*x + y*y < 1

count = sc.parallelize(range(0, NUM_SAMPLES)) \
            .filter(inside).count()
print("Pi is roughly %f" % (4.0 * count / NUM_SAMPLES))
```

The transformation is '**filter**', which helps to create a new dataset containing the (x,y) within the circle. The action is '**count**', which counts the number of data points in the filtered dataset.

# Question 2

(1) There are 698 stations whose longitude is between -74.04 and -73.94



(2) There are 11885 bikes available in region 71

```
1  SELECT sum(num_bikes_available) FROM `e6893-big-data-analytics.NYC_Citi_Bike.
   citi_bike_data` where region_id = 71
```

**Query results**

| Row | f0_ |
|-----|-----|
| 1 | 11885 |

(3) The largest capacity is 79 and station 445, 422, and 501 attain the largest capacity.

SELECT max(capacity) FROM `e6893-big-data-analytics.NYC_Citi_Bike.citi_bike_data`

**Query results**

| Row | f0_ |
|-----|-----|
| 1 | 79 |

**▶ RUN**   **⬇ SAVE ▾**   **👤 SHARE ▾**   **🕐 SCHEDULE ▾**   **⚙ MORE ▾**   ✅ This query will p

```
1  SELECT station_id FROM `e6893-big-data-analytics.NYC_Citi_Bike.citi_bike_data`
   where capacity = (select max(capacity) from `e6893-big-data-analytics.
   NYC_Citi_Bike.citi_bike_data`)
2  |
```

Press Alt+F1 for Accessibility Options

## Query results

**⬇ SAVE RESULTS ▾**   **📊 EXPLORE DATA ▾**   ⇕

| JOB INFORMATION | RESULTS | JSON | EXECUTION DETAILS |

| Row | station_id |
|---|---|
| 1 | 445 |
| 2 | 422 |
| 3 | 501 |

**C REFRESH**   ⌃

# Question 3

```python
word_count.py ×
E6893 Big Data Analytics > word_count.py > ...
  1  import sys, pyspark, re, nltk
  2  from nltk.corpus import stopwords
  3  nltk.download('stopwords')
  4
  5  # preprocessing utility function
  6  stop_words = stopwords.words('english')
  7
  8  def preprocess(lines):
  9      def word_preprocess(w):
 10          return re.sub("[^A-Za-z0-9]+", "", w.lower())
 11      temp = [word_preprocess(word) for word in lines.split()]
 12      return [w for w in temp if w not in stop_words and w != ""]
 13
 14
 15  if len(sys.argv) != 3:
 16      raise Exception("Exactly 2 arguments are required: <inputUri> <outputUri>")
 17
 18  inputUri = sys.argv[1]
 19  outputUri = sys.argv[2]
 20
 21  sc = pyspark.SparkContext()
 22  lines = sc.textFile(inputUri)
 23
 24  PROCESS = False
 25  words = lines.flatMap(preprocess) if PROCESS else lines.flatMap(lambda l: l.split())
 26
 27  wordCount = words.map(lambda w: (w, 1)).reduceByKey(lambda c1, c2: c1 + c2)
 28  wordCount_sorted = wordCount.map(lambda x: (x[1], x[0])).sortByKey(False)
 29  wordCount_sorted.collect()
 30
 31  print()
 32  print(">>>>>>>>>>> Top 10 frequent words without text preprocessing >>>>>>>>>>>")
 33  print(wordCount_sorted.take(10))
 34  print(">>>>>>>>>>> End >>>>>>>>>>> \n")
 35
 36  wordCount_sorted.saveAsTextFile(outputUri)
```

The code above shows the overall workflow of the program. In this exercise, the transformation operations are **'flatMap', 'map'**. The action operations are **'reduceByKey', 'sortByKey', 'collect', 'take'**.

(1) As the result shows, the top 10 frequent words are 'the', 'and', 'of', 'to', 'I', 'a', 'you', 'in', 'is', 'my'.

```
TERMINAL  GITLENS  JUPYTER  PROBLEMS  OUTPUT  DEBUG CONSOLE                                    [>] zsh - E6893 Big Data Analytics  + ∨  □  🗑  ∨  ⤬

(base) Yi ~/Desktop/E6893 Big Data Analytics >>
(base) Yi ~/Desktop/E6893 Big Data Analytics >>
(base) Yi ~/Desktop/E6893 Big Data Analytics >> gcloud dataproc jobs submit pyspark word_count.py --cluster=shakes-count --region=us-east1 -- gs://6893_course_data/shakes
_text/ gs://6893_course_data/count_result/
Job [ae918338be4a4e629e20f1fc128625da] submitted.
Waiting for job output...
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data]   Package stopwords is already up-to-date!
22/09/15 16:58:09 INFO org.apache.spark.SparkEnv: Registering MapOutputTracker
22/09/15 16:58:09 INFO org.apache.spark.SparkEnv: Registering BlockManagerMaster
22/09/15 16:58:10 INFO org.apache.spark.SparkEnv: Registering OutputCommitCoordinator
22/09/15 16:58:10 INFO org.apache.spark_project.jetty.util.log: Logging initialized @4660ms to org.spark_project.jetty.util.log.Slf4jLog
22/09/15 16:58:10 INFO org.apache.spark_project.jetty.server.Server: jetty-9.4.z-SNAPSHOT; built: unknown; git: unknown; jvm 1.8.0_332-b09
22/09/15 16:58:10 INFO org.apache.spark_project.jetty.server.Server: Started @4838ms
22/09/15 16:58:10 INFO org.apache.spark_project.jetty.server.AbstractConnector: Started ServerConnector@461f2d5b{HTTP/1.1, (http/1.1)}{0.0.0.0:35533}
22/09/15 16:58:11 INFO org.apache.hadoop.yarn.client.RMProxy: Connecting to ResourceManager at shakes-count-m/10.142.0.8:8032
22/09/15 16:58:11 INFO org.apache.hadoop.yarn.client.AHSProxy: Connecting to Application History server at shakes-count-m/10.142.0.8:10200
22/09/15 16:58:11 INFO org.apache.hadoop.conf.Configuration: resource-types.xml not found
22/09/15 16:58:11 INFO org.apache.hadoop.yarn.util.resource.ResourceUtils: Unable to find 'resource-types.xml'.
22/09/15 16:58:11 INFO org.apache.hadoop.yarn.util.resource.ResourceUtils: Adding resource type - name = memory-mb, units = Mi, type = COUNTABLE
22/09/15 16:58:11 INFO org.apache.hadoop.yarn.util.resource.ResourceUtils: Adding resource type - name = vcores, units = , type = COUNTABLE
22/09/15 16:58:13 INFO org.apache.hadoop.yarn.client.api.impl.YarnClientImpl: Submitted application application_1663254897326_0008
22/09/15 16:58:20 INFO org.apache.hadoop.mapred.FileInputFormat: Total input files to process : 1

>>>>>>>>>>> Top 10 frequent words without text preprocessing >>>>>>>>>>>
[(620, 'the'), (427, 'and'), (396, 'of'), (367, 'to'), (326, 'I'), (256, 'a'), (193, 'you'), (190, 'in'), (185, 'is'), (170, 'my')]
>>>>>>>>>>> End >>>>>>>>>>>

22/09/15 16:58:29 INFO org.apache.spark_project.jetty.server.AbstractConnector: Stopped Spark@461f2d5b{HTTP/1.1, (http/1.1)}{0.0.0.0:0}
Job [ae918338be4a4e629e20f1fc128625da] finished successfully.
done: true
driverControlFilesUri: gs://dataproc-staging-us-east1-413280864426-jy7elnnn/google-cloud-dataproc-metainfo/1a55e1e0-b04a-4288-9cf6-6f862bde861f/jobs/ae918338be4a4e629e20
f1fc128625da/
driverOutputResourceUri: gs://dataproc-staging-us-east1-413280864426-jy7elnnn/google-cloud-dataproc-metainfo/1a55e1e0-b04a-4288-9cf6-6f862bde861f/jobs/ae918338be4a4e629e
20f1fc128625da/driveroutput
jobUuid: 56675615-ac17-3cf7-8aa5-731cf9fcd1fd
placement:
  clusterName: shakes-count
  clusterUuid: 1a55e1e0-b04a-4288-9cf6-6f862bde861f
pysparkJob:
  args:
  - gs://6893_course_data/shakes_text/
  - gs://6893_course_data/count_result/
  mainPythonFileUri: gs://dataproc-staging-us-east1-413280864426-jy7elnnn/google-cloud-dataproc-metainfo/1a55e1e0-b04a-4288-9cf6-6f862bde861f/jobs/ae918338be4a4e629e20f1
fc128625da/staging/word_count.py
reference:
  jobId: ae918338be4a4e629e20f1fc128625da
  projectId: e6893-big-data-analytics
status:
```

(2) Setting the **PROCESS** variable to **True**, the program will first do preprocessing such as converting all capital letters to lowercase and filtering out stopwords in English. The top 10 after preprocessing becomes **'macb', 'haue', 'thou', 'enter', 'shall', 'macbeth', 'thee', 'vpon', 'macd', 'yet'**.