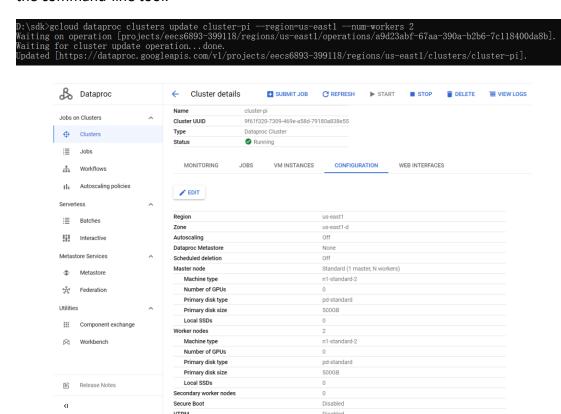# EECS6893 HW0

## 1. Warm-up exercises

### (1) Pi calculation

Create a computing cluster named "cluster-pi" and set the number of workers to 2.

The cluster is created through the API console, while the cluster is updated through the command-line tool.

```
D:\sdk>gcloud dataproc clusters update cluster-pi --region=us-east1 --num-workers 2
Waiting on operation [projects/eecs6893-399118/regions/us-east1/operations/a9d23abf-67aa-390a-b2b6-7c118400da8b].
Waiting for cluster update operation...done.
Updated [https://dataproc.googleapis.com/v1/projects/eecs6893-399118/regions/us-east1/clusters/cluster-pi].
```



n1-standrad-2 machines are used to create a cluster with 2 workers.

The java code of this program is:

```
1.  /*
2.   * Licensed to the Apache Software Foundation (ASF) under one or more
3.   * contributor license agreements.  See the NOTICE file distributed with
4.   * this work for additional information regarding copyright ownership.
5.   * The ASF licenses this file to You under the Apache License, Version 2.0
```

```java
6.    * (the "License"); you may not use this file except in compliance wi
      th
7.    * the License.  You may obtain a copy of the License at
8.    *
9.    *     http://www.apache.org/licenses/LICENSE-2.0
10.   *
11.   * Unless required by applicable law or agreed to in writing, softwar
      e
12.   * distributed under the License is distributed on an "AS IS" BASIS,
13.   * WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or im
      plied.
14.   * See the License for the specific language governing permissions an
      d
15.   * limitations under the License.
16.   */
17.
18. package org.apache.spark.examples;
19.
20. import org.apache.spark.api.java.JavaRDD;
21. import org.apache.spark.api.java.JavaSparkContext;
22. import org.apache.spark.sql.SparkSession;
23.
24. import java.util.ArrayList;
25. import java.util.List;
26.
27. /**
28.  * Computes an approximation to pi
29.  * Usage: JavaSparkPi [partitions]
30.  */
31. public final class JavaSparkPi {
32.
33.   public static void main(String[] args) throws Exception {
34.     SparkSession spark = SparkSession
35.       .builder()
36.       .appName("JavaSparkPi")
37.       .getOrCreate();
38.
39.     JavaSparkContext jsc = new JavaSparkContext(spark.sparkContext())
    ;
40.
41.     int slices = (args.length == 1) ? Integer.parseInt(args[0]) : 2;

42.     int n = 100000 * slices;
43.     List<Integer> l = new ArrayList<>(n);
```

```
44.      for (int i = 0; i < n; i++) {
45.        l.add(i);
46.      }
47.
48.      JavaRDD<Integer> dataSet = jsc.parallelize(l, slices);
49.
50.      int count = dataSet.map(integer -> {
51.        double x = Math.random() * 2 - 1;
52.        double y = Math.random() * 2 - 1;
53.        return (x * x + y * y <= 1) ? 1 : 0;
54.      }).reduce((integer, integer2) -> integer + integer2);
55.
56.      System.out.println("Pi is roughly " + 4.0 * count / n);
57.
58.      spark.stop();
59.    }
60. }
```

In this program, the RDD action is:

reduce(func): this action aggregates the elements of the RDD.

In this program, the 'reduce' action counts the number of the generated points located inside the circle.

The RDD transformation is:

map(func): this applys a function to every element of an RDD and return a new result RDD

In this program, the 'map' transformation is used to randomly generate points and decide whether the points are inside the circle.

Submit the job and runs the Spark program:

The program shows the result:

Pi is roughly 3.141894711418947

## (2) Word count



A single-node cluster is used in this part (and problem 3).

The Python code of 'word count' is:

```
1.  import pyspark
2.  import sys
```

```
3.
4.  if len(sys.argv) != 3:
5.    raise Exception("Exactly 2 arguments are required: <inputUri> <outp
      utUri>")
6.
7.  inputUri = sys.argv[1]
8.  outputUri = sys.argv[2]
9.
10. sc = pyspark.SparkContext()
11. lines = sc.textFile(sys.argv[1])
12. words = lines.flatMap(lambda line: line.split())
13. wordCounts = words.map(lambda word: (word, 1)).reduceByKey(lambda cou
    nt1, count2: count1 + count2)
14. wordCounts.saveAsTextFile(sys.argv[2])
```

In this program, the RDD transformations are:

1) map(func): Return a new distributed dataset formed by passing each element of the source through a function func.

    In this program, 'map' maps each word to a key-value pair, where the word is the key and the value is 1.

2) faltMap(func): Similar to map, but each input item can be mapped to 0 or more output items (so func should return a Seq rather than a single item).

    In this program, 'flatMap' split the lines to words.

3) reduceByKey(func, [numPartitions]): When called on a dataset of (K, V) pairs, returns a dataset of (K, V) pairs where the values for each key are aggregated using the given reduce function func, which must be of type (V,V) => V. Like in groupByKey, the number of reduce tasks is configurable through an optional second argument.

    In this program, 'reduceByKey' counts the number of occurrences of each word.

In this program, the RDD actions are:

1) saveAsTextFile(path): Write the elements of the dataset as a text file (or set of text files) in a given directory in the local filesystem, HDFS or any other Hadoop-

supported file system. Spark will call toString on each element to convert it to a line of text in the file.

In this program, 'saveAsTextFile' saves the words and the number of them in the output address entered in the command.

Submit the job:



Check the output:

## 2. NYC Bike expert

(1) Get the number of stations with longitude between -73.94 and -74.04

The following query is used:



The number of stations with longitude between -73.94 and -74.04 is 698.

(2) Total number of bikes available in region_id 71

The following query is used:



Total number of bikes available in region_id 71 is 11884.

(3) List all the station_id of the stations that have the largest capacity

The following query is used:



Stations 445, 422 and 501 have the max capacity of 79.

3. Understanding William Shakespeare

(1) Find top 10 frequent words without any text preprocessing.

The python code is shown as follows:

```python
import pyspark
import sys

if len(sys.argv) != 3:
    raise Exception("Exactly 2 arguments are required: <inputUri> <outputUri>")

inputUri = sys.argv[1]
outputUri = sys.argv[2]

sc = pyspark.SparkContext()
lines = sc.textFile(sys.argv[1])
words = lines.flatMap(lambda line: line.split())
wordCounts = words.map(lambda word: (word, 1)).reduceByKey(lambda count1, count2: count1 + count2)
reverse=wordCounts.map(lambda x: (x[1], x[0]))
rsort=reverse.sortByKey(ascending=False)
sort=rsort.map(lambda x: (x[1], x[0]))
top_ten=sort.take(10)
# Convert top_ten list to an RDD
top_ten_rdd = sc.parallelize(top_ten)

# Save the RDD as text files
top_ten_rdd.saveAsTextFile(sys.argv[2])

sc.stop()
```

Submit the job and store the output to the bucket.



Print the output:



(2) ) Find top 10 frequent words by first filtering out stop words provided byNLTK package. The Natural Language Toolkit, or more commonly NLTK, is a suite of libraries and programs to conduct natural language processing in Python.

The python code is shown as follows:

```python
import pyspark
import sys
import nltk
from nltk.corpus import stopwords

stopwords = set(stopwords.words('english'))


if len(sys.argv) != 3:
    raise Exception("Exactly 2 arguments are required: <inputUri> <outputUri>")

inputUri = sys.argv[1]
outputUri = sys.argv[2]


sc = pyspark.SparkContext()
lines = sc.textFile(sys.argv[1])
words = lines.flatMap(lambda line: line.split())

filtered = words.filter(lambda word: word not in stopwords)

wordCounts = filtered.map(lambda word: (word, 1)).reduceByKey(lambda count1, count2: count1 + count2)
reverse=wordCounts.map(lambda x: (x[1], x[0]))
rsort=reverse.sortByKey(ascending=False)
sort=rsort.map(lambda x: (x[1], x[0]))
top_ten=sort.take(10)
```

```python
# Convert top_ten list to an RDD
top_ten_rdd = sc.parallelize(top_ten)

# Save the RDD as text files
top_ten_rdd.saveAsTextFile(sys.argv[2])

sc.stop()
```

Submit the job and store the output to the bucket:

```
D:\6893\wordcount>gcloud dataproc jobs submit pyspark fstopw.py --cluster=cluster-wordcount --region=us-east1 -- gs://6893_bucket_1/shakes.txt gs://6893_bucket_1/output_filtered
/
Job [2879c3aa14f344fea5fc23dba6f7aeb5] submitted.
Waiting for job output...
23/09/21 19:58:11 INFO SparkEnv: Registering MapOutputTracker
23/09/21 19:58:11 INFO SparkEnv: Registering BlockManagerMaster
23/09/21 19:58:11 INFO SparkEnv: Registering BlockManagerMasterHeartbeat
23/09/21 19:58:11 INFO SparkEnv: Registering OutputCommitCoordinator
23/09/21 19:58:12 INFO DefaultNoHARMFailoverProxyProvider: Connecting to ResourceManager at cluster-wordcount-m.c.eecs6893-399118.internal./10.142.0.6:8032
23/09/21 19:58:12 INFO AHSProxy: Connecting to Application History server at cluster-wordcount-m.c.eecs6893-399118.internal./10.142.0.6:10200
23/09/21 19:58:13 INFO Configuration: resource-types.xml not found
23/09/21 19:58:13 INFO ResourceUtils: Unable to find 'resource-types.xml'.
23/09/21 19:58:13 INFO YarnClientImpl: Submitted application application_1695151409962_0022
23/09/21 19:58:14 INFO DefaultNoHARMFailoverProxyProvider: Connecting to ResourceManager at cluster-wordcount-m.c.eecs6893-399118.internal./10.142.0.6:8030
23/09/21 19:58:17 INFO GoogleCloudStorageImpl: Ignoring exception of type GoogleJsonResponseException; verified object already exists with desired state.
23/09/21 19:58:19 INFO FileInputFormat: Total input files to process : 1
23/09/21 19:58:33 INFO GoogleCloudStorageFileSystem: Successfully repaired 'gs://6893_bucket_1/output_filtered/' directory.
Job [2879c3aa14f344fea5fc23dba6f7aeb5] finished successfully.
done: true
driverControlFilesUri: gs://dataproc-staging-us-east1-656379034426-gcwzypcd/google-cloud-dataproc-metainfo/432fb230-c37e-49ec-b284-26f7a4c3f8c5/jobs/2879c3aa14f344fea5fc23dba6f7
aeb5/
driverOutputResourceUri: gs://dataproc-staging-us-east1-656379034426-gcwzypcd/google-cloud-dataproc-metainfo/432fb230-c37e-49ec-b284-26f7a4c3f8c5/jobs/2879c3aa14f344fea5fc23dba6
f7aeb5/driveroutput
jobUuid: 972a3434-0ee5-306b-8316-48d6684e90c8
placement:
  clusterName: cluster-wordcount
  clusterUuid: 432fb230-c37e-49ec-b284-26f7a4c3f8c5
pysparkJob:
  args:
  - gs://6893_bucket_1/shakes.txt
  - gs://6893_bucket_1/output_filtered/
  mainPythonFileUri: gs://dataproc-staging-us-east1-656379034426-gcwzypcd/google-cloud-dataproc-metainfo/432fb230-c37e-49ec-b284-26f7a4c3f8c5/jobs/2879c3aa14f344fea5fc23dba6f7ae
b5/staging/fstopw.py
reference:
  jobId: 2879c3aa14f344fea5fc23dba6f7aeb5
  projectId: eecs6893-399118
status:
  state: DONE
  stateStartTime: '2023-09-21T19:58:36.044630Z'
statusHistory:
- state: PENDING
  stateStartTime: '2023-09-21T19:58:05.478219Z'
- state: SETUP_DONE
  stateStartTime: '2023-09-21T19:58:05.514148Z'
- details: Agent reported job success
  state: RUNNING
  stateStartTime: '2023-09-21T19:58:05.711027Z'
yarnApplications:
- name: fstopw.py
  progress: 1.0
  state: FINISHED
  trackingUrl: http://cluster-wordcount-m.c.eecs6893-399118.internal.:8088/proxy/application_1695151409962_0022/
```

Print the output:

```
D:\6893\wordcount>gsutil cat gs://6893_bucket_1/output_filtered/*
('I', 326)
('And', 169)
('Macb.', 137)
('The', 131)
('haue', 114)
('That', 80)
('To', 79)
('Enter', 73)
('But', 61)
('thou', 61)
```