

1)

```
#!/usr/bin/env 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])
temp = lines.flatMap(Lambda line: [(line.split())[1]]).filter(Lambda word: word[0].isdigit())
words = temp.map(Lambda word: '-'.join([str((6*(int(word[0:2])/6))),(str(((6*(int(word[0:2])/6))+int(6))))))
wordCounts = words.map(Lambda word: (word, 1)).reduceByKey(Lambda count1, count2: count1 + count2)
wordCounts.coalesce(1).saveAsTextFile(sys.argv[2])
```

The above is the code used for the assignment. The python file for the same is attached with the submission zip file.

The output obtained by running the Dataproc job is:

```
('0-6', 501103)
('12-18', 499906)
('6-12', 499409)
('18-24', 499583)
```

Cluster name: lab4-cluster

Region: us-central1

Output path: gs://naveenrd_lab4/output/output.txt

The screenshot shows the Google Cloud Platform interface for a Dataproc job. The job ID is `bcdb1bb93aee4332afd5dabc17d0053`. The job status is "Succeeded". The start time is "Mar 6, 2021, 6:18:56 PM" and the elapsed time is "29 sec". The region is "us-central1" and the cluster is "lab4-cluster". The job type is "PySpark". The main python file is located at `gs://dataproc-staging-us-central1-147548493474-cporjya/google-cloud-dataproc-metainfo/1299398d-6fa6-47b4-a523-4fcaa1d2d171/jobs/bcdb1bb93aee4332afd5dabc17d0053/staging/lab4.py`.

The "Job output" section shows the following log entries:

```
21/03/06 12:48:59 INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat: Logging initialized
21/03/06 12:48:59 INFO org.spark_project.jetty.server.Server: jetty-9.3.z-SNAPSHOT, build timestamp: unknown, git hash: unknown
21/03/06 12:48:59 INFO org.spark_project.jetty.server.Server: Started @2697ms
21/03/06 12:48:59 INFO org.spark_project.jetty.server.AbstractConnector: Started ServerConnector@1b8d823b(HTTP/1.1,[http/1.1]){0.0.0.0:4040}
21/03/06 12:48:59 WARN org.apache.spark.scheduler.FairSchedulableBuilder: Fair Scheduler configuration file not found so jobs will be scheduled in FIFO order. To use fair scheduling, configure pools
21/03/06 12:49:01 INFO org.apache.hadoop.yarn.client.RMProxy: Connecting to ResourceManager at lab4-cluster-m/10.128.0.2:8032
21/03/06 12:49:01 INFO org.apache.hadoop.yarn.client.AHSProxy: Connecting to Application History server at lab4-cluster-m/10.128.0.2:10200
21/03/06 12:49:03 INFO org.apache.hadoop.yarn.client.api.impl.YarnClientImpl: Submitted application application_1615032111174_0006
21/03/06 12:49:10 INFO org.apache.hadoop.mapred.FileInputFormat: Total input files to process : 1
21/03/06 12:49:23 INFO org.spark_project.jetty.server.AbstractConnector: Stopped Spark@1b8d823b(HTTP/1.1,[http/1.1]){0.0.0.0:4040}
```

Job output is complete

2)

- a) **HDFS (Hadoop Distributed File System):** HDFS is used for storage permissions in Hadoop cluster. It is designed to work on a distributed file system design. HDFS is capable of handling larger size data with high volume velocity and variety makes Hadoop work more efficient and reliable with easy access to all its components. HDFS stores the data in the form of the block where the size of each data block is 128MB in size which is configurable means you can change it according to your requirement in *hdfs-site.xml* file in your Hadoop directory. HDFS provides high availability and fault tolerance. In HDFS data is stored in distributed format across various nodes. It provides scalability to scaleup or scale down nodes as per requirement. Using HDFS there is no fear of data loss because it provides replication. When a node fails the data is transferred across the nodes to restore.
- b) **Hive:** It is a data warehouse and ETL (Extraction, Transform & Load) tool which provides an SQL-like interface between the user and the HDFS. It facilitates reading, writing and handling wide datasets that are stored in distributed storage and queried by SQL syntax. Hive is built to manage and query only structured data which is residing under tables.
- c) **Pig:** Pig is a high-level platform or tool which is used to process the large datasets. It provides a high-level of abstraction for processing over the MapReduce. It provides a high-level scripting language, known as *Pig Latin* which is used to develop the data analysis codes. First, to process the data which is stored in the HDFS, the programmers will write the scripts using the Pig Latin Language. Internally *Pig Engine* (a component of Apache Pig) converted all these scripts into a specific map and reduce task. But these are not visible to the programmers in order to provide a high-level of abstraction. Pig Latin and Pig Engine are the two main components of the Apache Pig tool. The result of Pig always stored in the HDFS.
- d) **YARN (Yet Another Resource Negotiator):** YARN architecture separates resource management layer from the processing layer. YARN also allows different data processing engines like graph processing, interactive processing, stream processing as well as batch processing to run and process data stored in HDFS (Hadoop Distributed File System) thus making the system much more efficient. Through its various components, it can dynamically allocate various resources and schedule the application processing. For large volume data processing, it is quite necessary to manage the available resources properly so that every application can leverage them.