```python
1  # ML1020 - Assignment 2 (DAG for wordcount)
2
3  """"Example Airflow DAG that creates a Cloud Dataproc
   cluster, runs the Hadoop
4  wordcount example, and deletes the cluster.
5
6  This DAG relies on three Airflow variables
7  https://airflow.apache.org/docs/apache-airflow/stable/
   concepts/variables.html
8  * gcp_project - Google Cloud Project to use for the
   Cloud Dataproc cluster.
9  * gce_zone - Google Compute Engine zone where Cloud
   Dataproc cluster should be
10   created.
11 * gcs_bucket - Google Cloud Storage bucket to use for
   result of Hadoop job.
12   See https://cloud.google.com/storage/docs/creating-
   buckets for creating a
13   bucket.
14 """
15
16 import datetime
17 import os
18
19 from airflow import models
20 from airflow.providers.google.cloud.operators import
   dataproc
21 from airflow.utils import trigger_rule
22
23 # Output file for Cloud Dataproc job.
24 # If you are running Airflow in more than one time zone
25 # see https://airflow.apache.org/docs/apache-airflow/
   stable/timezone.html
26 # for best practices
27 output_file = os.path.join(
28     models.Variable.get('gcs_bucket'), 'wordcount',
29     datetime.datetime.now().strftime('%Y%m%d-%H%M%S'
   )) + os.sep
```

```python
30 # Path to Hadoop wordcount example available on every
   Dataproc cluster.
31 WORDCOUNT_JAR = (
32     'file:///usr/lib/hadoop-mapreduce/hadoop-mapreduce-
   examples.jar'
33 )
34 # Arguments to pass to Cloud Dataproc job.
35 #input_file = 'gs://pub/shakespeare/rose.txt'
36 input_file = 'gs://ml1020-bucket/asn1/*.txt'
37 wordcount_args = ['wordcount', input_file, output_file]
38
39 HADOOP_JOB = {
40     "reference": {"project_id": models.Variable.get('
   gcp_project')},
41     "placement": {"cluster_name": 'composer-hadoop-
   tutorial-cluster-{{ ds_nodash }}'},
42     "hadoop_job": {
43         "main_jar_file_uri": WORDCOUNT_JAR,
44         "args": wordcount_args,
45     },
46 }
47
48 CLUSTER_CONFIG = {
49     "master_config": {
50         "num_instances": 1,
51         "machine_type_uri": "n1-standard-2"
52     },
53     "worker_config": {
54         "num_instances": 2,
55         "machine_type_uri": "n1-standard-2"
56     },
57 }
58
59 yesterday = datetime.datetime.combine(
60     datetime.datetime.today() - datetime.timedelta(1),
61     datetime.datetime.min.time())
62
63 default_dag_args = {
```

```python
64        # Setting start date as yesterday starts the DAG
   immediately when it is
65        # detected in the Cloud Storage bucket.
66        'start_date': yesterday,
67        # To email on failure or retry set 'email' arg to
   your email and enable
68        # emailing here.
69        'email_on_failure': False,
70        'email_on_retry': False,
71        # If a task fails, retry it once after waiting at
   least 5 minutes
72        'retries': 1,
73        'retry_delay': datetime.timedelta(minutes=5),
74        'project_id': models.Variable.get('gcp_project'),
75        'location': models.Variable.get('gce_region'),
76
77 }
78
79
80 with models.DAG(
81          'wordcount_assignment',
82          # Continue to run DAG once per day
83          schedule_interval=datetime.timedelta(days=1),
84          default_args=default_dag_args) as dag:
85
86      # Create a Cloud Dataproc cluster.
87      create_dataproc_cluster = dataproc.
   DataprocCreateClusterOperator(
88          task_id='create_dataproc_cluster',
89          # Give the cluster a unique name by appending
   the date scheduled.
90          # See https://airflow.apache.org/docs/apache-
   airflow/stable/macros-ref.html
91          cluster_name='composer-hadoop-tutorial-cluster
   -{{ ds_nodash }}',
92          cluster_config=CLUSTER_CONFIG,
93          region=models.Variable.get('gce_region'))
94
```

```
95      # Run the Hadoop wordcount example installed on
    the Cloud Dataproc cluster
96      # master node.
97      run_wordcount = dataproc.DataprocSubmitJobOperator
    (
98          task_id='run_dataproc_hadoop',
99          job=HADOOP_JOB)
100
101     # Delete Cloud Dataproc cluster.
102     delete_dataproc_cluster = dataproc.
    DataprocDeleteClusterOperator(
103          task_id='delete_dataproc_cluster',
104          cluster_name='composer-hadoop-tutorial-cluster
    -{{ ds_nodash }}',
105          region=models.Variable.get('gce_region'),
106          # Setting trigger_rule to ALL_DONE causes the
    cluster to be deleted
107          # even if the Dataproc job fails.
108          trigger_rule=trigger_rule.TriggerRule.ALL_DONE
    )
109
110     # Define DAG dependencies.
111     create_dataproc_cluster >> run_wordcount >>
    delete_dataproc_cluster
```