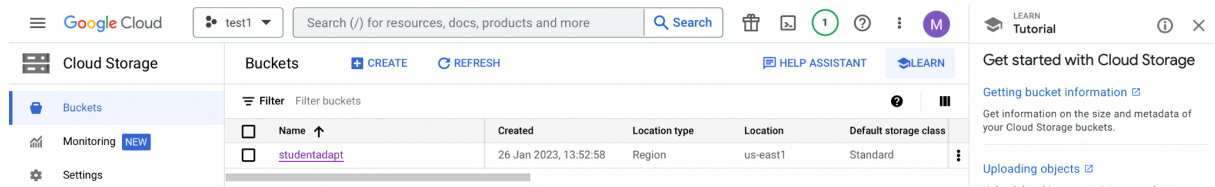# Data Engineering Final Assignment

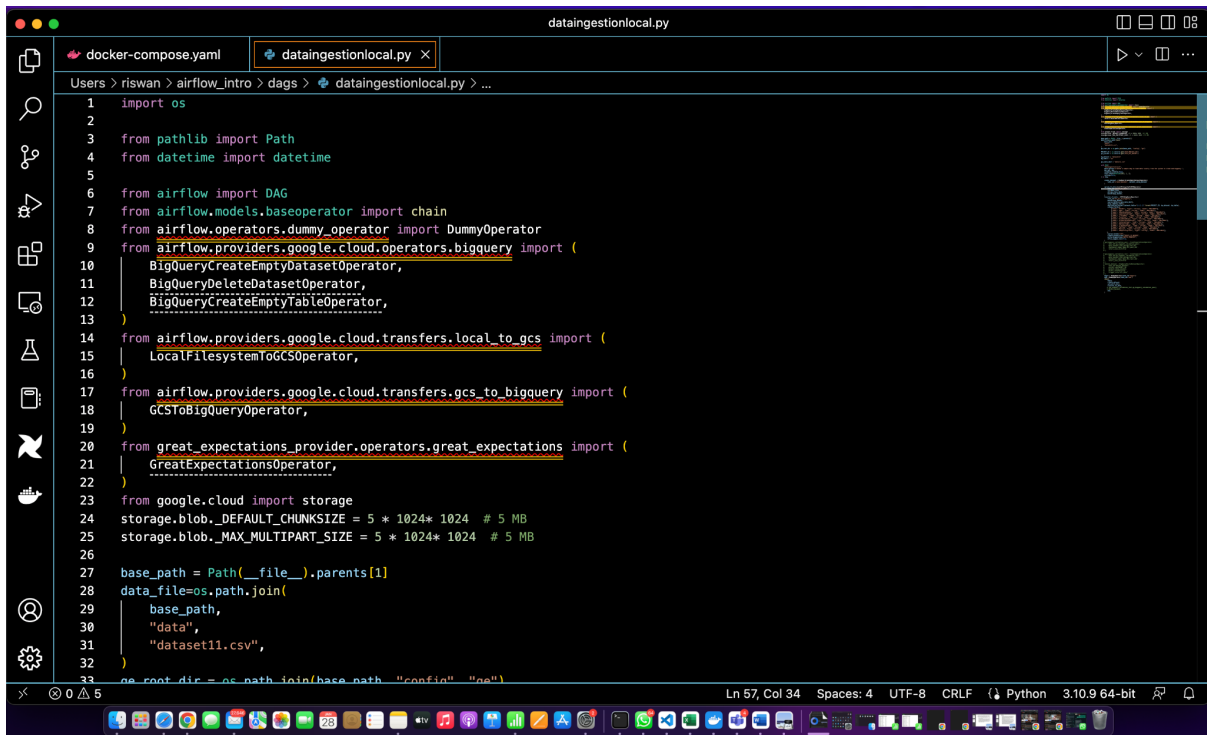1. Created a new project on Google Cloud Platform(GCP) project-id `test1-375701`
2. Created a bucket with the name "studentadapt" in cloud storage under the same project



3. Created a new dataset table in bigquery and named it "dataset11"
4. Created and downloaded service account key .json file and saved it locally in the airflow docker folder
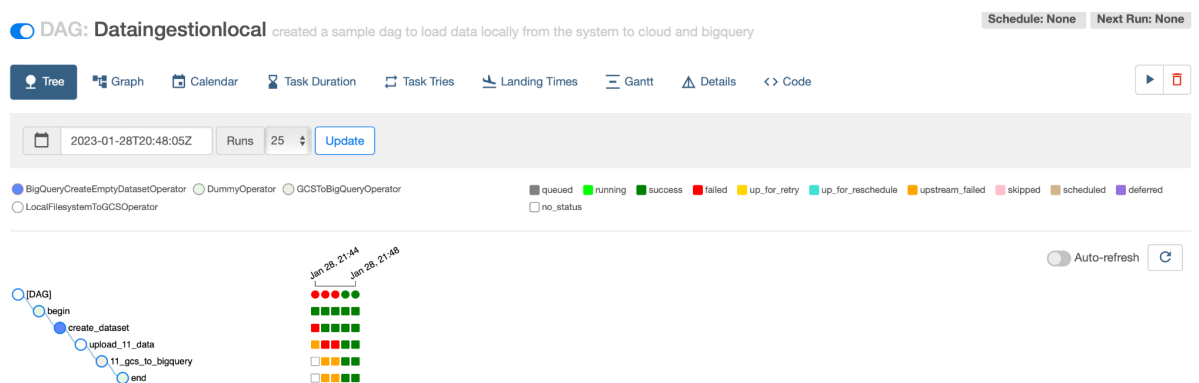5. Updated docker-compose.yaml file environment variables.

```
AIRFLOW__CORE__SQL_ALCHEMY_CONN: postgresql+psycopg2://airflow:airflow@postgres/airflow
AIRFLOW__CORE__FERNET_KEY: ''
AIRFLOW__CORE__DAGS_ARE_PAUSED_AT_CREATION: 'true'
AIRFLOW__CORE__LOAD_EXAMPLES: 'false'
AIRFLOW__API__AUTH_BACKEND: 'airflow.api.auth.backend.basic_auth'
_PIP_ADDITIONAL_REQUIREMENTS: ${_PIP_ADDITIONAL_REQUIREMENTS:-}
GOOGLE_APPLICATION_CREDENTIALS: /users/riswan/airflow_intro/.google/dtest1-375701-4df4e3d02696.json
AIRFLOW_CONN_GOOGLE_CLOUD_DEFAULT: 'google-cloud-platform://?extra__google_cloud_platform__key_path=/opt/airflow/.google/test1-3757
GCP_PROJECT_ID: 'test1-375701'
GCP_GCS_BUCKET: 'studentadapt'
GCP_BIGQUERY_DATASET: 'warehousede'
AIRFLOW_CONN_POSTGRES_DEFAULT: postgresql+psycopg2://airflow:airflow@postgres/postgres
```

6. Created a DAG python file under the dags folder named as `dataingestionlocal.py` (the dag file is present in airflow_intro>dags)

```
import os

from pathlib import Path
from datetime import datetime

from airflow import DAG
from airflow.models.baseoperator import chain
from airflow.operators.dummy_operator import DummyOperator
from airflow.providers.google.cloud.operators.bigquery import (
    BigQueryCreateEmptyDatasetOperator,
    BigQueryDeleteDatasetOperator,
    BigQueryCreateEmptyTableOperator,
)
from airflow.providers.google.cloud.transfers.local_to_gcs import (
    LocalFilesystemToGCSOperator,
)
from airflow.providers.google.cloud.transfers.gcs_to_bigquery import (
    GCSToBigQueryOperator,
)
from great_expectations_provider.operators.great_expectations import (
    GreatExpectationsOperator,
)
from google.cloud import storage
storage.blob._DEFAULT_CHUNKSIZE = 5 * 1024* 1024   # 5 MB
storage.blob._MAX_MULTIPART_SIZE = 5 * 1024* 1024   # 5 MB

base_path = Path(__file__).parents[1]
data_file=os.path.join(
    base_path,
    "data",
    "dataset11.csv",
)
ge_root_dir = os.path.join(base_path, "config", "ge")
```

7. Initialised docker-compose up yaml file (inside folder airflow_intro>)

8. Made sure that all the docker containers in airflow_intro folder is running

9. Opened and logged in airflow user interface

10.     Triggered the DAG dataingestionlocal

11. Result show as follows



12.     It shows that the process has begun, data set created in cloud storage, the data uploaded to cloud storage, the data transferred from cloud storage to bigquery.

### 13.   The data in google cloud bucket is shown as follows



### 14.   Data reflected in BigQuery:



### 15.   Running a sample query to make sure the data has reached the datawarehouse.

16. Created a blank report in Google Looker Studio and connected it with BigQuery using in-built function of google connectors.

17. The project name test1 was selected as shown in bigquery and the dataset was loaded.

18. Created 2 graphs

19. Graph1 shows quantitative educational levels of students who have adopted online education. Graph was returned as

below:



20. Graph2 shows categorical data of students' adaptability to online education based on the type of institution.