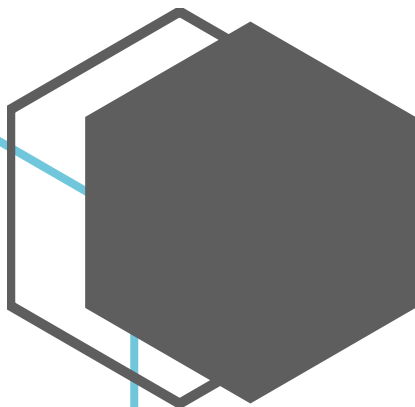# CSCI 5410

## Assignment 4 – Part B

**Name: Benny Daniel Tharigopala**
**Banner ID: B00899629**
**GitLab URL: https://git.cs.dal.ca/benny/csci5410_B00899629_Benny_Tharigopala**

# Event-driven serverless application using GCP ML

## Creating a project on Google Cloud Platform (GCP)



## Creating a service account on GCP



## A new service account is created

# Generating a key to establish a connection to GCP with JAVA SDK



# Creating a Cloud Function with a Cloud Storage bucket as a Trigger

# New Cloud Function on GCP

| | | Environment | Name ↑ | Region | Trigger | Runtime | Memory allocated | Executed function | Last deployed | Authentication ❓ | Actions |
|---|---|---|---|---|---|---|---|---|---|---|---|
| ☐ | ✓ | 1st gen | generateVector | us-central1 | HTTP | Python 3.9 | 256 MB | hello_world | Jul 9, 2022, 7:36:26 PM | | ⋮ |

Cloud Functions    Functions    ➕ CREATE FUNCTION    ↻ REFRESH

Filter   Filter functions

# Snapshot of Cloud Storage dashboard with no Buckets



# Driver Code for Creating Buckets

# Buckets for Source, Train and Test Files

| | Name ↑ | Created | Location type | Location | Default storage class | Last modified | Public access | Access control | Protection | Lifecycle |
|---|---|---|---|---|---|---|---|---|---|---|
| ☐ | sourcedatab00899629 | Jul 9, 2022, 6:55:05 PM | Multi-region | us | Standard | Jul 9, 2022, 6:55:05 PM | Subject to object ACLs | Fine-grained | None | None |
| ☐ | testdatab00899629 | Jul 9, 2022, 7:03:31 PM | Multi-region | us | Standard | Jul 9, 2022, 7:03:31 PM | Subject to object ACLs | Fine-grained | None | None |
| ☐ | traindatab00899629 | Jul 9, 2022, 7:03:30 PM | Multi-region | us | Standard | Jul 9, 2022, 7:03:30 PM | Subject to object ACLs | Fine-grained | None | None |

# Empty Bucket for Train data

**traindatab00899629**

| Location | Storage class | Public access | Protection |
|---|---|---|---|
| us (multiple regions in United States) | Standard | Subject to object ACLs | None |

OBJECTS    CONFIGURATION    PERMISSIONS    PROTECTION    LIFECYCLE

Buckets > traindatab00899629

UPLOAD FILES    UPLOAD FOLDER    CREATE FOLDER    MANAGE HOLDS    DOWNLOAD    DELETE

Filter by name prefix only ▾    Filter objects and folders    Show deleted data

| ☐ | Name | Size | Type | Created | Storage class | Last modified | Public access |
|---|---|---|---|---|---|---|---|

No rows to display

# Empty Bucket for Test data

**testdatab00899629**

| Location | Storage class | Public access | Protection |
|---|---|---|---|
| us (multiple regions in United States) | Standard | Subject to object ACLs | None |

OBJECTS    CONFIGURATION    PERMISSIONS    PROTECTION    LIFECYCLE

Buckets > testdatab00899629

UPLOAD FILES    UPLOAD FOLDER    CREATE FOLDER    MANAGE HOLDS    DOWNLOAD    DELETE

Filter by name prefix only ▾    Filter objects and folders    Show deleted data

| ☐ | Name | Size | Type | Created | Storage class | Last modified | Public access |
|---|---|---|---|---|---|---|---|

No rows to display

# Upload "Train" Files to the Source Bucket



# Source Bucket with Files

# Automatically generated CSV file



# Enabling API for Using Managed Notebooks



# Creating a new Notebook on Vertex AI

## Peek at the Dataframe



## Visualize the Plot for the Training Dataset

```
[9]: plt.plot(BeforeTrainOutPut, TrainOutput , label = "Actual Output",color="red")
     plt.plot( TrainOutput,BeforeTrainOutPut , label = "Predicted Output",color="green")
     plt.legend()
     plt.show()
```

## Visualize the cluster for levenshtein distance



## Upload Files to the "Test" Bucket

# Empty Bucket for Test data



# Automatically Generated CSV file for Test Data

## Empty Bucket for Test data



Test Data Model Prediction



Test dataset Clusters

**Code Blocks**

## Create Buckets (Driver Class) [1-2]

```java
import java.io.File;
import java.io.IOException;

public class CreateBuckets
{
    public static void main(String[] args)
    {
        String projectId = Credentials.PROJECT_ID;
        String sourceBucket = "sourcedatab00899629";
        String trainBucket = "traindatab00899629";
        String testBucket = "testdatab00899629";
        File parentDirectory = new File("./src/data/Train");
        File[] trainDataset = parentDirectory.listFiles();

        File testParentDirectory = new File("./src/data/Test");
```

```
        File[] testDataset = testParentDirectory.listFiles();

        try {

CreateBucketWithStorageClassAndLocation.createBucketWithStorageClassAndLocation(projectId
,sourceBucket,trainDataset);
CreateBucketWithStorageClassAndLocation.createBucketWithStorageClassAndLocation(projectId
,trainBucket);
CreateBucketWithStorageClassAndLocation.createBucketWithStorageClassAndLocation(projectId
,sourceBucket);
        } catch (IOException e) {
            System.out.println(e);
        }
    }
}
```

## Create Buckets and Upload Files (Driver Class) [1-2]

```
import com.google.auth.Credentials;
import com.google.auth.oauth2.GoogleCredentials;
import com.google.cloud.storage.BlobId;
import com.google.cloud.storage.BlobInfo;
import com.google.cloud.storage.Storage;
import com.google.cloud.storage.StorageOptions;
import java.io.File;
import java.io.FileInputStream;
import java.io.IOException;
import java.nio.file.Files;
import com.google.cloud.storage.BucketInfo;
import com.google.cloud.storage.Bucket;


public class CreateBucketWithStorageClassAndLocation {
    public static void createBucketWithStorageClassAndLocation(String projectId, String
bucketName,  File[] dataset) throws IOException
    {
        Credentials credentials = GoogleCredentials.fromStream(new
FileInputStream("src/csci5410-assignment-4-0356bf0945f3.json"));
        Storage storage =
StorageOptions.newBuilder().setCredentials(credentials).setProjectId(projectId).build().g
etService();
        Bucket bucket = storage.create(BucketInfo.newBuilder(bucketName).build());
        System.out.println("Created bucket " + bucket.getName());
int count=0;
        for(File file : dataset)
        {
            if(count<10) {
                BlobId blobId = BlobId.of(bucketName, file.getName());
                BlobInfo blobInfo = BlobInfo.newBuilder(blobId).build();
                storage.create(blobInfo, Files.readAllBytes(file.toPath()));
                System.out.println("File - " + file.getName() + " has been uploaded to
bucket - " + bucketName);
                count++;
            }
            else {
                break;
            }
        }
```

```
    }

    public static void createBucketWithStorageClassAndLocation(String projectId, String
bucketName) throws IOException
    {
        Credentials credentials = GoogleCredentials.fromStream(new
FileInputStream("src/csci5410-assignment-4-0356bf0945f3.json"));
        Storage storage =
StorageOptions.newBuilder().setCredentials(credentials).setProjectId(projectId).build().g
etService();
        Bucket bucket = storage.create(BucketInfo.newBuilder(bucketName).build());
        System.out.println("Created bucket " + bucket.getName());
    }
}
```

# Cloud Function to extract words and Levenshtein distance [3-5]

```
    file = event
    bucket = storage_client.get_bucket(file['bucket'])
    uploadbucket = storage_client.get_bucket('traindatab00899629');
    print('traindatab00899629')
    blob = bucket.blob(file['name'])
    contents = blob.download_as_string()

    decodedstring = contents.decode(encoding="utf-8", errors="ignore")
    decodedstring = decodedstring.replace("\n", " ")
    words = decodedstring.split(" ")
    length = len(words)
    uploadCsv = 'trainVector.csv'
    StoreInCsv = ""
    blob = uploadbucket.blob(uploadCsv)
```

```
def distance(s1,s2):
    if len(s1) > len(s2):
        s1,s2 = s2,s1
    l1 = len(s1) + 1
    l2 = len(s2) + 1
    dp = {}
    for i in range(l1):
        dp[i,0] = i
    for j in range(l2):
        dp[0,j] = j
    for i in range(1,l1):
        for j in range(1,l2):
            cost = 0 if s1[i - 1] == s2[j - 1] else 1
            dp[i,j] = min(dp[i,j - 1] + 1, dp[i - 1,j] + 1,dp[i - 1,j - 1] + cost)
    return dp[l1 - 1,l2 - 1]
```

```
for i in range(length - 1):
        s1 = re.sub('[^A-Za-z]+', '', words[i])
        s2 = re.sub('[^A-Za-z]+', '', words[i  + 1])
        if (s1.lower() not in stopWords) and (s2.lower() not in stopWords):
            distance = distance(s1.lower(), s2.lower())
```

```
        StoreInCsv += s1 + "," + s2 + "," + str(distance) + '\n'
    print(StoreInCsv)
```
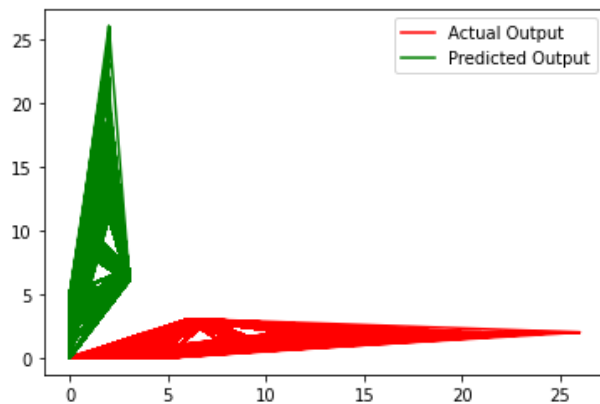
## K-Means Cluster and Plots [6-7]

```python
client_storage = storage.Client()
bucketNameTrain = "traindatab00899629"
trainDataSet = "trainVector.csv"
TrainBucket = client_storage.get_bucket(bucketNameTrain)
TrainBlob = TrainBucket.blob(trainDataSet)
TrainBlobData = TrainBlob.download_as_string()
decodedstring = TrainBlobData.decode(encoding="utf-8", errors="ignore")
dfTrain = pd.read_csv(StringIO(decodedstring.lower()),sep=",")
BeforeTrainOutPut = dfTrain.iloc[1:,-1:]
print(BeforeTrainOutPut)
```

```python
kmeans = KMeans(n_clusters=4, init='k-means++', max_iter=300)
TrainOutput = kmeans.fit_predict(BeforeTrainOutPut)
print(TrainOutput)

plt.plot(BeforeTrainOutPut, TrainOutput , label = "Actual Output")
plt.plot( TrainOutput,BeforeTrainOutPut , label = "Predicted Output")
plt.legend()
plt.show()
```

```python
plt.scatter(BeforeTrainOutPut[TrainOutput == 0],BeforeTrainOutPut[TrainOutput == 0],color =
'red')
plt.scatter(BeforeTrainOutPut[TrainOutput == 2],BeforeTrainOutPut[TrainOutput == 2],color =
'blue')
plt.scatter(BeforeTrainOutPut[TrainOutput == 4],BeforeTrainOutPut[TrainOutput == 4],color =
'green')
plt.scatter(BeforeTrainOutPut[TrainOutput == 6],BeforeTrainOutPut[TrainOutput == 6],color =
'yellow')
plt.show();
```

```python
bucketNameTest = "testdatab00899629"
testDataSet = "testVector.csv"
TestBucket = client_storage.get_bucket(bucketNameTest)
TestBlob = TestBucket.blob(testDataSet)
TestBlobData = TestBlob.download_as_string()
decodedstring = TestBlobData.decode(encoding="utf-8", errors="ignore")
dfTest = pd.read_csv(StringIO(decodedstring.lower()),sep=",")
BeforeTestOutPut = dfTest.iloc[1:,-1:]

TestOutput = kmeans.fit_predict(BeforeTestOutPut)
```

```
plt.plot(BeforeTestOutPut, TestOutput , label = "Actual Output")
plt.plot( TestOutput,BeforeTestOutPut , label = "Predicted Output")
plt.show()
```

```
cluster_centres = kmeans.cluster_centers_

plt.scatter(BeforeTestOutPut[TestOutput == 0],BeforeTestOutPut[TestOutput == 0],color =
'red')
plt.scatter(BeforeTestOutPut[TestOutput == 2],BeforeTestOutPut[TestOutput == 2],color =
'blue')
plt.scatter(BeforeTestOutPut[TestOutput == 4],BeforeTestOutPut[TestOutput == 4],color =
'green')
plt.scatter(BeforeTestOutPut[TestOutput == 6],BeforeTestOutPut[TestOutput == 6],color =
'yellow')

plt.title("Test dataset Clusters")
plt.xlabel("Levenshtein distance")
plt.ylabel("Levenshtein distance")
plt.legend()
plt.show();
```

## Citations

[1] "Cloud Storage client libraries," *Google Cloud*, May 12, 2017. https://cloud.google.com/storage/docs/reference/libraries (accessed Jul. 06, 2022).

[2] "Cloud Storage," *Google Cloud*, May 26, 2019. https://cloud.google.com/storage/ (accessed Aug. 13, 2016).

[3] "Cloud Functions," *Google Cloud*, Sep. 11, 2019. https://cloud.google.com/functions (accessed Jul. 06, 2022).

[4] "Google Cloud Storage Triggers | Cloud Functions Documentation," *Google Cloud*, Dec. 14, 2016. https://cloud.google.com/functions/docs/calling/storage (accessed Jul. 06, 2022).

[5] Vatsal, "Text Similarity w/ Levenshtein Distance in Python," *Medium*, May 15, 2015. https://towardsdatascience.com/text-similarity-w-levenshtein-distance-in-python-2f7478986e75 (accessed Jul. 07, 2022).

[6] R. Python, "K-Means Clustering in Python: A Practical Guide – Real Python," *realpython.com*, Dec. 12, 2012. https://realpython.com/k-means-clustering-python/ (accessed Nov. 12, 2016).

[7] "Create a user-managed notebooks instance | Vertex AI Workbench," *Google Cloud*, May 16, 2016. https://cloud.google.com/vertex-ai/docs/workbench/user-managed/create-new (accessed Jul. 04, 2022).