# Introduction to Big Data
# Graded Assignment

Name - Avijeet Palit
Roll - 21f1005675
Date - 22/11/2024

## Objective

Convert the Spark MLib code from the Databricks decision trees example to use the CrossValidator autotuner. Analyze and report the best-performing model parameters.

## Detailed Steps

### Step 1: Create a GCS Bucket

- Navigate to the Cloud Storage Section:
- Log in to your Google Cloud Platform (GCP) Console.
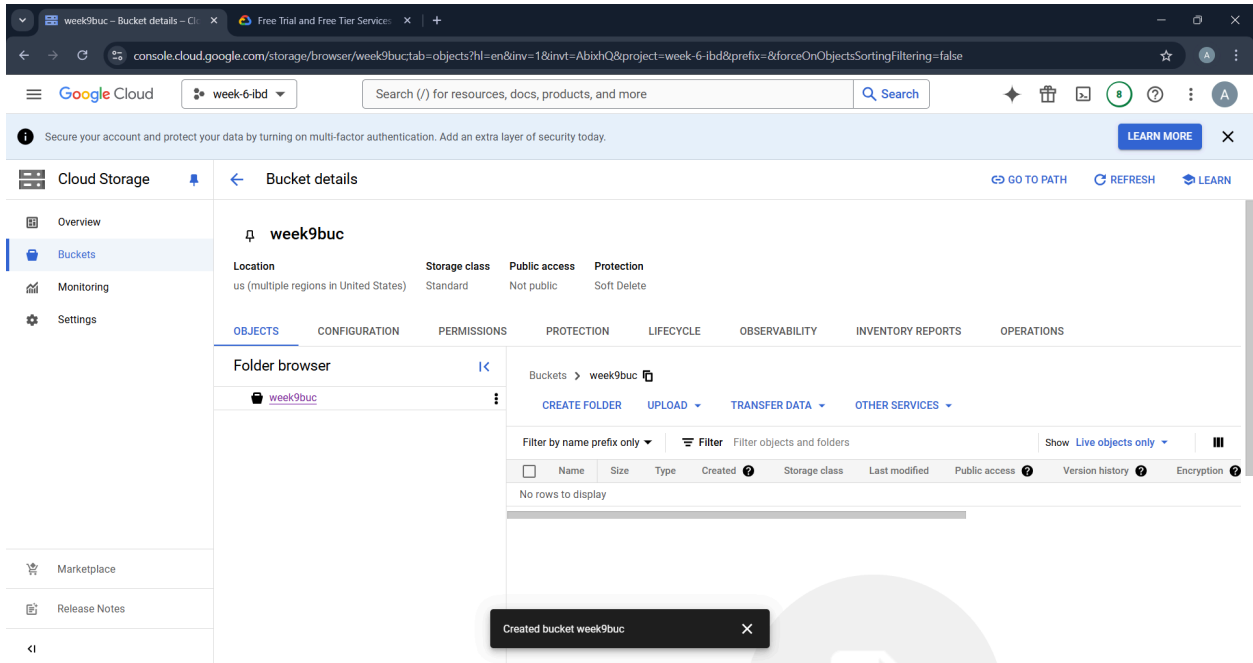- Click on Cloud Storage from the navigation menu.
- Create a New Bucket:



- Provide a name for the bucket (e.g., "week9buc").
- Set Location and Permissions:
- Select the desired bucket location (e.g., Regional/Multiregional).

## Step 2: Generate or Use MNIST Data

- Obtain MNIST Dataset:
- Either generate the MNIST Train and Test datasets or download them from a reliable source.

```python
import numpy as np
import pandas as pd
from tensorflow.keras.datasets import mnist

# Step 1: Load the MNIST dataset
(X_train, y_train), (X_test, y_test) = mnist.load_data()

# Step 2: Flatten the 28x28 images into a single 784-length vector
X_train_flat = X_train.reshape(X_train.shape[0], -1)
X_test_flat = X_test.reshape(X_test.shape[0], -1)

# Step 3: Normalize the data to [0, 1]
X_train_flat = X_train_flat.astype('float32') / 255.0
X_test_flat = X_test_flat.astype('float32') / 255.0

# Step 4: Create Pandas DataFrames
# Combine the features and labels for both train and test datasets
train_data = pd.DataFrame(X_train_flat)
train_data['label'] = y_train

test_data = pd.DataFrame(X_test_flat)
test_data['label'] = y_test

# Step 5: Save the datasets to CSV files
train_data.to_csv('mnist_train.csv', index=False)
test_data.to_csv('mnist_test.csv', index=False)

print("MNIST Train and Test datasets saved as 'mnist_train.csv' and 'mnist_test.csv'")
```

- Upload Data to the GCS Bucket:

**Step 3: Create a Dataproc Cluster**
- Go to the Dataproc Section:
- In the GCP Console, navigate to the Dataproc page.
- Create a Cluster:
- Click on Create Cluster.
- Configure Cluster Settings:
  - Set Cluster Type to "Standard".
  - Choose appropriate machine types for the master and worker nodes based on your data size and computational needs.
  - Enable the Component Gateway for easier job monitoring and debugging.
  - Ensure Spark Compatibility

**Step 4: Write and Modify the Spark Code**
- Fetch Base Code:
  - Download the sample decision trees code from the Databricks example provided in the assignment.

- Modify the Code to Use CrossValidator:
  - Use CrossValidator to tune the model parameters (e.g., maxDepth, maxBins).
  - Update Model Evaluation Metrics:
  - Adjust the code to output performance metrics for each model.



```
In [5]: pip install pyspark

        Requirement already satisfied: pyspark in /usr/lib/spark/python (3.5.1)
        Requirement already satisfied: py4j==0.10.9.7 in /opt/conda/miniconda3/lib/python3.11/site-packages (from pyspark) (0.10.9.7)
        WARNING: Running pip as the 'root' user can result in broken permissions and conflicting behaviour with the system package mana
        ger. It is recommended to use a virtual environment instead: https://pip.pypa.io/warnings/venv
        Note: you may need to restart the kernel to use updated packages.

In [6]: import numpy as np
        from pyspark.sql import SparkSession
        from pyspark import SparkContext
        from pyspark.sql.types import *
        from pyspark.ml.feature import VectorAssembler
        from pyspark.ml.tuning import CrossValidator, ParamGridBuilder
        from pyspark.ml.classification import DecisionTreeClassifier
        from pyspark.ml.evaluation import MulticlassClassificationEvaluator
        import matplotlib.pyplot as plt

In [7]: fields = [StructField('label', IntegerType(), False)]
        for i in range(1,29):
            for j in range(1,29):
                fields.append(StructField(str(i)+'x'+str(j),IntegerType(),True))
        Schema = StructType(fields)

In [9]: test_set = spark.read.load("gs://week9buc/mnist_test.csv", format="csv", header=True, schema = Schema)
        train_set = spark.read.load("gs://week9buc/mnist_train.csv", format="csv", header=True, schema = Schema)

In [10]: features = train_set.columns[1:]
         assembler = VectorAssembler(inputCols=features, outputCol="features")
```

jupyter **week9** (unsaved changes)

File   Edit   View   Insert   Cell   Kernel   Widgets   Help

```
cv = CrossValidator(estimator = model, estimatorParamMaps = grid,
                    evaluator = MulticlassClassificationEvaluator(labelCol="label", predictionCol="prediction"),
                    numFolds = 5)
```

In [*]: `cvModel = cv.fit(train_set)`

```
24/11/29 13:23:23 WARN DAGScheduler: Broadcasting large task binary with size 1219.1 KiB
24/11/29 13:23:26 WARN DAGScheduler: Broadcasting large task binary with size 1108.4 KiB
24/11/29 13:23:43 WARN DAGScheduler: Broadcasting large task binary with size 1149.0 KiB
24/11/29 13:23:47 WARN DAGScheduler: Broadcasting large task binary with size 1039.0 KiB
24/11/29 13:24:10 WARN DAGScheduler: Broadcasting large task binary with size 1170.8 KiB
24/11/29 13:24:15 WARN DAGScheduler: Broadcasting large task binary with size 1036.1 KiB
24/11/29 13:24:41 WARN DAGScheduler: Broadcasting large task binary with size 1009.4 KiB
24/11/29 13:24:48 WARN DAGScheduler: Broadcasting large task binary with size 1153.9 KiB
24/11/29 13:24:56 WARN DAGScheduler: Broadcasting large task binary with size 1292.7 KiB
24/11/29 13:24:59 WARN DAGScheduler: Broadcasting large task binary with size 1338.4 KiB
24/11/29 13:25:02 WARN DAGScheduler: Broadcasting large task binary with size 1019.7 KiB
24/11/29 13:25:15 WARN DAGScheduler: Broadcasting large task binary with size 1219.1 KiB
24/11/29 13:25:17 WARN DAGScheduler: Broadcasting large task binary with size 1467.0 KiB
24/11/29 13:25:18 WARN DAGScheduler: Broadcasting large task binary with size 1684.6 KiB
24/11/29 13:25:20 WARN DAGScheduler: Broadcasting large task binary with size 1853.7 KiB
24/11/29 13:25:21 WARN DAGScheduler: Broadcasting large task binary with size 1966.6 KiB
24/11/29 13:25:23 WARN DAGScheduler: Broadcasting large task binary with size 1474.5 KiB
24/11/29 13:25:39 WARN DAGScheduler: Broadcasting large task binary with size 1149.0 KiB
24/11/29 13:25:42 WARN DAGScheduler: Broadcasting large task binary with size 1358.6 KiB
24/11/29 13:25:44 WARN DAGScheduler: Broadcasting large task binary with size 1545.5 KiB
24/11/29 13:25:46 WARN DAGScheduler: Broadcasting large task binary with size 1693.9 KiB
24/11/29 13:25:48 WARN DAGScheduler: Broadcasting large task binary with size 1807.5 KiB
24/11/29 13:25:50 WARN DAGScheduler: Broadcasting large task binary with size 1383.3 KiB
24/11/29 13:26:11 WARN DAGScheduler: Broadcasting large task binary with size 1170.8 KiB
```

**Results**
- Best Parameters after CrossValidator Autotuning:
- List the parameters that yielded the best performance (e.g., maxDepth = 10, maxBins = 32).

In [9]:
```
est = cvModel.bestModel
print(f'Best tree depth {est.getMaxDepth()}, Best bin size {est.getMaxBins()}')
```

Best tree depth 16, Best bin size 2

In [10]:
```
testresult = est.transform(test_set)
```

In [11]:
```
print(f'Total test points {testresult.count()}')
wrong = testresult.where(testresult.label!=testresult.prediction)
print(f'Wrongly classified count {wrong.count()}')
```

Total test points 9999

```
24/11/29 14:15:48 WARN DAGScheduler: Broadcasting large task binary with size 1606.9 KiB
[Stage 544:============================>                           (1 + 1) / 2]
```

Wrongly classified count 1124

- Misclassified Data Points:
- Analyze the misclassified data points to identify potential patterns or weaknesses in the model.

```
In [12]: temp = wrong.limit(5).select("label","prediction","features").collect()

         for r in temp:
             plt.figure()
             row = np.array(r.features).reshape((28,28))
             plt.imshow(row,cmap = 'gray')
             plt.title(f'True {r.label} Predicted {r.prediction}')
```



True 4 Predicted 8.0