


Microsoft Azure

 databricks

Search data, notebooks, recents, and more...

CTRL + P

DATABRICKS\_WORKSPACE

5

New

Workspace

Recents

Catalog

Workflows

Compute

Data Engineering

Job Runs

Machine Learning

Playground

Experiments

Features

Models

Serving

Partner Connect

Workspace

sanjaymass6677@gmail.com

input\_data.csv

Untitled Notebook 2024-12-07 19:30:33

Untitled Notebook 2024-12-07 22:20:46

Untitled Notebook 2024-12-07 22:20:46

Python

File Edit View Run Help Last edit was 5 minutes ago

Run all

cluster001

Schedule

Share

11:28 PM (1s)

1

Python

# Databricks notebook containing ETL, ML training, and deployment steps.

# Step 1: Initialize Spark Session

from pyspark.sql import SparkSession

from pyspark.ml.feature import VectorAssembler

from sklearn.linear\_model import LinearRegression

from sklearn.model\_selection import train\_test\_split

from sklearn.metrics import mean\_squared\_error

import joblib

import numpy as np

import pandas as pd

from azure.storage.blob import BlobServiceClient

import io

import os

spark = SparkSession.builder \

.appName("ETL ML Training and Deployment") \

.getOrCreate()

# Step 2: ETL Process (Extract, Transform, Load)

# Example: Extracting data from an external CSV file or direct data input

new\_data = [

(1.5, 2.3, 3.1, 4.0, 5.2),

(1.4, 2.2, 3.0, 4.1, 5.1),

(1.7, 2.5, 3.3, 4.3, 5.4),

(1.8, 2.6, 3.4, 4.4, 5.5),

]

# Define the column names for the data

columns = ['feature1', 'feature2', 'feature3', 'feature4', 'feature5']

# Step 3: Load data into a PySpark DataFrame

new\_data\_df = spark.createDataFrame(new\_data, columns)

# Step 4: Feature Engineering (Transform Data)

# Using VectorAssembler to combine features into a single vector column

assembler = VectorAssembler(inputCols=columns, outputCol="features")

assembled\_data = assembler.transform(new\_data\_df)

# Step 5: Machine Learning - Train a Model (Linear Regression)

# For simplicity, we will use random data for training as an example

X = np.random.rand(100, 5) # Random features

y = np.random.rand(100) # Random target values

# Split the data for training and testing

X\_train, X\_test, y\_train, y\_test = train\_test\_split(X, y, test\_size=0.2, random\_state=42)

# Train a linear regression model

model = LinearRegression()

model.fit(X\_train, y\_train)

# Evaluate the model

y\_pred = model.predict(X\_test)

mse = mean\_squared\_error(y\_test, y\_pred)

print(f"Mean Squared Error (MSE): {mse}")

# Step 6: Save the trained model to Azure Blob Storage

# Azure Blob Storage details

connection\_string = "DefaultEndpointsProtocol=https;AccountName=dbstoragefordb01;

AccountKey=Vhki9bLFbWinGmKsc+0S/jVrJ+1ldHwrgWIyr7PbmMUwa04XwljMTLXsaerhPy4MvtzVXFhraQk+AST6oOH+Q==;

EndpointSuffix=core.windows.net"

container\_name = "container01fordb"

model\_blob\_name = "linear\_regression\_model.pkl" # Model filename in the container

# Initialize BlobServiceClient

blob\_service\_client = BlobServiceClient.from\_connection\_string(connection\_string)

# Get BlobClient to upload the model file

blob\_client = blob\_service\_client.get\_blob\_client(container=container\_name, blob=model\_blob\_name)

# Save the model as a .pkl file

with open("/dbfs/tmp/linear\_regression\_model.pkl", "wb") as model\_file:

joblib.dump(model, model\_file)

# Upload the saved model file to Azure Blob Storage

with open("/dbfs/tmp/linear\_regression\_model.pkl", "rb") as model\_file:

blob\_client.upload\_blob(model\_file, overwrite=True)

Generate (Ctrl + I)

```

print(f"Model uploaded successfully to Azure Blob Storage as {model_blob_name}")

# Step 7: Predicting with the trained model on new data

# Generate predictions (in a real-world scenario, this would use new data)
def generate_predictions(model, assembled_data):
    # Convert assembled data to pandas and get features
    pandas_df = assembled_data.select('features').toPandas()
    features = np.array(pandas_df['features'].tolist())

    # Use the trained model to predict
    predictions = model.predict(features)

    # Convert predictions back to a Spark DataFrame
    prediction_df = spark.createDataFrame(pd.DataFrame({'prediction': predictions}))
    return prediction_df

# Generate predictions using the trained model
prediction_df = generate_predictions(model, assembled_data)

# Show predictions
prediction_df.show()

# Step 8: Save the predictions to Azure Blob Storage as a CSV file

# Convert the Spark DataFrame to Pandas to write as CSV
prediction_pandas_df = prediction_df.toPandas()

# Save the predictions to a CSV file in Azure Blob Storage
predictions_blob_name = "predictions.csv"
predictions_blob_client = blob_service_client.get_blob_client(container=container_name,
    blob=predictions_blob_name)

# Write the CSV to Blob Storage
with io.BytesIO() as output:
    prediction_pandas_df.to_csv(output, index=False)
    output.seek(0)
    predictions_blob_client.upload_blob(output, overwrite=True)

print(f"Predictions saved to Azure Blob Storage as {predictions_blob_name}")

```

#### ► (2) Spark Jobs

- assembled\_data: pyspark.sql.dataframe.DataFrame
- new\_data\_df: pyspark.sql.dataframe.DataFrame = [feature1: double, feature2: double ... 3 more fields]
- prediction\_df: pyspark.sql.dataframe.DataFrame = [prediction: double]

Mean Squared Error (MSE): 0.10266535721410726

Model uploaded successfully to Azure Blob Storage as linear\_regression\_model.pkl

```

+-----+
|      prediction|
+-----+
|1.1930116386969973|
| 1.190164935977604|
|1.2246805286998022|
|1.2379174252448029|
+-----+

```

Predictions saved to Azure Blob Storage as predictions.csv



2

[Shift+Enter] to run and move to next cell  
 [Ctrl+Shift+P] to open the command palette  
 [Esc H] to see all keyboard shortcuts