# SPARK ML

1.Terraform Apply:



```
MINGW64:/c/Users/YuriiHordiichuk/Desktop/m08_sparkml_python_azure-mas...   —

azurerm_databricks_workspace.bdcc: Still creating... [01m50s elapsed]
azurerm_databricks_workspace.bdcc: Still creating... [02m00s elapsed]
azurerm_databricks_workspace.bdcc: Still creating... [02m10s elapsed]
azurerm_databricks_workspace.bdcc: Still creating... [02m20s elapsed]
azurerm_databricks_workspace.bdcc: Still creating... [02m30s elapsed]
azurerm_databricks_workspace.bdcc: Still creating... [02m40s elapsed]
azurerm_databricks_workspace.bdcc: Still creating... [02m50s elapsed]
azurerm_databricks_workspace.bdcc: Still creating... [03m00s elapsed]
azurerm_databricks_workspace.bdcc: Still creating... [03m10s elapsed]
azurerm_databricks_workspace.bdcc: Still creating... [03m20s elapsed]
azurerm_databricks_workspace.bdcc: Creation complete after 3m29s [id=/s
ons/316ab800-22b5-40ab-be41-8595de4fb2d4/resourceGroups/rg-dev-westeuro
roviders/Microsoft.Databricks/workspaces/dbw-dev-westeurope-j2ta]
Releasing state lock. This may take a few moments...

Apply complete! Resources: 5 added, 0 changed, 0 destroyed.

Outputs:

resource_group_name = "rg-dev-westeurope-j2ta"

AzureAD+YuriiHordiichuk@EPPLWROW0218 MINGW64 ~/Desktop/m08_sparkml_pyth
master/terraform
$
```

2.Baseline model result:



```
▶          ✓     12:53 PM (2s)

    model = mlflow.pyfunc.load_model(f"models:/{model_name}@production")

    # Sanity-check: This should match the AUC logged by MLflow
    print(f'AUC: {roc_auc_score(y_test, model.predict(X_test))}')
  ▶ (4) Spark Jobs

Downloading artifacts: 100% |████████████████████████| 9/9 [00:00<00:00, 21.4

AUC: 0.8540300975814177
```

Features importance:

```
feature_importances = pd.DataFrame(model.feature_importances_, index=X_train.columns.tolist(), columns=['importance'])
feature_importances.sort_values('importance', ascending=False)
```
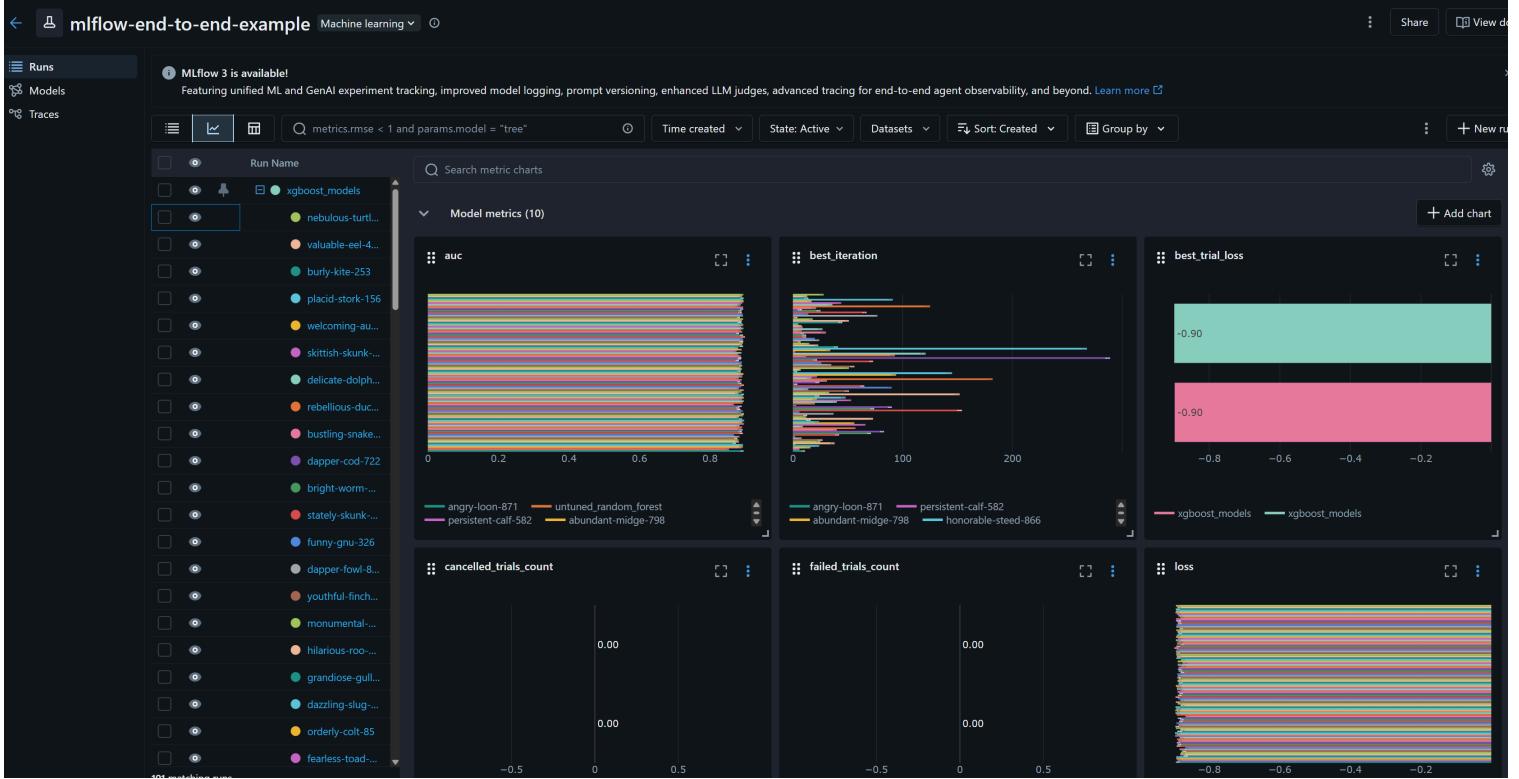
|  | importance |
|---|---|
| alcohol | 0.160192 |
| density | 0.117415 |
| volatile_acidity | 0.093136 |
| chlorides | 0.086618 |
| residual_sugar | 0.082544 |
| free_sulfur_dioxide | 0.080473 |
| pH | 0.080212 |
| total_sulfur_dioxide | 0.077798 |
| sulphates | 0.075780 |
| citric_acid | 0.071857 |
| fixed_acidity | 0.071841 |
| is_red | 0.002134 |

Model available in models tab:



3. New model xgboost_models experiments:

Version 2 with better results added:



4. Server thee model. I rewrote this cell to successfully log the model:

```
import mlflow
from mlflow.models.signature import infer_signature

model_name = "wine_quality"
X_sample = X_train.iloc[:10]

# Load the model to infer signature
model = mlflow.pyfunc.load_model(f"models:/{model_name}@production")
y_pred = model.predict(X_sample)
signature = infer_signature(X_sample, y_pred)

with mlflow.start_run():
    mlflow.sklearn.log_model(
        sk_model=model._model_impl,  # Access the underlying sklearn model
        artifact_path="my_model",
        registered_model_name="MyRegisteredModel1",
        signature=signature,
        input_example=X_sample
    )
```

▸ (8) Spark Jobs

▾ (1) MLflow run

Logged 1 run ⧉ to an experiment ⧉ in MLflow. Learn more ⧉

Downloading artifacts: 100% �preview▬▬▬▬▬▬▬ 9/9 [00:00<00:00, 30.95it/s]

/databricks/python/lib/python3.11/site-packages/mlflow/types/utils.py:435: UserWarning: Hint: Inferred schema contains integer column(s). Integer columns in Python cannot represen
issing values. If your input data contains missing values at inference time, it will be encoded as floats and will cause a schema enforcement error. The best way to avoid this pro
m is to infer the model schema based on a realistic data sample (training dataset) that includes missing values. Alternatively, you can declare integer columns as doubles (float64
henever these columns may have missing values. See `Handling Integers With Missing Values <https://www.mlflow.org/docs/latest/models.html#handling-integers-with-missing-values>`_
more details.
  warnings.warn(

Uploading artifacts: 100% ▬▬▬▬▬▬▬ 11/11 [00:00<00:00, 14.40it/s]

Successfully registered model 'dbw_dev_westeurope_j2ta.default.myregisteredmodel1'.

I created a served endpoint. Access token and model URL obtained:

Serving endpoints ›

## test-model-name

⋮  🔔 Alerts  Permissions  Edit  Stop

⊘ Ready  https://adb-7405616284574485.5.azuredatabricks.net/serving-endpoints/test-model-name/invocations  ⧉

**About this endpoint**

Created by  👤 Yurii Hordiichuk

Route  Not enabled

optimization

**Gateway**                                          Edit AI Gateway

Usage monitoring: ⓘ  system.serving.endpoint_usage      Rate limits:  Not enabled

Dimension table: ⓘ  system.serving.served_entities

Inference tables:  Not enabled

**Tags**

Add

**Serverless usage policy** ⓘ

✎

**Description**

No description

✎

## Active configuration

| Entity | Version | Name | State | Compute | Traffic (%) |
|---|---|---|---|---|---|
| ■ dbw_dev_westeurope_j2ta.default.myregisteredmodel1 | Version 1 | myregisteredmodel1-1 | ⊘ Ready | **CPU, Custom** 0-4 concurrency (0-4 DBU) | 100 |

5. The model prediction from the endpoint agree with the result from local model:

The model predictions from the endpoint should agree with the results from locally evaluating the model.

▶  ✓  Just now (<1s)                                                                                      60

```python
# Model serving is designed for low-latency predictions on smaller batches of data
num_predictions = 5
served_predictions = score_model(X_test[:num_predictions])
model_evaluations = model.predict(X_test[:num_predictions])

# Compare the results from the deployed model and the trained model
pd.DataFrame({
    "Model Prediction": model_evaluations,
    "Served Model Prediction": np.array(served_predictions['predictions'], dtype=np.float32),
})
```

|   | Model Prediction | Served Model Prediction |
|---|---|---|
| 0 | 0.000269 | 0.000269 |
| 1 | 0.003019 | 0.003019 |
| 2 | 0.020137 | 0.020137 |
| 3 | 0.044904 | 0.044904 |
| 4 | 0.019414 | 0.019414 |

Request appeared on the chart: