Module

Going back to the session

Go into module 4 folder and run :

```
None
python create_onnx_model.py
```

## What This Script Does

- Loads original bank marketing data plus newly labeled customer data
- Combines datasets and trains a fresh RandomForest model with improved performance
- Converts the sklearn model to ONNX format for optimized inference serving
- Registers the ONNX model in MLflow's Model Registry with proper versioning and metadata
- Makes the model ready for deployment to Cloudera AI Inference Service
- Outputs JSON file with performance metrics (F1 score, accuracy) and MLflow registration details

Ensure the script runs successfully.

Then go to the main CAI page

**CLOUDERA**
Data Platform

✕

⌂ Home

▤ Data Warehouse

◔ Operational Database

⟐ Machine Learning

⊙ Data Hub Clusters
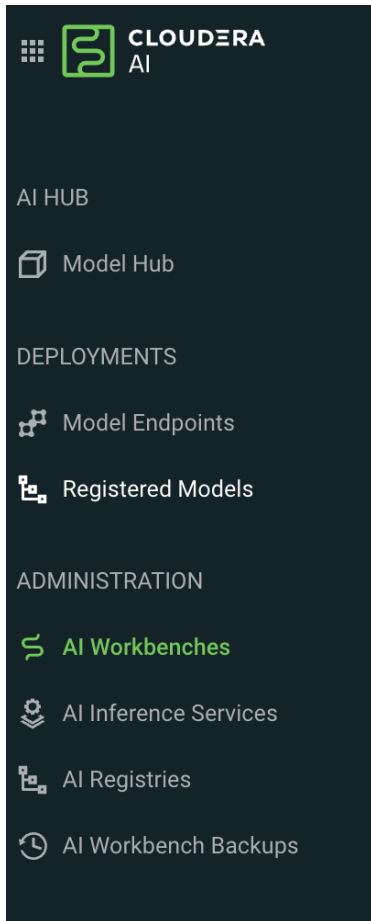
⊙ Data Catalog

⊐ Replication Manager

⊚ Observability

⊟ Management Console

Check to make sure you model got registered (select registered models):

Once there,
1. pick the right model registry
2. Look for your model (identified w/ you username)



Click on your model and you can see some important model information such as model id, metrics, parameters, model version and more.

Go ahead and click on the deploy button. And select the environment and inference services.



Fill in the following information - name "bank-marketing-campaign-your user name" and pick the instance type show below. Select 12 vcpus and 48 GB of ram for memory :

Description

## Served Model Builder

Select the model and version to deploy, then use the Traffic slider to set the traffic split. The first version always defaults to 100% traffic.

| * Model | * Version | Traffic |
|---------|-----------|---------|
| BankingCampaignPredictor_ONNX_ozarate | 1 | 100 |
| | | 0          100 |

## Resource Profile

Configure minimum resource requests for your endpoint here. For a multi-replica endpoint, the requested resources will be allocated to each replica.
Please note that you must not change GPU allocation for optimized models.

Instance Type ⊙

| m5.4xlarge | 16 CPU | 64 GiB | ⌄ |

GPU ⊙

0

| * CPU ⊙ | * Memory ⊙ |
|---------|-----------|
| 12                                    vCPU | 48                                    Gi |

Endpoint Autoscale Range ⊙

0                                              10          1    -    3

> Advanced Options