# Assignment – CodePro Leadscoring MLOps – Vaibhav Jain

## Model Experimentation - lead_scoring_model_experimentation.db

MLFlow UI Baseline model

# MLFlow UI Baseline model – One model with all artifacts

Lead_scoring_Baseline_model_01 > Light Gradient Boosting Machine

## Light Gradient Boosting Machine

Date: 2023-03-20 22:47:41

Source: 💻 ipykernel_launcher.py

User: root

Status: UNFINISHED

Lifecycle Stage: active

Parent Run: ad25adcaff104ec7a45fb0f72dd9dae4

▸ Description    Edit

▸ Parameters (20)

▸ Metrics (8)

▸ Tags (5)

▾ Artifacts

▾ 📁 model
    📄 MLmodel
    📄 conda.yaml
    📄 model.pkl
    📄 python_env.yaml
    📄 requirements.txt
📄 Holdout.html

Full Path:/home/Assignment/02_training_pipeline/mlruns/1/1df56bbb35e7410a814259bd400d4359/artifacts/model 📋

[Register Model]

## MLflow Model

The code snippets below demonstrate how to make predictions using the logged model. You can also register it to the model registry to version control

### Model schema

Input and output schema for your model. Learn more

| Name | Type |
|------|------|
| No schema. See MLflow docs for how to include input and output schema with your model. | |

### Make Predictions

Predict on a Spark DataFrame:

```
import mlflow
logged_model = 'runs:/1df56bbb35e7410a814259bd400d4359/model'

# Load model as a Spark UDF. Override result_type if the model does not return double values.
loaded_model = mlflow.pyfunc.spark_udf(spark, model_uri=logged_model, result_type='double')

# Predict on a Spark DataFrame.
columns = list(df.columns)
df.withColumn('predictions', loaded_model(*columns)).collect()
```

MLFlow UI Tuned model (Features dropped)



**Experiments** ➕ ◀

Search Experiments

Default ✏ 🗑
Lead_scoring_Baselin... ✏ 🗑
Tuned_model_exp01 ✏ 🗑

**Tuned_model_exp01** 📋                                    Share

ℹ Track machine learning training runs in experiments. Learn more                    ✕

Experiment ID: 2

▼ Description  Edit

**Vaibhav Jain - Tuned_model_exp01**

🔄 Refresh | Compare | Delete | ⬇ Download CSV | ↓ Start Time ⌄ | All time ⌄

☰ ▦ | ⚙ Columns | Only show differences ⬤ ❓ | 🔍 metrics.rmse < 1 and params.model = "tree" | Search | ⇶ Filter | Clear

Showing 12 matching runs

| | ↓ Start Time | Duration | Run Name | User | Source | Version | Models | AUC | Accuracy | F1 | C | CPU Jobs | Categorical Feat |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| ☐ | ⊟ 8 minutes ago | | Session Initialized 5a44 | root | ☐ ipykernel... | - | - | - | - | - | - | -1 | 4 |
| ☐ | 3 seconds ago | | Light Gradient Boosting Machine | root | ☐ ipykernel... | - | 📄 sklearn | 0.82 | 0.738 | 0.761 | - | - | - |
| ☐ | 5 minutes ago | | Light Gradient Boosting Machine | root | ☐ ipykernel... | - | 📄 sklearn | 0.821 | 0.739 | 0.762 | - | - | - |
| ☐ | 6 minutes ago | | Naive Bayes | root | ☐ ipykernel... | - | 📄 sklearn | 0.734 | 0.67 | 0.723 | - | - | - |
| ☐ | 6 minutes ago | | Linear Discriminant Analysis | root | ☐ ipykernel... | - | 📄 sklearn | 0.773 | 0.7 | 0.728 | - | - | - |
| ☐ | 6 minutes ago | | Ridge Classifier | root | ☐ ipykernel... | - | 📄 sklearn | 0 | 0.7 | 0.728 | - | - | - |
| ☐ | 6 minutes ago | | Logistic Regression | root | ☐ ipykernel... | - | 📄 sklearn | 0.784 | 0.71 | 0.74 | 1.0 | - | - |
| ☐ | 6 minutes ago | | Decision Tree Classifier | root | ☐ ipykernel... | - | 📄 sklearn | 0.817 | 0.736 | 0.758 | - | - | - |
| ☐ | 6 minutes ago | | Extra Trees Classifier | root | ☐ ipykernel... | - | 📄 sklearn | 0.817 | 0.737 | 0.758 | - | - | - |
| ☐ | 6 minutes ago | | Random Forest Classifier | root | ☐ ipykernel... | - | 📄 sklearn | 0.818 | 0.737 | 0.759 | - | - | - |
| ☐ | 6 minutes ago | | Extreme Gradient Boosting | root | ☐ ipykernel... | - | 📄 sklearn | 0.821 | 0.738 | 0.762 | - | - | - |
| ☐ | 6 minutes ago | | Light Gradient Boosting Machine | root | ☐ ipykernel... | - | 📄 sklearn | 0.821 | 0.739 | 0.762 | - | - | - |

Load more

# MLFlow UI – Best Model with all artifacts



## Best Model Artifacts

# Best Model Artifacts

# Data Pipeline - Airflow UI

My User – Vaibhav Jain



Pipeline Graph view

Pipeline Grid view

DAG: Lead_Scoring_Data_Engineering_Pipeline  DAG to run data pipeline for lead scoring

Schedule: @daily ⓘ    Next Run: 2023-03-19, 00:00:00

▦ Grid    ⫘ Graph    📅 Calendar    ⧗ Task Duration    ⇄ Task Tries    ⚓ Landing Times    ☰ Gantt    ⚠ Details    <> Code    🔍 Audit Log    ▶  🗑

19/03/2023 04:45:08 PM  📅    25  ⌄    All Run Types ⌄    All Run States ⌄    **Clear Filters**

deferred  failed  queued  running  scheduled  skipped  success  up_for_reschedule  up_for_retry  upstream_failed  no_status

Auto-refresh  ⬤                                                              →≡

Duration

00:02:13

00:01:06

00:00:00

building_db
checking_raw_data_schema
loading_data
mapping_city_tier
mapping_categorical_vars
mapping_interactions
checking_model_inputs_schema

DAG
**Lead_Scoring_Data_Engineering_Pipeline**

**DAG Details**

| DAG Runs Summary | |
|---|---|
| Total Runs Displayed | 2 |
| ■ Total success | 2 |
| First Run Start | 2023-03-19, 16:41:54 UTC |
| Last Run Start | 2023-03-19, 16:41:54 UTC |
| Max Run Duration | 00:02:13 |
| Mean Run Duration | 00:02:13 |
| Min Run Duration | 00:02:13 |
| DAG Summary | |
| Total Tasks | 7 |
| PythonOperators | 7 |

# Training Pipeline – Airflow UI



Pipeline Grid view

MLFlow – Production DB



Model staged as Production manually via MLFlow UI

# MLFlow - Production DB - Model Artifacts



ml*flow* 1.26.1    **Experiments**    **Models**                                    GitHub    Docs

Lead_scoring_mlflow_production  >  run_LightGB

## run_LightGB

Date: 2023-03-20 23:25:44              Source: 🖥 airflow              User: root

Duration: 5.0s                         Status: FINISHED              Lifecycle Stage: active

▸ Description    Edit

▸ Parameters (20)

▸ Metrics (1)

▸ Tags

▾ Artifacts

| ▾ 📁 models | Full Path:/home/Assignment/02_training_pipeline/mlruns/1/efd27af4b24f404e9c4cb96c55cff8e6/artifacts/models 🗍 | ⊘ LightGBM, v1 ⬈ |
|---|---|---|
| 📄 MLmodel | | Registered on 2023/03/20 |
| 📄 conda.yaml | |
| 📄 model.pkl | ### MLflow Model |
| 📄 python_env.yaml | The code snippets below demonstrate how to make predictions using the logged model. This model is also registered to the model registry. |
| 📄 requirements.txt | |

### Model schema

Input and output schema for your model. Learn more

| Name | Type |
|---|---|
| No schema. See MLflow docs for how to include input and output schema with your model. | |

### Make Predictions

Predict on a Spark DataFrame:

```
import mlflow
logged_model = 'runs:/efd27af4b24f404e9c4cb96c55cff8e6/models'

# Load model as a Spark UDF. Override result_type if the model does not return double values.
loaded_model = mlflow.pyfunc.spark_udf(spark, model_uri=logged_model, result_type='double')

# Predict on a Spark DataFrame.
columns = list(df.columns)
df.withColumn('predictions', loaded_model(*columns)).collect()
```

# Inference Pipeline – Airflow UI

Airflow – 3 dags ready