

MLflow

Introduction to Experiment Tracking and Model Management

MLflow is an open-source platform to manage the ML lifecycle, including experimentation, reproducibility, deployment, and a central model registry.

MLflow currently offers four components:

- MLflow Tracking-Record and query experiments: code, data, config, and results
- MLflow projects-Package data science code in a format to reproduce runs on any platform
- MLflow models- Deploy Machine learning models in diverse serving environments
- Model Registry- Store, annotate, discover, and manage models in a central repository

MLflow: Tracking of Experiments, logging or recording all the Experiments.

Login is important for Organization of Production pipeline properly.

MLflow interface is helping to track for-

- What kind of Algorithm
- What kind of Hyper parameters
- What kind of scores we get for the model.

MLflow is an open source platform for managing the end-to-end machine learning lifecycle. It tackles four primary functions:

- Tracking experiments to record and compare parameters and results ([MLflow Tracking](#)).
- Packaging ML code in a reusable, reproducible form in order to share with other data scientists or transfer to production ([MLflow Projects](#)).
- Managing and deploying models from a variety of ML libraries to a variety of model serving and inference platforms ([MLflow Models](#)).
- Providing a central model store to collaboratively manage the full lifecycle of an MLflow Model, including model versioning, stage transitions, and annotations ([MLflow Model Registry](#)).

Concepts

MLflow Tracking is organized around the concept of *runs*, which are executions of some piece of data science code. Each run records the following information:

Code Version:

Git commit hash used for the run, if it was run from an [MLflow Project](#).

Start & End Time

Start and end time of the run

Source

Name of the file to launch the run, or the project name and entry point for the run if run from an [MLflow Project](#).

Parameters

Key-value input parameters of your choice. Both keys and values are strings.

Metrics

Key-value metrics, where the value is numeric. Each metric can be updated throughout the course of the run (for example, to track how your model's loss function is converging), and MLflow records and lets you visualize the metric's full history.

Artifacts

Output files in any format. For example, you can record images (for example, PNGs), models (for example, a pickled scikit-learn model), and data files (for example, a [Parquet](#) file) as artifacts.

Introduction to Experiment Tracking

Terminologies:

1. Experiment
2. Run
3. Metadata (i.e. Tags, Parameters, Metrics)
4. Artifacts (i.e. Output files associated with experiment runs)

What do you want to track for each Experiment Run?

1. Training and Validation Data Used
2. Hyperparameters
3. Metrics
4. Models

Why Track?

Organization Optimization Reproducibility

Tool - MLFlow

MLFlow helps you to organize your experiments into runs.

MLFlow keeps track of:

- Tags
- Parameters
- Metrics
- Models
- Artifact
- Source code, Start and End Time, Authors etc..

Run below mentioned commands to install mlflow on your system:

- `pip install mlflow`
- `mlflow ui --backend-store-uri sqlite:///mlflow.db`

Introduction to MLFlow

Step 1 - Import MLFlow

```
import mlflow
```

Step 2 - Set the tracker and experiment

```
mlflow.set_tracking_uri(DATABASE_URI)
```

```
mlflow.set_experiment("EXPERIMENT_NAME")
```

Step 3 - Start a experiment run

```
with mlflow.start_run():
```

Step 4 - Logging the metadata

```
mlflow.set_tag(KEY, VALUE)
```

```
mlflow.log_param(KEY, VALUE) mlflow.log_metric(KEY, VALUE)
```

Step 5 - Logging the model and other files (2 ways)

Way 1 - `mlflow.<FRAMEWORK>.log_model(MODEL_OBJECT, artifact_path="PATH")`

Way 2 - `mlflow.log_artifact(LOCAL_PATH, artifact_path="PATH")`

Below are screenshots of Experiment tracking and Model Management of diamond price prediction models:

The screenshot displays the MLflow web interface for experiment tracking. The top navigation bar includes the MLflow logo, version 1.29.0, and tabs for 'Experiments' and 'Models'. On the left, the 'Experiments' sidebar shows a search bar and a list of experiments: 'Default' and 'Diamond Price Prediction' (selected). The main panel is titled 'Diamond Price Prediction' and shows 'Experiment ID: 1'. It includes a description field, a 'Description Edit' button, and a table of runs. The table has columns for 'Created', 'Duration', 'Run Name', 'Metrics' (mean_ab_error, mean_sqr_error, mean_srt_error), 'Parameters' (data-path), and 'Tags' (dev, algo). Three runs are listed, all created 8, 10, and 10 minutes ago. The first run is 'hilarious' with a duration of 40.8s and metrics (366.9, 573210.8, 757.1). The second run is 'nebulous' with a duration of 2.3min and metrics (273.8, 308573.3, 555.5). The third run is 'glamorous' with a duration of 18.3s and metrics (411.8, 686410.3, 828.5). The interface also features buttons for 'Refresh', 'Compare', 'Delete', 'Download CSV', and a search bar.

	Created	Duration	Run Name	Metrics	Parameters	Tags			
				mean_ab_error	mean_sqr_error	mean_srt_error	data-path	dev	algo
<input type="checkbox"/>	8 minutes ago	40.8s	hilarious	366.9	573210.8	757.1	data/diamonds...	Rinku Soni	Decision Tree
<input type="checkbox"/>	10 minutes ago	2.3min	nebulous	273.8	308573.3	555.5	data/diamonds...	Rinku Soni	Random For...
<input type="checkbox"/>	10 minutes ago	18.3s	glamorous	411.8	686410.3	828.5	data/diamonds...	Rinku Soni	KNN

Comparing 3 Runs from 1 Experiment

Visualizations

Parallel Coordinates Plot Scatter Plot Box Plot Contour Plot

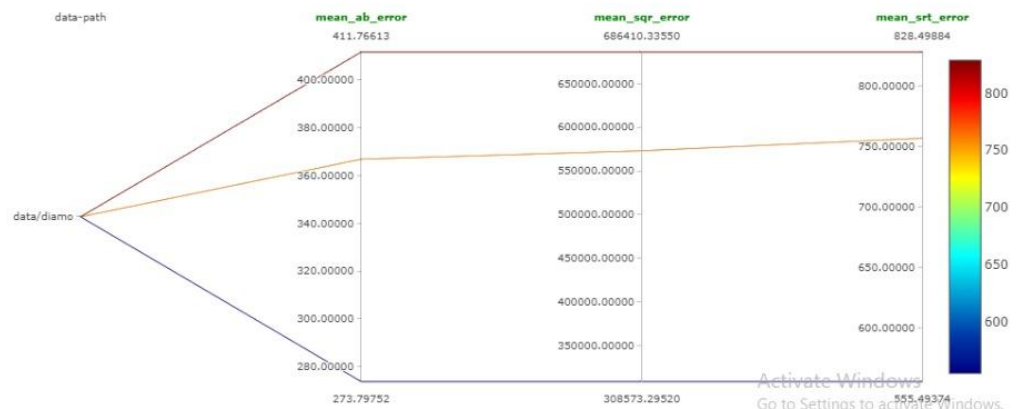
Parameters:

data-path X

Metrics:

mean_ab_error X mean_sqr_error X
mean_srt_error X

Clear All



mlflow 1.29.0 Experiments Models GitHub Docs

Registered Models > Diamond_Price_Prediction_Model

Diamond_Price_Prediction_Model

Created Time: 2022-10-01 19:31:28 Last Modified: 2022-10-01 20:04:28

> Description Edit

> Tags

▼ Versions All Active 2 Compare

<input type="checkbox"/>	Version	Registered at	Created by	Stage	Description
<input type="checkbox"/>	Version 3	2022-10-01 19:40:49		Archived	
<input type="checkbox"/>	Version 2	2022-10-01 19:40:36		Staging	
<input type="checkbox"/>	Version 1	2022-10-01 19:31:29		Production	

PREFECT ORCHESTRATE ML PIPELINES

Managing Machine Learning Workflows using the tool Prefect 2.0

Why Prefect?

- Python-based open source tool
- Manage ML Pipelines
- Schedule and Monitor the flow
- Gives observability into failures
- Native dask integration for scaling (Dask is used for parallel computing)

Creating and activating a Virtual Environment

In order to install prefect, create a virtual environment:

```
$ python -m venv mlops
```

Enter the Virtual Environment using below mentioned command:

```
$ .\mlops\Scripts\activate
```

Installing Prefect 2.0

Now install Prefect:

```
$ pip install prefect
```

OR if you have Prefect 1, upgrade to Prefect 2 using this command:

```
$ pip install -U prefect
```

OR to install a specific version:

```
$ pip install prefect==2.4
```

Check Prefect Version

Check the prefect version:

```
$ prefect version
```

Running Prefect Dashboard

```
$ prefect orion start
```

```
-----
|_ \ \ | | / | | | / \ | \ | \ | | | | |
| / / | | | ( | | | ( | / | ( | : |
| | | \ | | | \ | | | \ | | \ | \ |
Configure Prefect to communicate with the server with:
    prefect config set PREFECT_API_URL=http://127.0.0.1:4200/api
View the API reference documentation at http://127.0.0.1:4200/docs
Check out the dashboard at http://127.0.0.1:4200/
```

Note - In Windows OS, if your path contains spaces, it will generate error (as mentioned below) when you try to run prefect orion.

Prefect O/p screens:

This screenshot shows the Prefect interface for a flow named 'auburn-magpie'. The 'Logs' tab is selected, displaying a list of log entries for task runs. The logs show the creation and execution of tasks: 'load_data', 'split_data', 'seperating_numerical', 'seperating_categorical', 'encoding', and 'concat_df'. Each entry includes a timestamp and a status (e.g., 'Completed').

Flow Runs / auburn-magpie

Logs Task Runs Sub Flow Runs Parameters

Level: all

Oct 2nd, 2022

- INFO Created task run 'load_data-2ff00c39-0' for task 'load_data' 10:04:54 PM
- INFO Executing 'load_data-2ff00c39-0' immediately... 10:04:54 PM
- INFO Created task run 'split_data-b2f518fa-0' for task 'split_data' 10:04:58 PM
- INFO Executing 'split_data-b2f518fa-0' immediately... 10:04:58 PM
- INFO Created task run 'seperating_numerical-08a2c8b5-0' for task 'seperating_numerical' 10:05:01 PM
- INFO Executing 'seperating_numerical-08a2c8b5-0' immediately... 10:05:01 PM
- INFO Created task run 'seperating_categorical-81c18760-0' for task 'seperating_categorical' 10:05:02 PM
- INFO Executing 'seperating_categorical-81c18760-0' immediately... 10:05:02 PM
- INFO Created task run 'encoding-bf9aef86-0' for task 'encoding' 10:05:03 PM
- INFO Executing 'encoding-bf9aef86-0' immediately... 10:05:03 PM
- INFO Created task run 'concat_df-ec37f627-0' for task 'concat_df' 10:05:05 PM
- INFO Executing 'concat_df-ec37f627-0' immediately... 10:05:05 PM

Flow Run ID: ce45494e-5ba0-40d9-8c3d-174286746314

Flow ID: eb4557bc-989d-4081-b320-f9cd2ab270d2

Created: 2022/10/02 10:04:49 PM

Updated: 2022/10/03 12:09:13 AM

This screenshot shows the Prefect interface for the same flow 'auburn-magpie', but with the 'Task Runs' tab selected. It displays a list of completed task runs, each with a unique ID, status, and completion time. The tasks listed are 'find_best_model', 'rescale_num_data', 'get_scaler', 'concat_df', 'encoding', and 'seperating_categorical'.

Flow Runs / auburn-magpie

Logs Task Runs Sub Flow Runs Parameters

All states Search by run name Newest to Oldest

- find_best_model-d603fb17-0
Completed 2h 3m 2022/10/02 10:05:18 PM
- rescale_num_data-ec0e7e66-1
Completed 1s 2022/10/02 10:05:17 PM
- get_scaler-9c4a96f8-1
Completed 1s 2022/10/02 10:05:16 PM
- concat_df-ec37f627-1
Completed 1s 2022/10/02 10:05:14 PM
- encoding-bf9aef86-1
Completed 1s 2022/10/02 10:05:13 PM
- seperating_categorical-81c18760-1
Completed 1s 2022/10/03 10:05:11 PM

Flow Run ID: ce45494e-5ba0-40d9-8c3d-174286746314

Flow ID: eb4557bc-989d-4081-b320-f9cd2ab270d2

Created: 2022/10/02 10:04:49 PM

Updated: 2022/10/03 12:09:13 AM

Flow Runs

Flows

Deployments

Work Queues

Blocks

Notifications

Settings

Flow Runs / auburn-magpie

Logs

Task Runs

Sub Flow Runs

Parameters

```
{  
  "path": "../data/diamonds.csv"  
}
```

Completed

2h 4m

main

2022/10/02 10:04:53 PM (4s late)

Flow Run ID

ce45494e-5ba0-40d9-8c3d-174286746314

Flow ID

eb4557bc-989d-4081-b320-f9cd2ab270d2

Created

2022/10/02 10:04:49 PM

Updated

2022/10/03 12:09:13 AM

Flow Runs

Flows

Deployments


Work Queues

Blocks

Notifications

Settings

Flow Runs / meticulous-partridge / Radar



Activate Windows

Go to Settings to activate Windows.

Diamond Price Prediction > Comparing 61 Runs from 1 Experiment

Comparing 61 Runs from 1 Experiment

Visualizations

Parallel Coordinates Plot Scatter Plot Box Plot Contour Plot

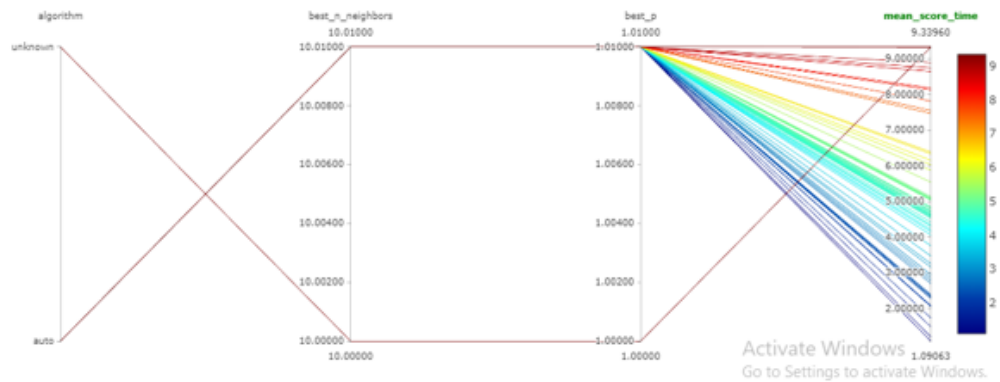
Parameters:

algorithm X best_n_neighbors X best_p X

Metrics:

mean_score_time X

Clear All



Diamond Price Prediction > Comparing 62 Runs from 1 Experiment

Comparing 62 Runs from 1 Experiment

Visualizations

Parallel Coordinates Plot Scatter Plot Box Plot Contour Plot

Parameters:

algorithm X best_n_neighbors X best_p X

Metrics:

training_mae X training_rmse X

Clear All

