# MLOps ML Flow Assignment

**Nitish Bhardwaj (B21AI056)**

## Introduction

This report presents an analysis of housing prices using two different machine learning models: Linear Regression and Random Forest Regression. The dataset used for this analysis is the Boston Housing dataset, which contains various features that influence housing prices in Boston.

## Data Acquisition

The dataset was sourced from a publicly available URL. It was loaded into a Pandas DataFrame for further analysis and modeling.

## Dataset Overview

The dataset consists of several features, with the target variable being medv, which represents the median value of owner-occupied homes in thousands of dollars. The features include variables such as crime rate, average number of rooms, accessibility to highways, and others.

## Data Preparation

The features (X) were separated from the target variable (y). The dataset was then split into training and testing sets, with 80% of the data used for training and 20% reserved for testing.

## Model Selection

Two models were chosen for the prediction task:

1. **Linear Regression**
2. **Random Forest Regressor**

These models were implemented using Scikit-learn, a popular machine learning library in Python.

## Model Training and Evaluation

An experiment was logged using MLflow, allowing for easy tracking of model performance metrics and parameters.

**Results**

- **Linear Regression**:
  - Mean Squared Error (MSE): 24.29
  - R-squared: 0.67
- **Random Forest**:
  - Mean Squared Error (MSE): 7.90
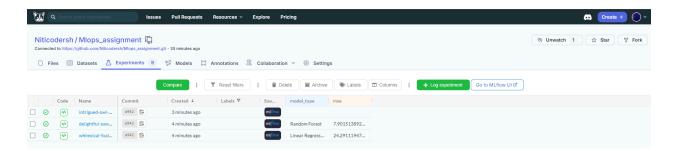  - R-squared: [Value]: 0.89

**Conclusion**

The performance of both models was evaluated using MSE and R-squared metrics. These metrics provide insight into the accuracy and reliability of the predictions made by each model. The use of MLflow facilitated the logging and tracking of experiments, ensuring reproducibility and ease of comparison between models.
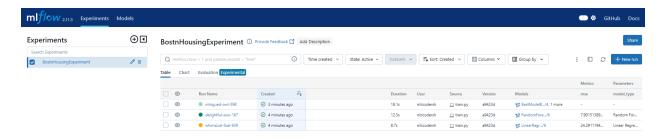
**ScreenShots**

Locally running train.py, MLflow experiments logs are created on DagsHub (a remote server set as the tracking uri): https://dagshub.com/Niticodersh/Mlops_assignment.mlflow/#/experiments
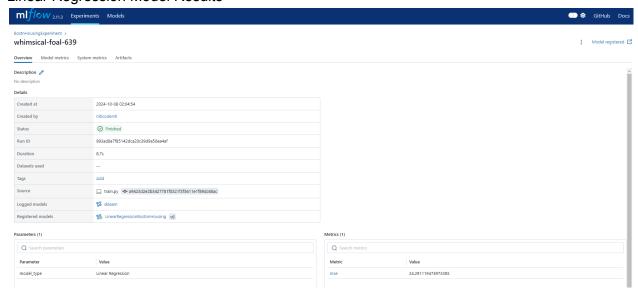
Screenshots of Experiments performed hosted at DagsHub:
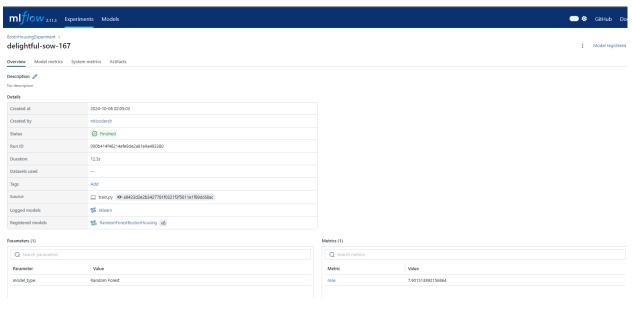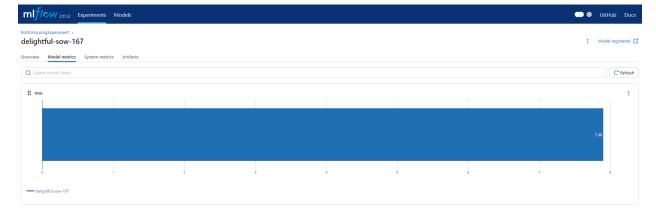


Screenshots of MLFlow UI hosted at DagsHub:



Linear Regression Model Results

# whimsical-foal-639

Overview   Model metrics   System metrics   Artifacts

🔍 Search metric charts                                    ⟳ Refresh

⠿ mse                                                        ⋮

24.29

0        5        10        15        20        25

— whimsical-foal-639

## Random Forest Model Results

**ml flow** 2.11.3   Experiments   Models                          GitHub   Do

# delightful-sow-167

Overview   Model metrics   System metrics   Artifacts

**Description** ✏️

No description

**Details**

| | |
|---|---|
| Created at | 2024-10-08 02:05:03 |
| Created by | niticodersh |
| Status | ✓ Finished |
| Run ID | 000b414f46214efe9de2a81e9a493380 |
| Duration | 12.3s |
| Datasets used | — |
| Tags | Add |
| Source | 💻 train.py  ⟠ a9423d2e2b3427781f0321f3f5611e1f89dc68ac |
| Logged models | 🔗 sklearn |
| Registered models | 🔗 RandomForestBostonHousing  v6 |

**Parameters (1)**

🔍 Search parameters

| Parameter | Value |
|---|---|
| model_type | Random Forest |

**Metrics (1)**

🔍 Search metrics

| Metric | Value |
|---|---|
| mse | 7.901513892156864 |

**ml flow** 2.11.3   Experiments   Models                    GitHub   Docs

# delightful-sow-167

Overview   Model metrics   System metrics   Artifacts

🔍 Search metric charts                                    ⟳ Refresh

⠿ mse                                                        ⋮

7.90

0      1      2      3      4      5      6      7      8

— delightful-sow-167

# Both Models MSE Comparisons Result

model_type

Random For

Linear Reg

mse
24.29112

24.00000

22.00000

20.00000

18.00000

16.00000

14.00000

12.00000

10.00000

8.00000

7.90151

24

22

20

18

16

14

12

10

8