# MLflow Model Tracking Report

## 1. Introduction

This report outlines the experiments conducted to compare the performance of Linear Regression and Random Forest models on the California housing dataset. The models were trained, logged, and evaluated using MLflow, a popular platform for managing the machine learning lifecycle.

## 2. Dataset

The dataset used for this experiment is the **California housing dataset**, which includes various features related to housing prices. The target variable to predict is the median house value.

## 3. Model Comparison

### Mean Squared Error (MSE) Values

The following table summarizes the Mean Squared Error (MSE) values obtained for each model:

| Model | MSE |
|---|---|
| Linear Regression | 0.5559 |
| Random Forest | 0.2554 |

### Analysis

- **Linear Regression**: The MSE of the Linear Regression model is **0.5559**. This indicates a relatively higher error in predictions, suggesting that this model does not capture the underlying patterns in the data effectively.
- **Random Forest**: The MSE of the Random Forest model is **0.2554**, significantly lower than that of Linear Regression. This indicates that the Random Forest model provides more accurate predictions by better capturing the complexities in the data.
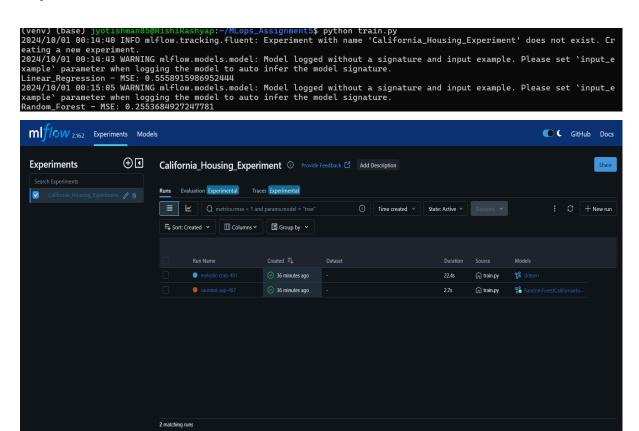
### Conclusion

Based on the MSE values, the **Random Forest** model is identified as the better-performing model for this task, demonstrating superior accuracy in predicting housing prices compared to the Linear Regression model.
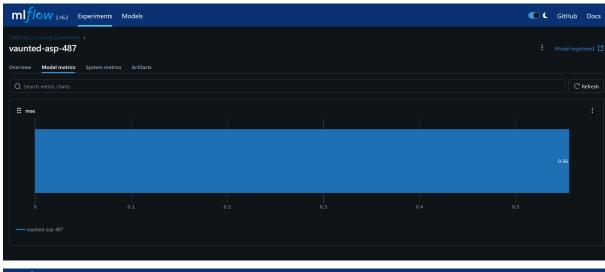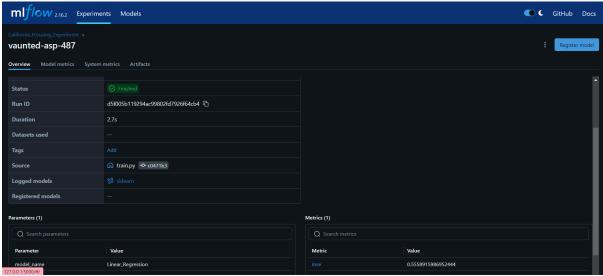
# 4. MLflow UI Screenshots

The MLflow UI was used to track the experiments, log metrics, and visualize model performance. Below are the relevant screenshots from the MLflow UI:
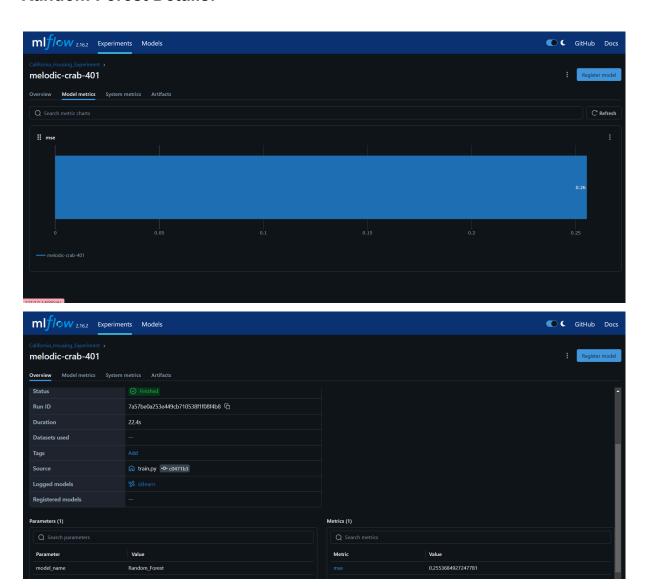
## Experiment Overview:

# Linear Regression Details:

**Random Forest Details:**





# 5. Model Registration

The Random Forest model, identified as the best performer, was registered in MLflow's Model Registry for future use. This allows for versioning and management of the model in production environments.

# 6. Final Remarks

The experiments conducted illustrate the power of using MLflow for managing machine learning workflows. The comparison of the two models highlights the effectiveness of the Random Forest algorithm for this particular dataset and task.