Robert Pavlik

D602 – Deployment

QBN1 – Task 2: Data Production Pipeline

Introduction

This project aimed to develop an end-to-end machine learning pipeline for predicting airport flight delays. I chose to do Miami, Florida Airport (MIA) using June 2024 data. The pipeline involved several steps: data importing, formatting, filtering, cleaning, and model training using a polynomial ridge regression model with MLflow for experiment tracking. The final code was integrated into a MLProject to enable reproducible runs and parameter tuning.

Requirements

A. Gitlab Repository – Task2_Branch
   • https://gitlab.com/wgu-gitlab-environment/student-repos/rpavli5/d602-deployment-task-2.git

B. Import and Format Script



C. Filter and Clean Script

D. Poly_Regressor Script

```
354    ∨ with mlflow.start_run(run_name = 'Test Data Run'):
355          # Log input parameters
356          mlflow.log_param( key: "alpha", parameters[0] / 10)
●           mlflow.log_param( key: "polynomial_order", parameters[1])
358
359          # Log model performance metrics
360          mlflow.log_metric( key: "Test Data Mean Squared Error", score)
361          mlflow.log_metric( key: "Test Data Average Delay", np.sqrt(score))
362
363          # Log the performance plot artifact
364          mlflow.log_artifact("Output/model_performance_test.jpg")
365
366          # Log additional artifacts
367          mlflow.log_artifact("Output/polynomial_regression.txt")
368          mlflow.log_artifact("Data/cleaned_data.csv")
369
370       mlflow.end_run()
371
372       logging.shutdown()
```

E. MLProject File

```
1     name: TestPipeline
2
3     conda_env: pipeline_env.yaml
4
5     entry_points:
6        main:
7           parameters:
8              data: {type: str, default: cleaned_data.csv}
9              num_alphas: {type: int, default: 20}
10          command: "python Main.py {data} {num_alphas}"
11       import_format:
12          parameters:
13             data: {type: str, default: cleaned_data.csv}
14          command: "python Steps/StepB_Import_Format.py {data}"
15       filter_clean:
16          command: "python Steps/StepC_Filter_Clean.py"
17       train_model:
18          parameters:
19             num_alphas: {type: int, default: 20}
20          command: "python Steps/poly_regressor_Python_1.0.0.py {num_alphas}"
21       run_all:
22          parameters:
23             data: {type: str, default: cleaned_data.csv}
24             num_alphas: {type: int, default: 20}
25          command: "python Main.py {data} {num_alphas}"
```

F. MLProject Pipeline Creation

The project was developed using Python and integrated with MLFlow to track model parameters, metrics, and artifacts. The overall goal was to create a reproducible pipeline that downloads, cleans, and filters data, and then runs a polynomial regression model to predict departure delays from a specified airport.

**Data Import and Formatting:**
I began by writing a Python script to import the raw flight data. This script reads data from a CSV file and reformats it to match the structure required by the regression model. Multiple iterations of the code were committed to GitLab, demonstrating a progression in functionality and robustness.

**Data Filtering and Cleaning:**
A separate script was developed to filter the dataset for departures from a chosen airport—in this

case, Miami International Airport (MIA). In addition to filtering by airport, I implemented further data cleaning steps, including:

- Dropping rows with missing or invalid delay values.

- Removing duplicate entries.

- Converting relevant columns from float to integer types where necessary.

Each stage was version-controlled and pushed to GitLab with incremental improvements.

**MLFlow Experiment Integration:**
The provided regression model (poly_regressor) was modified to integrate MLFlow for experiment tracking. In this script:

- Model parameters and performance metrics (mean squared error, average delay) were logged.

- Artifacts, such as the training log files and a performance plot, were saved.

- The MLFlow experiment was made reproducible by linking the run context to a YAML-based MLProject file.

Several iterations of this script were committed to GitLab, showing iterative improvements and troubleshooting.

**MLProject File Creation:**
To connect the data import/cleaning scripts with the MLFlow experiment, I created an MLProject YAML file that:

- Specifies the entry points for both the data processing script and the MLFlow experiment script.

- Defines parameters (including their types and default values) that can be passed during runtime.

- Enables seamless execution of the entire pipeline from the command line using mlflow run.

The pipeline is designed so that other analysts can easily extend the work to different airports by modifying the parameters and re-running the pipeline. The MLProject file and related scripts are version-controlled in the GitLab repository, ensuring full traceability of changes.

**Challenges Encountered and Solutions:**

- **MLFlow Run Context:**
  Managing the MLFlow run context when running via mlflow run was challenging. Initially, I attempted to manually retrieve or start an active run within the script, which led to errors such as "Run not found." This was resolved by removing explicit run management and relying on the MLFlow CLI to handle the run context automatically.

- **Parameter Passing:**
  Integrating the scripts to accept parameters via the command line initially resulted in Key Error: 'data' errors. I addressed this by updating the MLProject file to include a parameter

section for the entry points and ensuring that the Python scripts correctly accepted these parameters, specifying types and default values.

**Conclusion:**

By addressing these challenges through careful updates to the code, MLProject file, and environment management, the final MLProject pipeline was successfully executed using MLFlow. The pipeline now supports dynamic parameter substitution, reproducible Conda environments, and detailed experiment tracking via the MLFlow UI. A screenshot of the successful MLProject run is included below.