

MLOps Project Workflow Documentation

Travel Industry Machine Learning Operations

Project Overview

This document explains the complete MLOps workflow for the travel industry machine learning project developed by Hritikrai55. The project includes three main components:

- Flight Price Prediction
- Hotel Price Prediction
- Gender Classification

The project follows industry best practices for machine learning operations, including automated pipelines, experiment tracking, containerization, and scalable deployment.

1. Data Management and Preparation Stage

Dataset Structure

The project uses three main datasets stored in the `Datasets/` folder:

- **flights.csv** - Raw flight data for training flight price prediction models
- **hotels.csv** - Raw hotel data for training hotel price prediction models
- **users.csv** - User data for analytics and gender classification

Data Processing Workflow

Each module (`flight_price`, `hotel_price`, `gender_classification`) handles its own data processing:

- Data is loaded from CSV files in the `Datasets` folder
- Each module contains Jupyter notebooks for Exploratory Data Analysis (EDA)
- Data cleaning and feature engineering are performed within individual notebooks
- Processed data is used for model training and validation

Key Features of Data Management

- **Modular approach:** Each ML task has its own data processing pipeline
- **Version control:** All datasets are tracked in the repository

- **Reproducibility:** Consistent data processing across different environments
 - **Environment isolation:** Each module has its own requirements and dependencies
-

2. Model Development and Training Stage

Development Environment Setup

Each module follows a consistent structure for model development:

Environment Configuration

- **Virtual environments:** Each module has its own `requirements.txt` file
- **Environment variables:** `.env` files store configuration settings
- **Dependency isolation:** Separate Python environments prevent conflicts

Jupyter Notebook Development

Each module includes dedicated notebooks for model development:

Flight Price Module:

- `eda.ipynb` - Exploratory data analysis for flight data
- `flight_price_prediction.ipynb` - Model development and experimentation

Hotel Price Module:

- `hotel_price_prediction.ipynb` - Model development and experimentation

Gender Classification Module:

- `gender_classification.ipynb` - Model development and experimentation

Model Training Process

The training process follows these steps:

1. **Data Loading:** Raw data is loaded from CSV files
2. **Data Exploration:** EDA notebooks analyze data patterns and distributions
3. **Feature Engineering:** Features are created and transformed for model input
4. **Model Selection:** Different algorithms are tested and compared
5. **Model Training:** Selected models are trained on processed data
6. **Model Validation:** Performance is evaluated using validation datasets
7. **Model Serialization:** Trained models are saved as `.pkl` or `.joblib` files

Model Artifacts

Trained models are stored as serialized files:

- **Pickle files (.pkl):** Standard Python serialization format
 - **Joblib files (.joblib):** Optimized for NumPy arrays and scikit-learn models
 - **Model versioning:** Different versions can be maintained for experiments
-

3. Experiment Tracking with MLflow

MLflow Integration

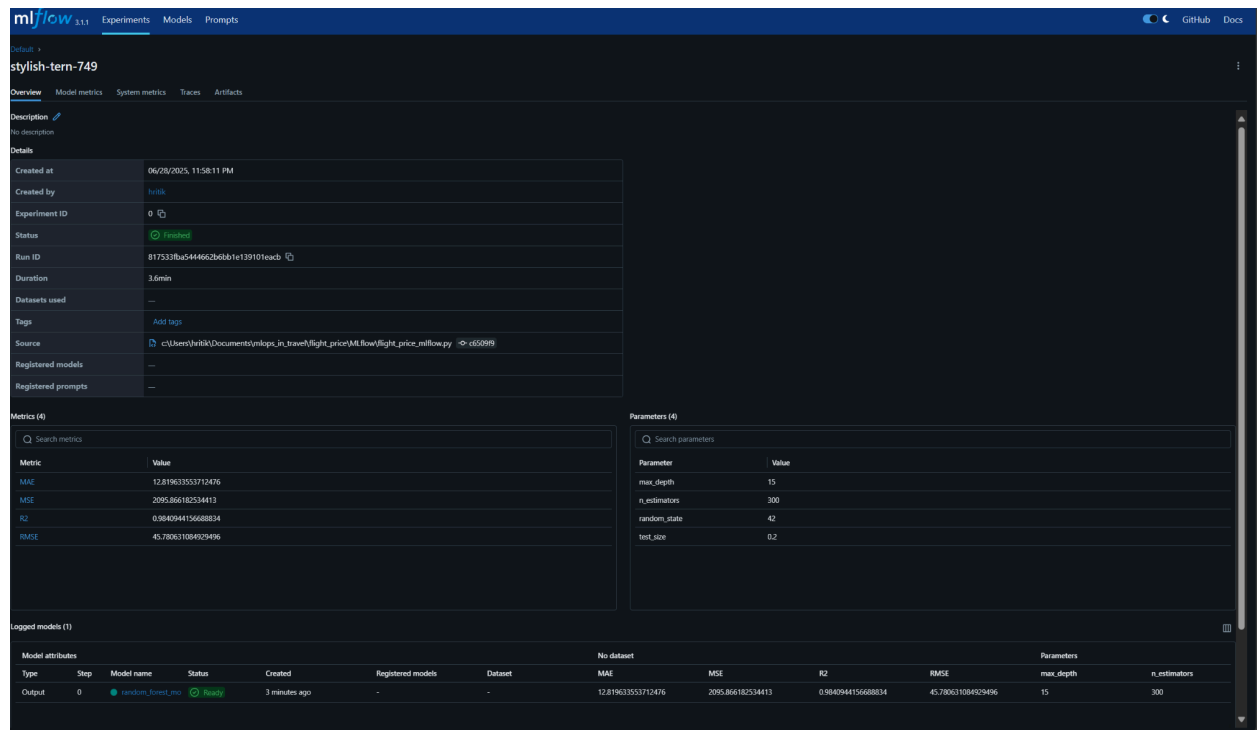
The project uses MLflow for experiment tracking, specifically in the flight price module:

MLflow Setup

- **Location:** `flight_price/MLflow/flight_price_mlflow.py`
- **Purpose:** Track experiments, parameters, metrics, and model artifacts
- **Integration:** Connected with model training workflows

Experiment Tracking Features

MLflow provides comprehensive experiment management:



Metrics Tracking

- Model performance metrics (accuracy, RMSE, MAE, etc.)
- Training and validation scores
- Model comparison across different runs

Parameter Logging

- Hyperparameters used in model training
- Data preprocessing parameters
- Model configuration settings

Artifact Management

- Model files and their versions
- Training datasets and preprocessed data
- Visualizations and plots from experiments

Model Registry

- Centralized model storage and versioning
- Model lifecycle management (staging, production, archived)
- Model lineage and metadata tracking

Benefits of MLflow Integration

- **Reproducibility:** All experiments can be replicated with logged parameters
- **Comparison:** Easy comparison between different model versions
- **History:** Complete audit trail of model development
- **Collaboration:** Team members can access and review experiments

4. Workflow Orchestration with Apache Airflow

Airflow Setup and Configuration

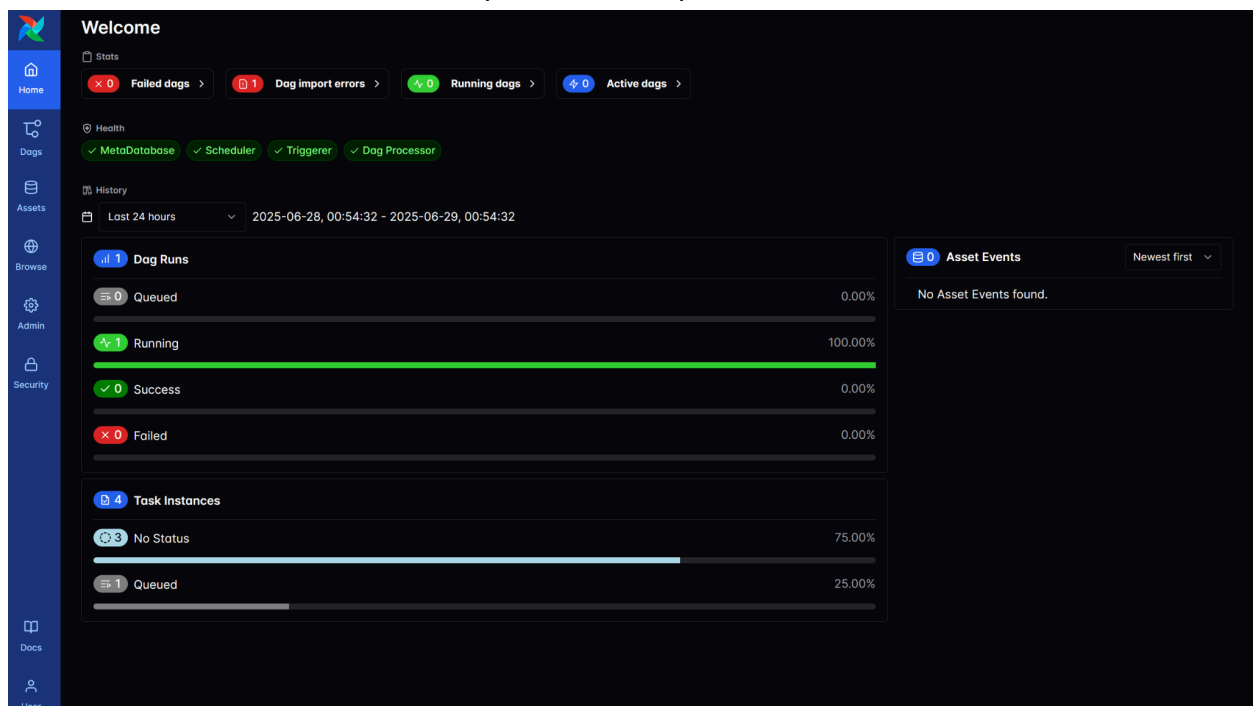
The project uses Apache Airflow for workflow orchestration:

Airflow Structure

- **Location:** `flight_price/airflow_home/`
- **Configuration:** `config/airflow.cfg` contains Airflow settings
- **Docker Compose:** `docker-compose.yaml` for easy Airflow deployment
- **DAGs:** `days/flight_price_dags.py` contains workflow definitions

DAG (Directed Acyclic Graph) Design

Airflow DAGs define the workflow steps and their dependencies:



Typical DAG Structure for ML Workflows

1. **Data Ingestion Task:** Load raw data from sources
2. **Data Validation Task:** Check data quality and consistency
3. **Data Preprocessing Task:** Clean and transform data
4. **Feature Engineering Task:** Create features for model training
5. **Model Training Task:** Train machine learning models
6. **Model Evaluation Task:** Validate model performance
7. **Model Registration Task:** Register approved models in MLflow
8. **Deployment Task:** Deploy models to production environment

Airflow Benefits for MLOps

- **Scheduling:** Automatic execution of ML pipelines on schedules
- **Dependency Management:** Ensures tasks run in correct order
- **Monitoring:** Real-time monitoring of pipeline execution
- **Error Handling:** Automatic retries and failure notifications
- **Scalability:** Can handle multiple concurrent workflows

Pipeline Automation Features

- **Automated Retraining:** Models can be retrained on new data automatically
 - **Data Pipeline Management:** Consistent data processing workflows
 - **Integration:** Seamless connection with MLflow for experiment tracking
 - **Notification System:** Alerts for pipeline success or failure
-

5. Containerization with Docker

Docker Implementation

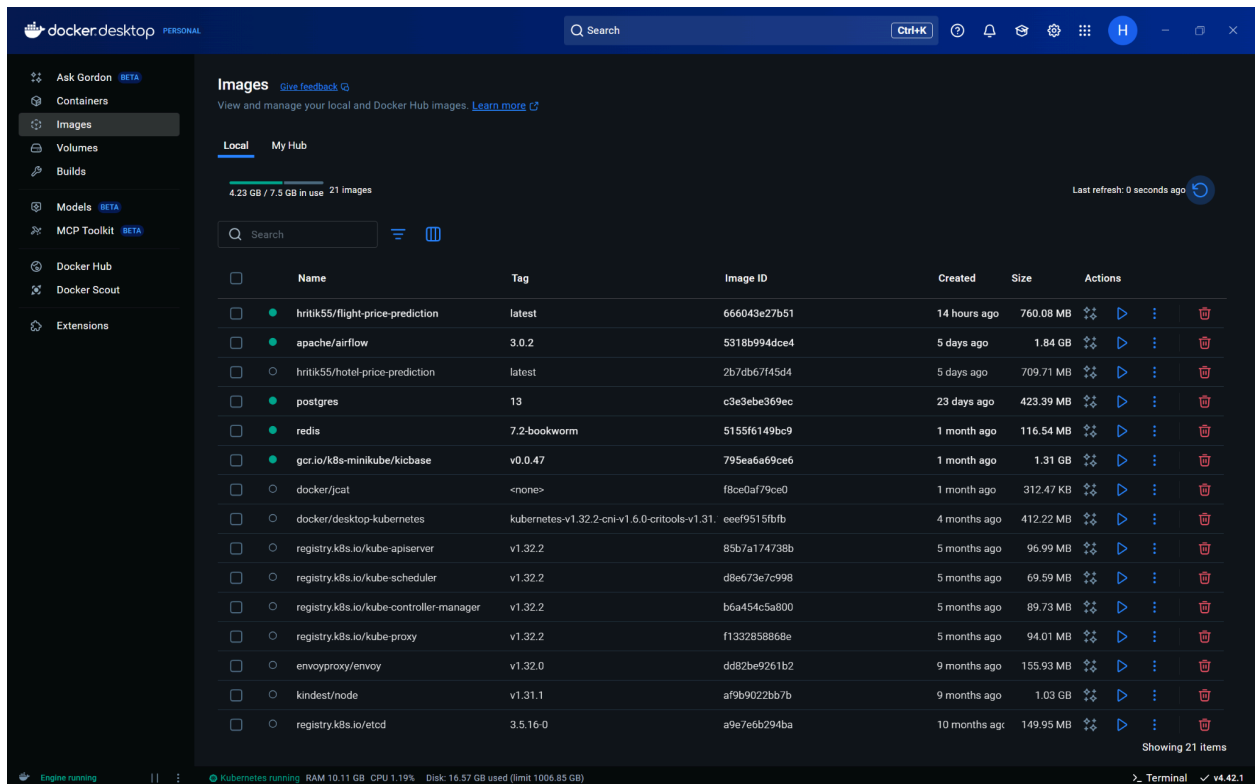
Each module includes Docker configuration for containerization:

Dockerfile Structure

Every module contains a **Dockerfile** that includes:

- **Base Image:** Python runtime environment
- **Dependencies:** Installation of required packages from requirements.txt
- **Application Code:** Copying source code into container

- **Entry Point:** Command to start the application



Container Benefits

- **Consistency:** Same environment across development, testing, and production
- **Portability:** Applications run the same way on any Docker-enabled system
- **Isolation:** Each application runs in its own container environment
- **Scalability:** Easy to scale applications horizontally

Container Management

Docker containers provide several advantages for MLOps:

Development Benefits

- **Environment Consistency:** Eliminates "works on my machine" problems
- **Quick Setup:** New developers can start quickly with Docker
- **Dependency Management:** All dependencies are contained within the image

Production Benefits

- **Deployment Simplicity:** Single container deployment process
- **Resource Management:** Controlled resource allocation per container
- **Service Isolation:** Individual services can be updated independently

6. Application Development and APIs

Web Application Structure

Each module includes a production-ready web application:

Flask Applications

- **app.py files:** Main application entry points for each module
- **Web Interface:** HTML templates for user interaction
- **API Endpoints:** RESTful APIs for model predictions
- **Production Ready:** Configured for production deployment

API Design

The applications provide REST APIs for model inference:

API Features

- **Prediction Endpoints:** Accept input data and return predictions
- **JSON Interface:** Standard JSON request/response format
- **Error Handling:** Proper error responses and status codes
- **Input Validation:** Validation of incoming prediction requests

Example API Workflow

1. **Request:** Client sends prediction data to API endpoint
2. **Validation:** Server validates input data format and values
3. **Processing:** Data is preprocessed using trained pipelines
4. **Prediction:** Trained model generates predictions
5. **Response:** Predictions are returned in JSON format

Web Interface

HTML templates provide user-friendly interfaces:

- **Forms:** Input forms for entering prediction parameters
- **Results Display:** Formatted display of prediction results
- **Interactive Elements:** User-friendly interface components

🌴 Hotel Price Predictor 🏨

Your smart travel companion for hotel recommendations

Enter your travel details and let us predict your hotel experience!

✈️ Travel Code	✈️ User Code
<input type="text" value="0"/>	<input type="text" value="0"/>
📅 Number of Days	📅 Total
<input type="text" value="1"/>	<input type="text" value="10.00"/>
💰 Price	
<input type="text" value="1000.00"/>	
<button>🔮 Predict My Stay!</button>	

🏨 Hotel Recommendation

✈️ Recommended Hotel: Hotel A
📍 Location: Florianopolis (SC)
💰 Recommended Price: ₹313.02

✈️ Happy Travels! | Powered by ML & Streamlit

Gender Classification Model

Username:

Charlotte Johnson

Usercode:

Enter the user id of traveller

Traveller_Age:

Enter the age of traveller

Company name:

4You

Predict

Flight Price Predictor

CONNECTINGDeploy

Your smart travel companion for the best flight deals

Enter your travel details and let us predict your journey cost!

From: Boarding City

To: Destination City

Aracaju

Campo_Grande

Class: Flight Class

Agency

economic

CloudFy

Week Number

Week Day

7

5

Day of Month

5

Predict My Flight Price!

Estimated Flight Price

Rs. 579.41

Bon voyage! Plan your trip with confidence.

7. Deployment and Scaling

Kubernetes Deployment

The project includes Kubernetes configurations for scalable deployment:

Deployment Manifests

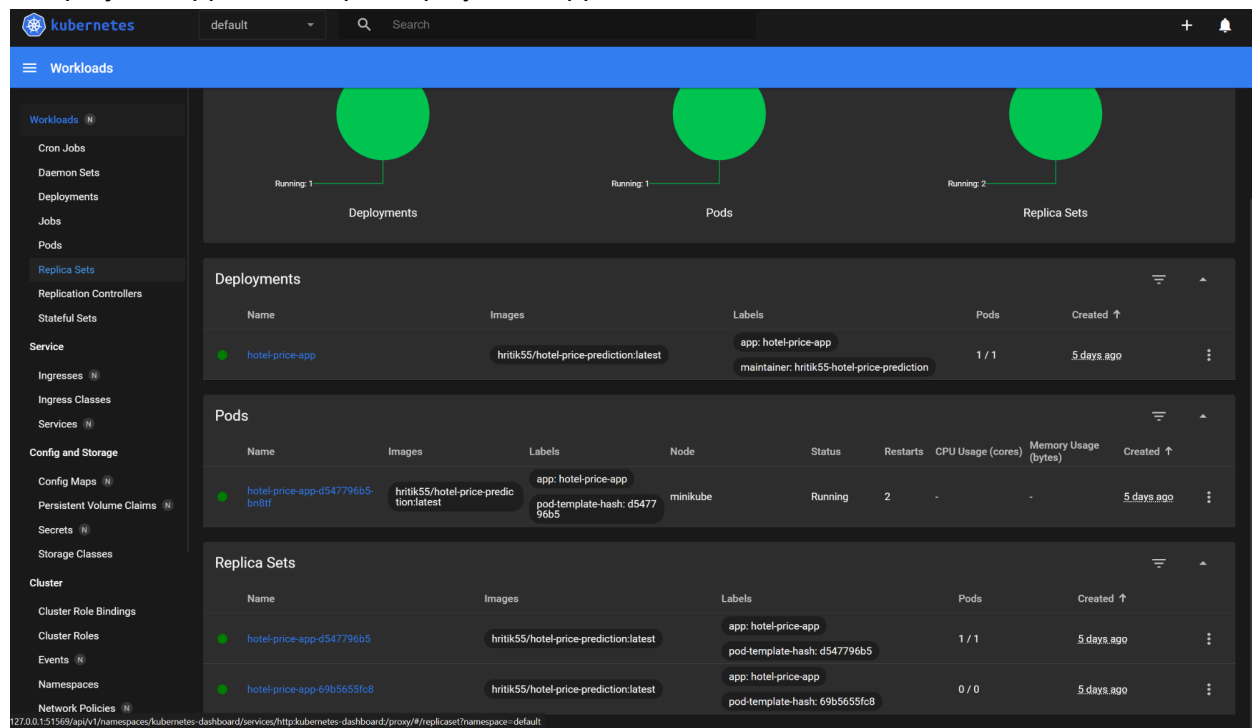
- **deployment.yaml** (flight_price): Kubernetes deployment for flight price service
- **hotel-price-app-deployment.yaml** (hotel_price): Kubernetes deployment for hotel price service

Kubernetes Benefits

- **Auto-scaling:** Automatic scaling based on demand
- **High Availability:** Multiple replicas ensure service availability
- **Load Balancing:** Traffic is distributed across multiple instances
- **Health Monitoring:** Automatic health checks and restarts

Deployment Strategies

The project supports multiple deployment approaches:



Docker Deployment

- **Single Container:** Deploy individual services as Docker containers
- **Docker Compose:** Multi-container applications with dependencies
- **Port Management:** Configurable port exposure for services

Production Deployment

- **Environment Variables:** Configuration through environment variables
- **Health Checks:** Application health monitoring endpoints
- **Logging:** Comprehensive logging for production monitoring

Scaling Considerations

- **Horizontal Scaling:** Multiple instances of the same service
- **Resource Management:** CPU and memory allocation per service
- **Load Distribution:** Traffic routing across multiple instances
- **Database Scaling:** Separate scaling for data storage components

8. Continuous Integration and Continuous Deployment (CI/CD)

Jenkins Integration

The project includes CI/CD automation with Jenkins:

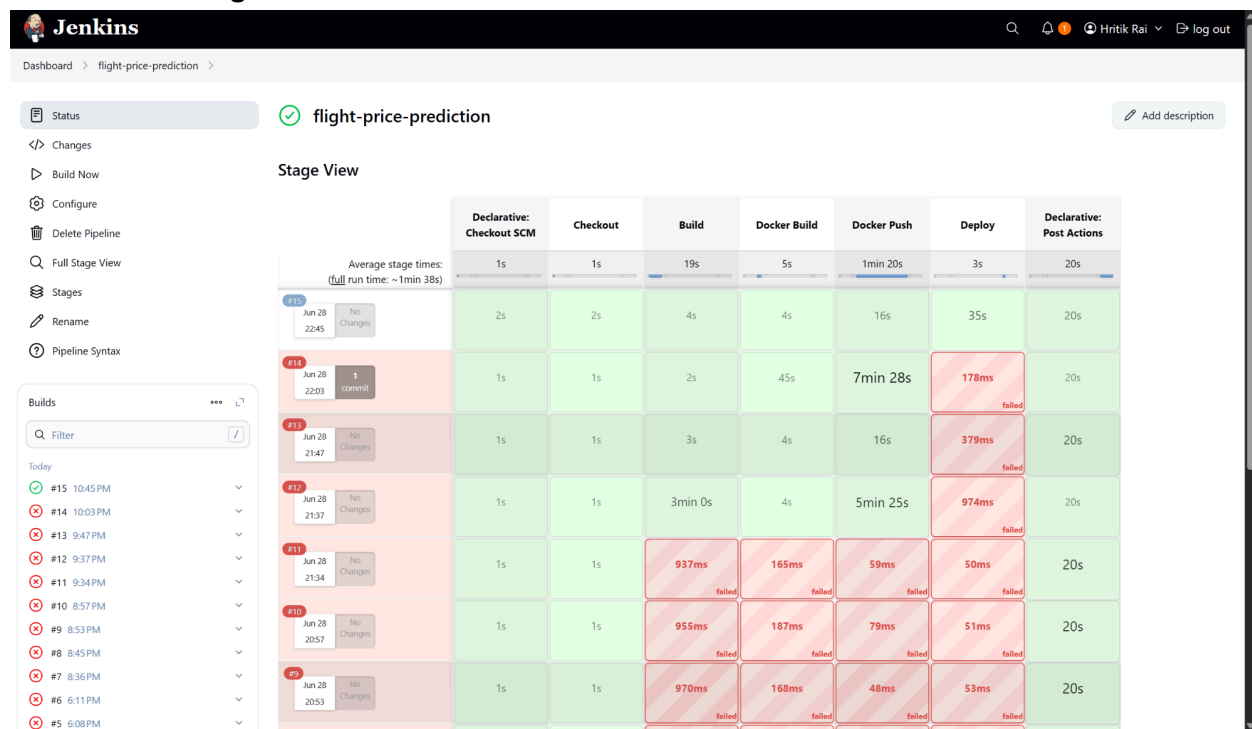
Jenkins Pipeline

- **Jenkinsfile:** Located in `flight_price/Jenkinsfile`
- **Automated Testing:** Run tests on code changes
- **Build Process:** Automatic building of Docker images
- **Deployment:** Automated deployment to staging and production

CI/CD Workflow

The continuous integration and deployment process includes:

Continuous Integration Ste



ps

1. **Code Commit:** Developers commit code changes to repository
2. **Trigger Build:** Jenkins automatically triggers build process
3. **Run Tests:** Automated tests verify code quality and functionality
4. **Build Container:** Docker images are built with new code
5. **Push Images:** Built images are pushed to container registry

Continuous Deployment Steps

1. **Deploy to Staging:** New versions are deployed to staging environment
2. **Integration Testing:** End-to-end tests verify functionality
3. **Production Deployment:** Approved changes are deployed to production
4. **Monitoring:** Production deployment is monitored for issues

Quality Assurance

- **Automated Testing:** Unit tests and integration tests
 - **Code Quality:** Linting and code style checks
 - **Security Scanning:** Container and dependency security checks
 - **Performance Testing:** Load testing for performance validation
-

9. Monitoring and Maintenance

Application Monitoring

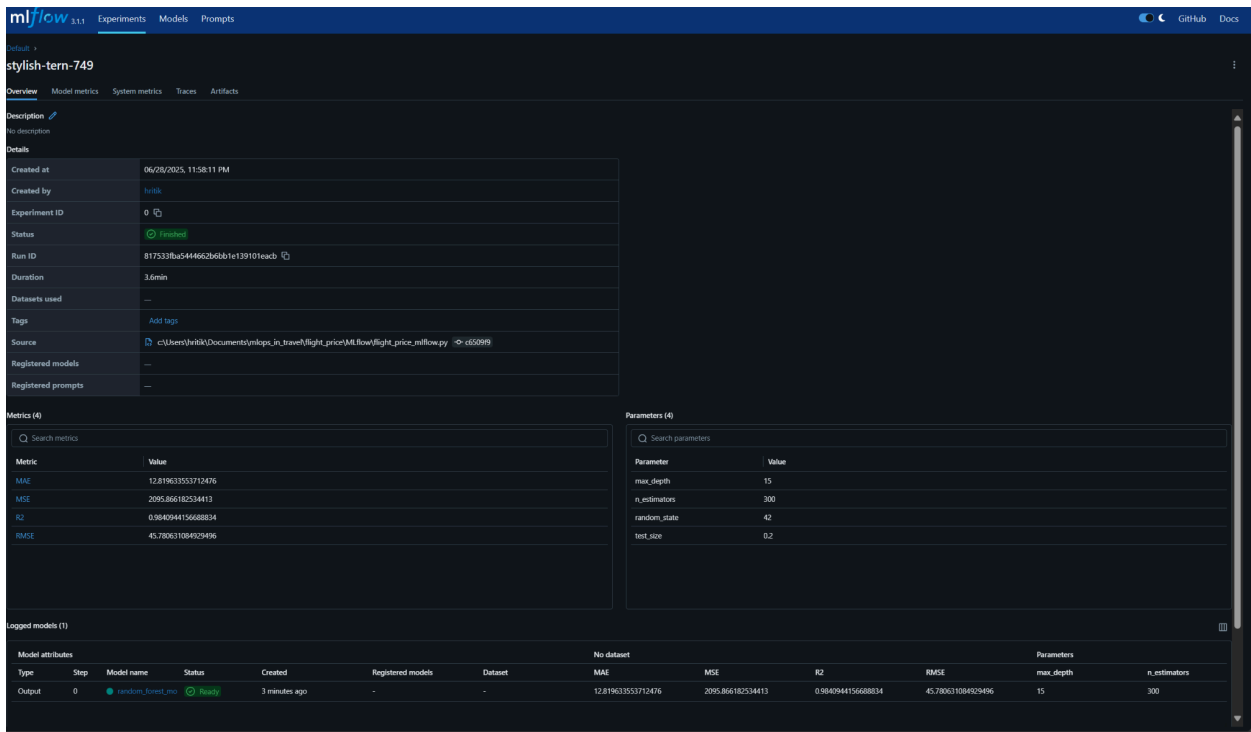
The MLOps workflow includes comprehensive monitoring:

Performance Monitoring

- **Model Performance:** Track prediction accuracy over time
- **API Response Times:** Monitor API latency and throughput
- **Resource Usage:** CPU, memory, and storage utilization
- **Error Rates:** Track application errors and failures

Model Monitoring

Specific monitoring for machine learning models:



Model Drift Detection

- **Data Drift:** Changes in input data distribution
- **Concept Drift:** Changes in the relationship between inputs and outputs
- **Performance Degradation:** Declining model accuracy over time

Retraining Triggers

- **Scheduled Retraining:** Regular model updates on new data
- **Performance-Based:** Retraining when performance drops below threshold
- **Data-Driven:** Retraining when significant data changes are detected

Maintenance Workflows

Regular maintenance ensures system reliability:

Automated Maintenance

- **Log Rotation:** Automatic cleanup of old log files
- **Backup Processes:** Regular backups of models and data
- **Health Checks:** Automated system health monitoring
- **Update Management:** Systematic updates of dependencies

10. Best Practices and Recommendations

Code Organization

The project follows MLOps best practices:

Modular Design

- **Separation of Concerns:** Each module handles specific functionality
- **Reusable Components:** Common functionality is shared across modules
- **Clear Structure:** Consistent directory structure across modules
- **Documentation:** Comprehensive documentation for each component

Security Considerations

- **Environment Variables:** Sensitive configuration stored in .env files
- **Access Control:** Proper authentication and authorization
- **Data Privacy:** Secure handling of user data and predictions
- **Container Security:** Regular security updates for container images

Performance Optimization

- **Model Optimization:** Efficient model architectures and algorithms
- **Caching:** Strategic caching of frequently accessed data
- **Resource Management:** Optimal allocation of computational resources
- **Load Balancing:** Even distribution of workload across instances

Troubleshooting Guidelines

Common issues and their solutions:

Model Performance Issues

- Check data quality and distribution changes
- Verify feature engineering pipeline consistency
- Review hyperparameter settings and model configuration
- Analyze prediction errors and edge cases

Deployment Issues

- Verify container configuration and dependencies
- Check environment variable settings
- Validate network connectivity and port configurations
- Review log files for error messages and stack traces

Conclusion

This MLOps project demonstrates a comprehensive approach to machine learning operations in the travel industry. The workflow covers the complete machine learning lifecycle from data management through deployment and monitoring. Key strengths include:

- **Modularity:** Clear separation of different ML tasks
- **Automation:** Comprehensive automation through Airflow and Jenkins
- **Scalability:** Kubernetes-based deployment for horizontal scaling
- **Monitoring:** MLflow integration for experiment tracking
- **Reproducibility:** Docker containerization ensures consistent environments
- **Best Practices:** Industry-standard MLOps practices throughout

The project serves as an excellent template for implementing production-grade MLOps workflows in any domain, with specific optimizations for travel industry use cases.

Next Steps for Enhancement

- Implement advanced monitoring and alerting systems
- Add automated testing pipelines for model validation
- Enhance security with comprehensive authentication systems
- Expand to additional ML tasks and models
- Implement advanced deployment strategies like blue-green deployments