# MLOps Assignment 1

Heart Disease MLOps Pipeline

Group Number: 80

| BITS ID | Name | GitHub Username |
| --- | --- | --- |
| 2024AB05123 | AISWARYA ANTONY | Aiswarya2020 |
| 2024AB05120 | GAUTHAM KRISHNA M | GauthamKrM |
| 2024AB05121 | GOWRI ANILKUMAR | gowrianair |
| 2024AB05127 | MERIL SARA JOSE | Meril2003 |
| 2024AB05099 | AMULYA K.M | Amulya2025123 |

## 1. SETUP AND INSTALLATION INSTRUCTIONS

### 1.1 Local Development

The local development setup begins with the installation of all required project dependencies. These dependencies are listed in the requirements.txt file and can be installed using the Python package manager.

```
pip install -r requirements.txt
```

After installing the dependencies, the machine learning model is trained. The training process includes downloading the dataset and training the model. The trained model is saved in the models/model directory.

```
python scripts/download_data.py
python src/train.py --n_estimators 100
```

Once the model training is completed, the API can be executed locally. Before starting the API, it must be ensured that the trained model exists in the models/model directory. The API is started using the following command:

```
python src/api.py
```

After execution, the API becomes accessible at http://0.0.0.0:8000.To validate the functionality of the application, automated tests are executed using the PyTest framework. These tests ensure that the individual components of the system function as expected.

```
pytest tests/
```

**1.2 Kubernetes Deployment Using Minikube**

The application can be deployed in a Kubernetes environment using Minikube. Prior to deployment, it is required that both kubectl and minikube are installed and running on the system.The Kubernetes configuration files are applied using the following command:

```
kubectl apply -f k8s/
```

This deployment creates multiple Kubernetes resources. These include the heart-disease-api Deployment, which runs the Docker image of the application, the heart-disease-service, which exposes the API using a LoadBalancer on port 80, and the Prometheus Deployment and Service used for monitoring purposes.

Since Minikube does not assign external IP addresses by default, LoadBalancer services are accessed using the minikube service command. The URLs for accessing the API service, Grafana service, and Prometheus service are obtained using the following commands:

```
minikube service heart-disease-service --url
```

```
minikube service grafana-service --url

minikube service prometheus-service --url
```

**1.3 Monitoring via Prometheus**

The application is instrumented for monitoring using the prometheus-fastapi-instrumentator library. This enables the collection of application-level metrics.

The metrics are exposed through the /metrics endpoint of the application. Prometheus scrapes these metrics based on the configuration defined in the k8s/monitoring.yaml file, enabling continuous monitoring of the application.

## 2. EDA AND MODELLING CHOICES

**2.1 Exploratory Data Analysis (EDA) Summary**

2.1.1 Dataset Quality and Structure
The dataset consists of 297 patient records with 14 clinical features and contains no missing values. All attributes were validated for numeric consistency during preprocessing, ensuring uniform and reliable input data for model training and CI/CD executions.

2.1.2 Target Class Distribution
The target variable shows a moderately balanced distribution, with 160 healthy cases (54%) and 137 heart disease cases (46%). This balance reduces the risk of majority-class bias and allows standard evaluation metrics such as accuracy, precision, recall, and ROC-AUC to remain meaningful without requiring additional class balancing techniques.

2.1.3 Statistical Characteristics
The average patient age is 54.5 years, with heart disease cases most concentrated in the 55–65 age range, indicating age as a significant clinical

factor. Certain clinical attributes, particularly cholesterol levels, exhibit wide variability, ranging from 126 mg/dl to 564 mg/dl, suggesting the presence of high-impact outliers that may influence model behavior.

2.1.4 Correlation and Predictive Features
Correlation analysis identifies thal, ca, and oldpeak as the strongest positive predictors of heart disease, with correlation values of 0.53, 0.46, and 0.42 respectively. In contrast, maximum heart rate (thalach) shows a negative correlation of -0.42, indicating that lower peak heart rates are often associated with higher disease risk.

## 2.2 Modeling Choices and Strategy

The modeling strategy is guided by insights from exploratory data analysis and the clinical importance of accurate heart disease prediction. Due to varying feature scales across clinical attributes, feature standardization using StandardScaler is applied for scale-sensitive models such as Logistic Regression and SVM. Categorical variables that are numerically encoded but represent distinct categories, such as chest pain type, are one-hot encoded to avoid unintended ordinal relationships.

Logistic Regression is selected as the baseline model due to its interpretability and strong linear relationships observed in the data. To capture non-linear interactions among clinical features, ensemble models such as Random Forest and Gradient Boosting are considered as advanced alternatives.Given the medical context, recall (sensitivity) is prioritized to minimize false negatives. Model robustness and generalizability are ensured using k-fold cross-validation, particularly due to the limited dataset size.

## 3.EXPERIMENT TRACKING SUMMARY
This project implements remote experiment tracking using MLflow hosted on DagsHub, rather than local tracking, to ensure persistence across CI/CD executions. The MLflow Tracking URI, username, and access token obtained

from the DagsHub repository are securely stored as GitHub Secrets, enabling automated logging during pipeline runs without exposing credentials.

### 3.1 Trigger Points for Experiment Logging

Experiment tracking is automatically triggered through GitHub Actions on both pull requests and direct pushes to the main branch. During pull requests, validation checks are executed before allowing merge, ensuring that only verified code and model executions proceed. After merging into main, the workflow is re-triggered to execute the full pipeline, including dataset download, model training, experiment logging, Docker image creation, deployment on Minikube, and monitoring. This design ensures that experiment logging occurs during both pre-merge validation and post-merge execution stages.

### 3.2 MLflow Logging Behavior

As part of the CI/CD pipeline, a lightweight model training smoke test is executed in the build-and-test job, where the dataset is downloaded and the training script runs with a minimal estimator configuration. This run is logged to the remote MLflow server using credentials supplied through GitHub Secrets. After merge, a second job performs full model training, which is also logged remotely. As a result, both smoke test runs and final production training runs are recorded with their parameters, metrics, model artifacts, and generated plots.

### 3.3 Artifact Storage and Model Lineage

All experiment outputs are stored persistently in the MLflow artifact store hosted on DagsHub. These artifacts include trained model binaries, evaluation plots, execution logs, and intermediate files generated during training. This artifact storage establishes a complete model lineage, ensuring that every deployed model version can be traced back to its corresponding training run, configuration, and performance metrics.

**3.4 Dashboard Visualization**

Experiment history and performance metrics are visualized using the MLflow dashboard provided by DagsHub. Each pipeline execution automatically appears in the dashboard, displaying run statistics, parameters, metrics, and associated artifacts. Since both pull request smoke tests and post-merge training runs are logged, the dashboard always reflects the most recent execution, enabling easy comparison and reproducibility of experiments.

**3.5 Validation and Deployment Continuity**

Before deployment, the pipeline enforces automated validation through flake8 linting, pytest-based unit testing, and model training smoke tests. Only after all checks pass is a pull request merged into main. Following successful experiment logging, the pipeline proceeds with deployment automation by building and pushing a Docker image to Docker Hub, deploying the API on Minikube, restarting Kubernetes deployments, and serving predictions through running pods. Monitoring metrics are visualized using Grafana, ensuring that only validated and traceable models are deployed.

## 4. UNIT TEST SUMMARY

The project employs pytest as its unit testing framework, with all test cases organized under the tests/ directory. The testing suite validates both model training logic and API service reliability, ensuring correctness before deployment. Dependencies required for executing tests are installed from requirements.txt, using the command pip install -r requirements.txt. Local test execution is supported via pytest tests/ or python3 -m pytest tests/ -v from the project root.

The test suite is structured to maintain modularity and reusability. Shared fixtures such as sample datasets and mock models are defined in conftest.py. Core pipeline functionality is verified through test_train.py, which tests data loading, model evaluation metrics, and the overall training flow in

train.py. Service layer validation is handled by test_api.py, which tests API endpoints and error handling logic within api.py.
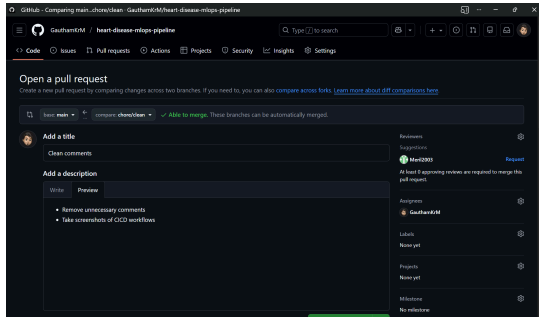
To enforce continuous validation, the testing framework is integrated into the GitHub Actions CI/CD workflow (mlops-pipeline.yml). Tests are automatically executed on every push and pull request targeting the main branch, preventing faulty code or models from progressing to deployment stages. This setup ensures automated regression testing, pipeline stability, and production readiness, following standard MLOps testing best practices.

## **5.ARCHITECTURE DIAGRAM**

# 6.CI/CD AND DEPLOYMENT WORKFLOW SCREENSHOT

## 6.1 Raising a PR:



## 6.2 Checks running on PR - before merge: checks success:



## 6.3  Actions triggered on merging the PR:



## 6.4 Mlflow - experiment tracking dashboard:

## 6.5  Running minikube:
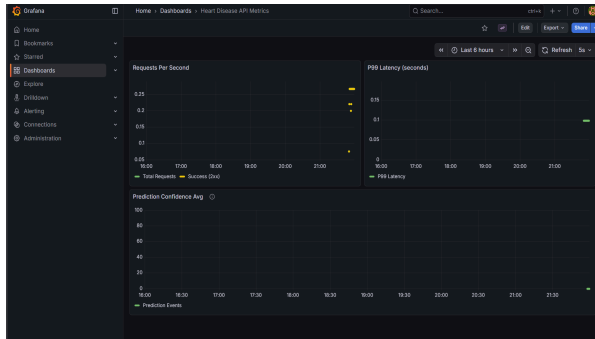


## 6.6  Restarting deployments for heart disease api:





## 6.7 API demo:

## 6.8 Grafana:



## 6.9  API image in dockerhub pushed from pipeline:



## 7.IMPORTANT LINKS

**GitHub Repository Link: *github-repo***

**Workflow Screenshots: *screenshots***

**CICD Workflow Screen Recording: *recording-link***