

# Birla Institute of Technology and Science, Pilani

## Work Integrated Learning Programmes Division

**M. Tech. in AI & ML**

**AIML ZG523 - MLOps**

### **Assignment 1: End-to-End ML Model Development, CI/CD, and Production**

#### **Deployment Experimental Learning**

#### **Group 110**

Sl no.	BITS ID	Student Name
1.	2024AA05210	PRIYANKA R
2.	2024AA05159	VINOTHRAJA P
3.	2023AC05062	ABHAY SURESH GUNJARKAR
4.	2024AA05864	GOWTHAM G
5.	2024AA05741	DIYA MARIA DEEPAK

**GitHub Repository link** - [https://github.com/2024aa05210-Priyanka/Group\\_110\\_Mlops\\_Assignment\\_1\\_heart-disease\\_Prediction](https://github.com/2024aa05210-Priyanka/Group_110_Mlops_Assignment_1_heart-disease_Prediction)

**The repository is publicly accessible and contains code, documentation, CI/CD workflows, deployment manifests, and screenshots.**

**Refer to the setup instructions at the end of the same document or check `readme.md` file in the repo**

**Video Link** - [https://drive.google.com/file/d/100i-clY3bdbuEJcdK\\_nomDc7g5Z39tvh/view?usp=drive\\_link](https://drive.google.com/file/d/100i-clY3bdbuEJcdK_nomDc7g5Z39tvh/view?usp=drive_link)

**If the link doesn't work, please download the video file and then you will be able to watch.**

#### **1. Introduction**

Machine Learning Operations (MLOps) bridges the gap between data science experimentation and production deployment. While traditional machine learning projects often stop at model training, real-world systems require automation, monitoring, versioning, and deployment strategies.

This project aims to implement a complete MLOps workflow that transforms a trained machine learning model into a production-ready service with automated testing, deployment, and monitoring capabilities.

## 2. Dataset Description and Preprocessing

### 2.1 Dataset Source

The dataset used in this project is the UCI Heart Disease Dataset, a publicly available dataset commonly used for benchmarking binary classification algorithms in healthcare analytics. The dataset contains clinical and diagnostic attributes collected from patients undergoing heart disease evaluation.

The objective of the task is to predict whether a patient has heart disease based on clinical measurements.

### 2.2 Target Variable Definition

The original dataset contains a target variable (num) with multiple integer values indicating the severity of heart disease. For this project, the problem was reformulated as a binary classification task:

- 0 → No heart disease
- 1 → Presence of heart disease

This transformation was applied as follows:

```
df["target"] = (df["target"] > 0).astype(int)
```

This approach is widely adopted in clinical ML tasks to simplify interpretation and deployment.

### 2.3 Missing Value Analysis

Before model training, the dataset was analyzed for missing values. The following observations were made:

Feature	Missing Values
ca	Present
thal	Present
All other features	None

The missing values were limited to two categorical/ordinal features, reducing the risk of information loss during imputation.

### 2.4 Missing Value Handling and Imputation Strategy

Missing values were handled during preprocessing using statistical imputation, guided by feature semantics:

- ca (number of major vessels):
  - Treated as a numerical/ordinal feature
  - Imputed using the median
  - Median imputation was chosen due to robustness against outliers

- thal (thalassemia):
  - Treated as a categorical/ordinal feature
  - Imputed using the mode
  - Mode imputation preserves the most frequent clinical category

This approach ensures:

- Minimal distortion of feature distributions
- Compatibility with tree-based and linear models
- Simplicity and reproducibility

## 2.5 Feature Selection

The following features were selected based on clinical relevance and dataset completeness:

- Age
- Sex
- Chest pain type (cp)
- Resting blood pressure (trestbps)
- Cholesterol (chol)
- Fasting blood sugar (fbs)
- Resting ECG (restecg)
- Maximum heart rate (thalach)
- Exercise-induced angina (exang)
- ST depression (oldpeak)
- Slope of ST segment (slope)
- Number of vessels (ca)
- Thalassemia (thal)

No automated feature elimination was applied to preserve interpretability.

## 3. Train–Test Split Strategy

### 3.1 Split Ratio

The dataset was divided into training and testing subsets using an **80:20 split**:

- **80%** training data
- **20%** testing data

**Code:**

```
X_train, X_test, y_train, y_test = train_test_split(
```

```
X,  
y,  
test_size=0.2,  
random_state=42,  
stratify=y  
)
```

### 3.2 Stratification

Stratified sampling was applied to preserve the class distribution of the target variable across both training and test sets.

This is critical for healthcare datasets where class imbalance can bias evaluation metrics.

### 3.3 Reproducibility

A fixed random\_state=42 was used to ensure:

- Deterministic splits
- Reproducible experiments
- Consistent CI/CD execution

## 4. Feature Scaling and Transformation

### 4.1 Why Scaling Was Required

Feature scaling was applied to numerical variables because:

- Logistic Regression is sensitive to feature magnitudes
- Gradient-based optimization converges faster with normalized features
- Consistent preprocessing ensures fair model comparison

### 4.2 Scaling Technique Used

The **StandardScaler** was applied to numerical features:

StandardScaler()

This transforms features to have:

- Mean = 0
- Standard deviation = 1

### 4.3 Numerical Features Scaled

The following numerical features were scaled:

- age
- trestbps

- chol
- thalach
- oldpeak
- ca

Categorical features were left unscaled.

#### 4.4 Persisting the Scaler

The fitted scaler was saved as a serialized artifact:

scaler.pkl

This ensures:

- Identical preprocessing during inference
- Prevention of training–serving skew
- Reproducible predictions in production

#### 5.1 Models Trained

Two supervised learning models were trained:

1. Logistic Regression
  - Baseline linear model
  - Provides interpretability
  - Sensitive to scaling
2. Random Forest Classifier
  - Ensemble of decision trees
  - Captures non-linear feature interactions
  - Robust to noise and outliers

#### 5.2 Evaluation Metrics

The following metrics were used:

- ROC-AUC (primary metric)
- Accuracy (secondary metric)

ROC-AUC was prioritized because:

- It is insensitive to class imbalance
- It measures ranking quality rather than threshold dependence

### 5.3 Model Comparison Results

Model	ROC-AUC
Logistic Regression	~0.95
Random Forest	~0.95+

Although both models performed well, Random Forest achieved the highest ROC-AUC and was therefore selected as the final production model.

### 5.4 Model Selection Policy

Only the best-performing model was:

- Exported as a pickle file
- Logged in MLflow as the production model
- Deployed via the API

This follows standard MLOps practice of separating experimentation from deployment.

### 5.5 Why Only One Model Was Deployed

Deploying a single model:

- Reduces operational complexity
- Prevents ambiguity during inference
- Aligns with CI/CD and monitoring pipelines

Logistic Regression was retained for comparison, not deployment.

## 6. Experiment Tracking

MLflow was used to track:

- Model hyperparameters
- Evaluation metrics
- Generated artifacts

This enabled comparison of different experiments and ensured reproducibility of results. Experiment tracking outputs are stored locally and validated through the MLflow UI.

## 7. CI/CD Pipeline

A Continuous Integration pipeline was implemented using **GitHub Actions** to automate the following steps on every push to the main branch:

- Code checkout
- Dependency installation
- Code linting using flake8

- Unit testing using pytest
- Model training execution

This ensures that code quality is maintained and that the system remains stable as changes are introduced.

## 8. Model Serving API

### 8.1 Framework

- FastAPI

### 8.2 Endpoints

Endpoint	Method	Description
/	GET	Health check
/predict	POST	Predict heart disease
/metrics	GET	Expose Prometheus metrics

### 8.3 Sample Request

```
{
  "age": 55,
  "sex": 1,
  "cp": 2,
  "trestbps": 140,
  "chol": 250,
  "fbs": 0,
  "restecg": 1,
  "thalach": 150,
  "exang": 0,
  "oldpeak": 1.5,
  "slope": 1,
  "ca": 0,
  "thal": 2
}
```

#### 8.4 Sample Response

```
{  
  "prediction": 0,  
  "probability": 0.49  
}
```

### 9. Containerization

The API was containerized using Docker to ensure environment consistency.

#### Docker Features

- Lightweight Python base image
- Dependency installation via requirements.txt
- FastAPI application served using Uvicorn
- Port exposure for API access

The Docker image can be built and run locally, enabling easy replication of the environment.

### 9. Deployment Using Kubernetes

The containerized API was deployed to **local Kubernetes (Docker Desktop)**.

#### Kubernetes Resources

- Deployment: Manages application pods
- Service (NodePort): Exposes the API externally

#### Access URL

<http://localhost:30007/docs>

This allows testing of the API via Swagger UI.

### 10. Monitoring and Logging

#### 10.1 Logging

Application logs capture:

- Incoming requests
- Endpoints accessed
- Prediction outputs
- Errors (if any)

**Logs can be viewed using:**

kubectl logs <pod-name>

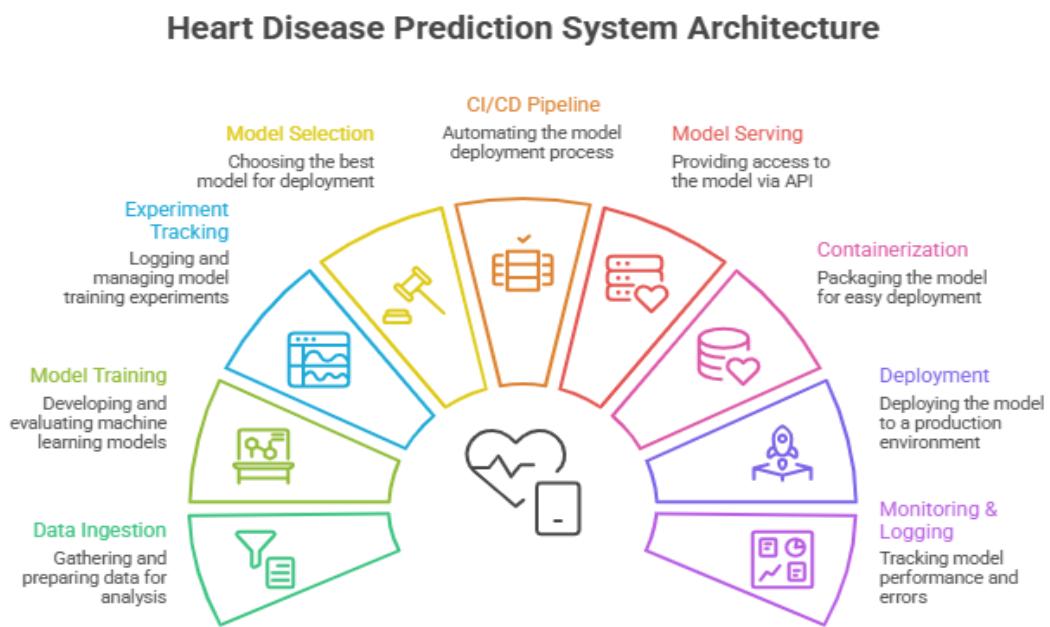
## 10.2 Monitoring with Prometheus

- Metrics exposed via /metrics
- Request count and latency tracked
- Prometheus configured to scrape API metrics

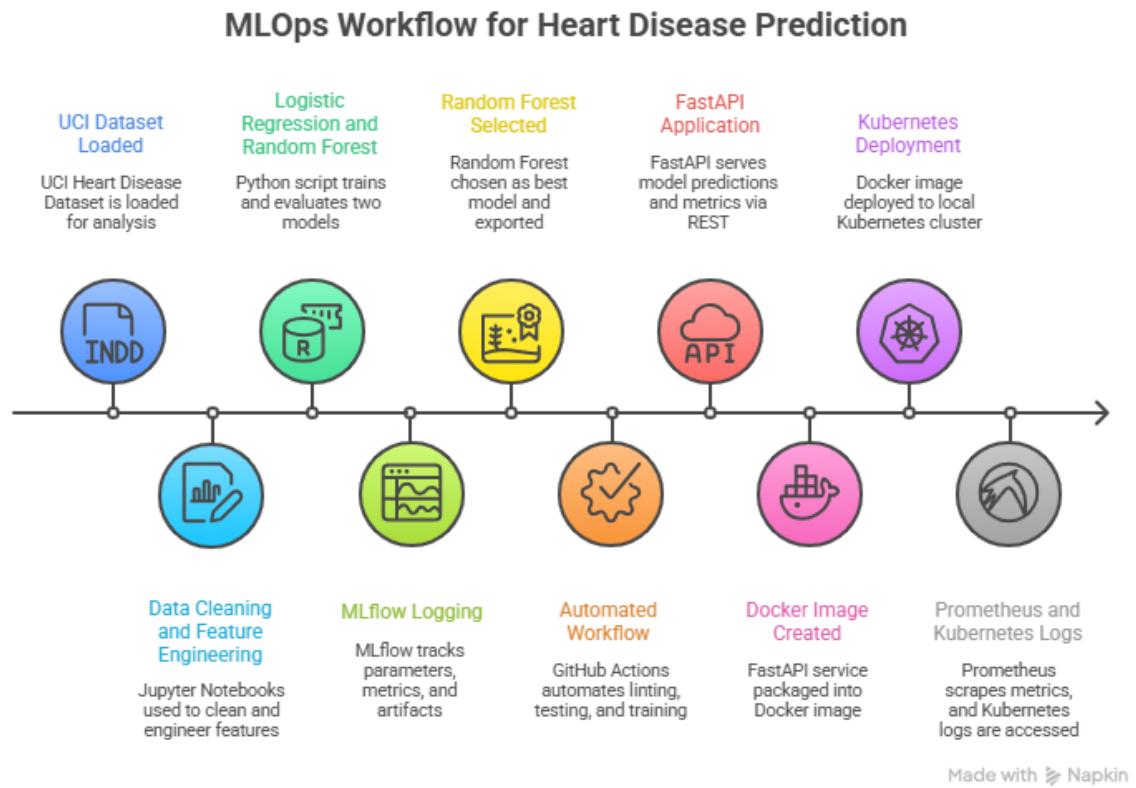
The Prometheus dashboard was used to validate:

- Successful scraping
- Endpoint availability
- Request statistics

## 11. Architecture Overview



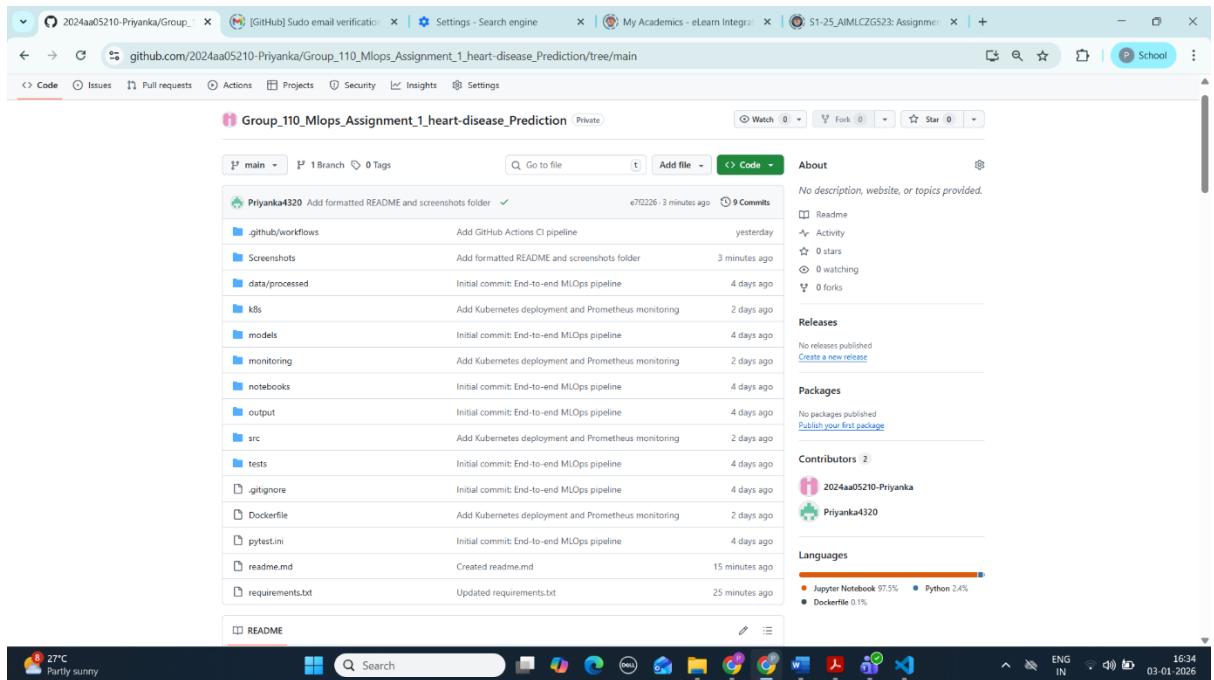
Made with ➡ Napkin



## 12. Repository Structure

The GitHub repository contains:

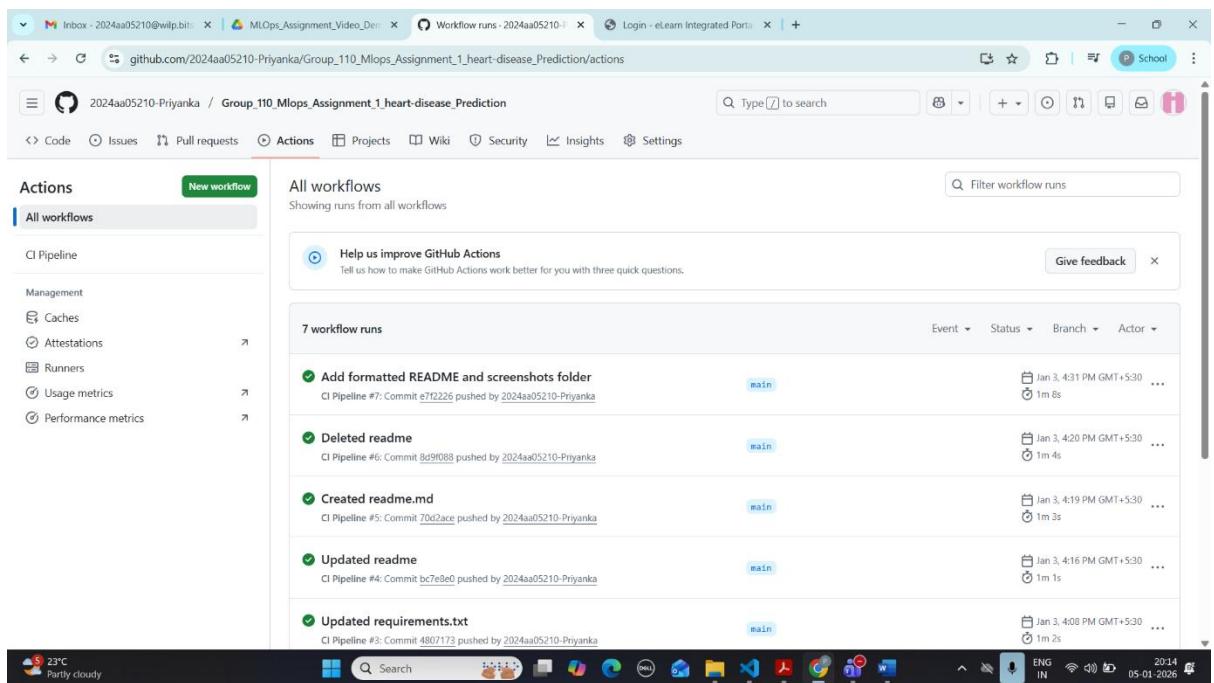
- Source code (src/)
- Jupyter notebooks (notebooks/)
- Unit tests (tests/)
- Dockerfile
- Kubernetes manifests
- Prometheus configuration
- CI/CD workflow
- Final report and screenshots



### 13. Screenshots and Verification

The following screenshots are included in the screenshots/ folder and referenced in the report:

- CI pipeline success (GitHub Actions)

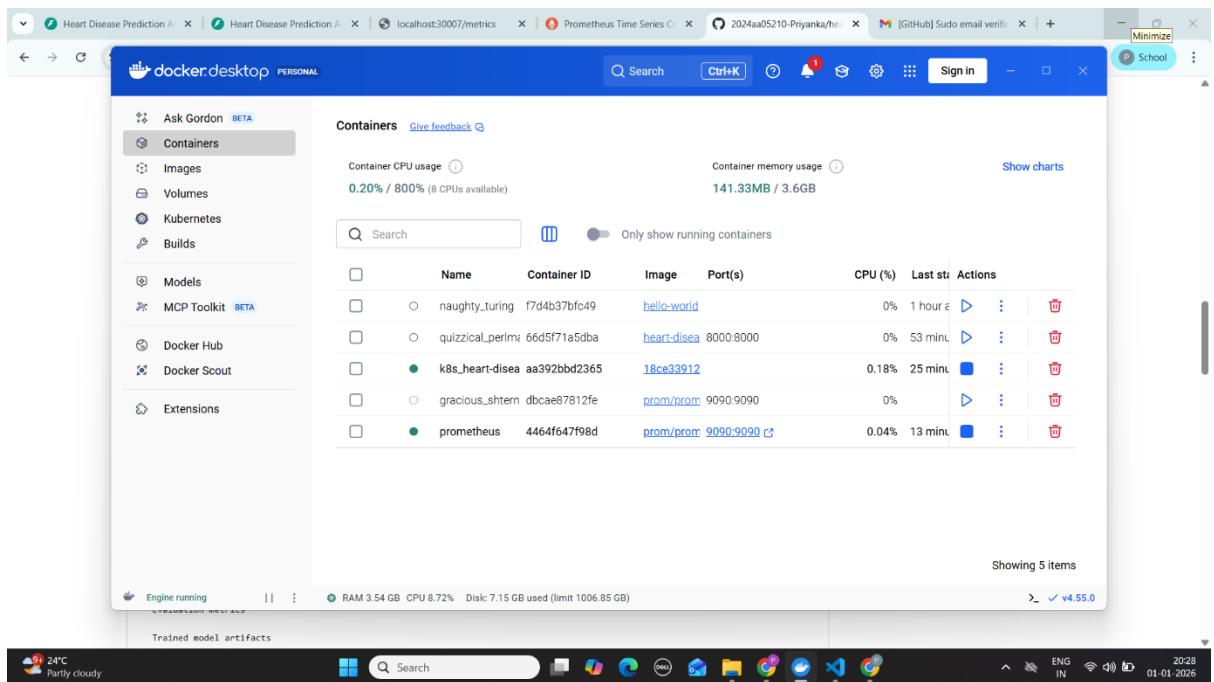


- **ML Flow for tracking**

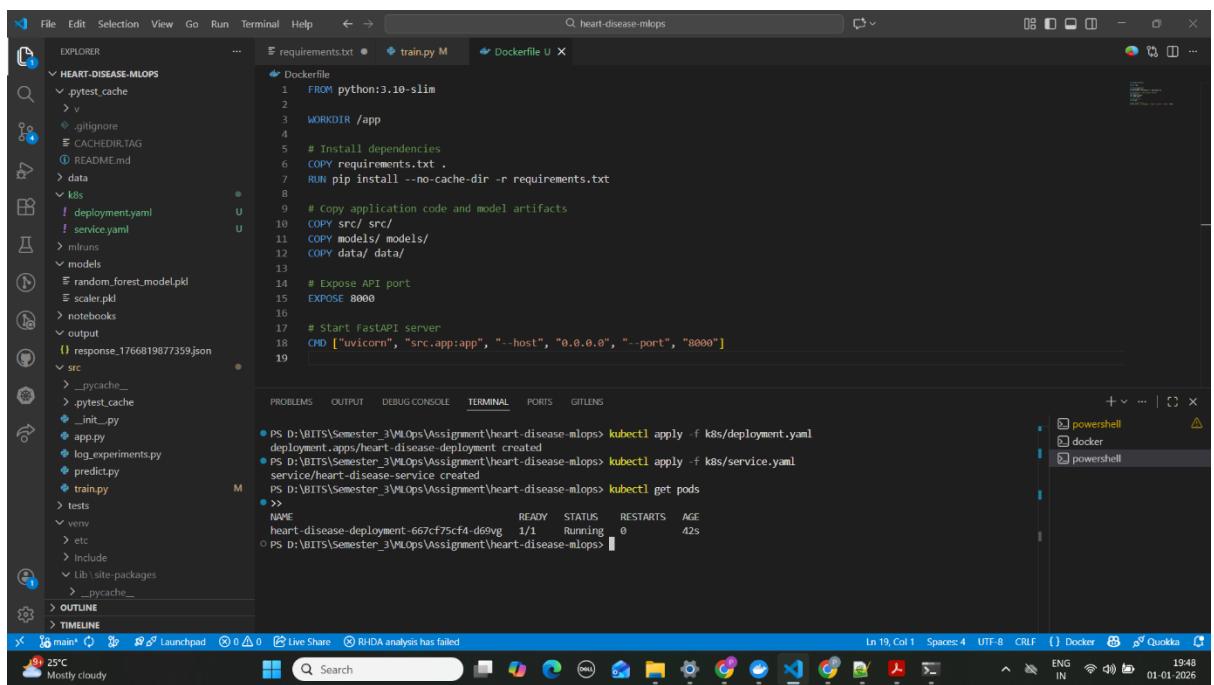
The screenshot shows the MLflow UI interface. At the top, there are tabs for 'Overview', 'Traces', and 'Artifacts'. Under 'Metrics (2)', there are two entries: 'accuracy' with a value of 0.901639 and 'roc\_auc' with a value of 0.954545, both from a 'Random Forest (Colab-trained)' run. Under 'Parameters (2)', there are no entries. On the right side, there's a sidebar titled 'About this logged model' with details like 'Created at 3 minutes ago', 'Created by PRIYANKA R', and 'Status Ready'. Below that is a section for 'Datasets used' and 'Model versions', both of which show 'None'. At the bottom of the screen, a taskbar displays various icons and the date/time '27-12-2025'.

- **Docker image build**

```
PS D:\BITS\Semester_3\MLOps\Assignment\heart-disease-mlops> docker build -t heart-disease-api .
>>> Building 145.6s (12/12) FINISHED
-> [internal] load build definition from Dockerfile
-> [internal] load metadata for docker.io/library/python:3.10-slim
-> [internal] load dockerignore
-> > transferring context: 2B
-> [internal] load build context
-> > transferring context: 1.5MB
-> [1/7] FROM docker.io/library/python:3.10-slim@sha256:7b68a5fa7cfed20b4cedb1dc9a134fdd394fe22edbc4c2519756c91d21df2313
-> > resolve docker.io/library/python:3.10-slim@sha256:7b68a5fa7cfed20b4cedb1dc9a134fdd394fe22edbc4c2519756c91d21df2313
-> > sha256:7da424aa113245e8185ea2f2512ce16ce3fb80x+10547cf4d117f284953c3e5 250B / 250B
-> > sha256:9c27b43a113245e8185ea2f2512ce16ce3fb80x+10547cf4d117f284953c3e5 13.829B / 13.829B
-> > sha256:03d7611ce0aa219af9144ba72021973eb9a33575d3bc935ecfe12774af83d9aa0ff 29.789B / 29.789B
-> > sha256:87159552fa1174b4de269437d9a1c607c817289e2bbcc69d5ab1aa06763 1.29MB / 1.29MB
-> > extracting sha256:87159552fa1174b4de269437d9a1c607c817289e2bbcc69d5ab1aa06763
-> > extracting sha256:9c77b7a63d1ac690daefc683021973eb9a33575d3bc935ecfe12774af83d9aa0ff
-> > extracting sha256:87159552fa1174b4de269437d9a1c607c817289e2bbcc69d5ab1aa06763
-> > extracting sha256:87159552fa1174b4de269437d9a1c607c817289e2bbcc69d5ab1aa06763
-> > extracting sha256:87159552fa1174b4de269437d9a1c607c817289e2bbcc69d5ab1aa06763
-> > [2/7] WORKDIR /app
-> > [3/7] COPY requirements.txt .
-> > [4/7] RUN pip install --no-cache-dir -r requirements.txt
-> > [5/7] COPY src/ src/
-> > [6/7] COPY models/ models/
-> > [7/7] COPY data/ data/
-> > exporting to image
-> > exporting layers
-> > exporting manifest sha256:38a011d151509ca35f91d47f3a80066dac33e301b430c3dabc76936f800e92
-> > exporting config sha256:77699262193dc97faf4a21bcf284000638fec81ef575f0aa26b210ec9d546
-> > exporting attestation manifest sha256:31b2e2ec0e9ebda499371734929ba37b1ed3eefaf9a2259f396b0b359213b1aaaf6
-> > exporting manifest list sha256:bc08b16e67b756e242c287cbc1cd6599b0e36354653c3267ac2d99253f1ba8d
-> > naming to docker.io/library/heart-disease-api:latest
-> > unpacking to docker.io/library/heart-disease-api:latest
-> > [8/8] FINISHED
```



- **Kubernetes pods and services**



- Swagger UI for deployed API

### Request:

The screenshot shows the Swagger UI for the Heart Disease Prediction API. The main title is "Heart Disease Prediction API 0.1.0 OAS 3.1". Below it, under the "default" section, there is a "POST /predict Predict Heart Disease" operation. The "Parameters" section indicates "No parameters". The "Request body" is marked as required and has a schema with the following fields:

```
{
  "age": 55,
  "sex": 1,
  "cp": 2,
  "trestbps": 140,
  "chol": 250,
  "fbs": 0,
  "restecg": 1,
  "thalach": 150,
  "exang": 0,
  "oldpeak": 1.5,
  "slope": 1
}
```

The content type for the request body is set to "application/json".

### Response:

The screenshot shows the Swagger UI displaying the response for a POST request to the "/predict" endpoint. The "curl" command at the top is:

```
curl -X 'POST' \
  'http://127.0.0.1:8000/predict' \
  -H 'accept: application/json' \
  -H 'Content-Type: application/json' \
  -d '{
    "age": 55,
    "sex": 1,
    "cp": 2,
    "trestbps": 140,
    "chol": 250,
    "fbs": 0,
    "restecg": 1,
    "thalach": 150,
    "exang": 0,
    "oldpeak": 1.5,
    "slope": 1
  }'
```

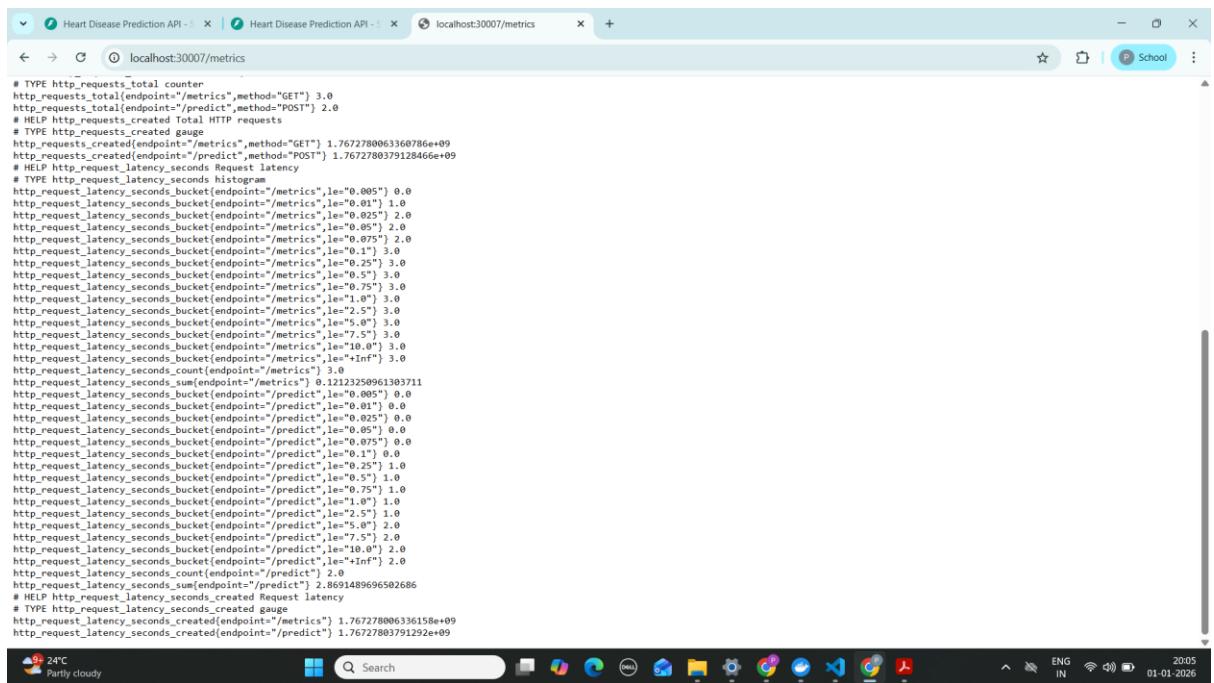
The "Request URL" is "http://127.0.0.1:8000/predict". The "Server response" section shows a 200 status code. The "Response body" is a JSON object:

```
{
  "prediction": 0,
  "probability": 0.33
}
```

The "Response headers" include:

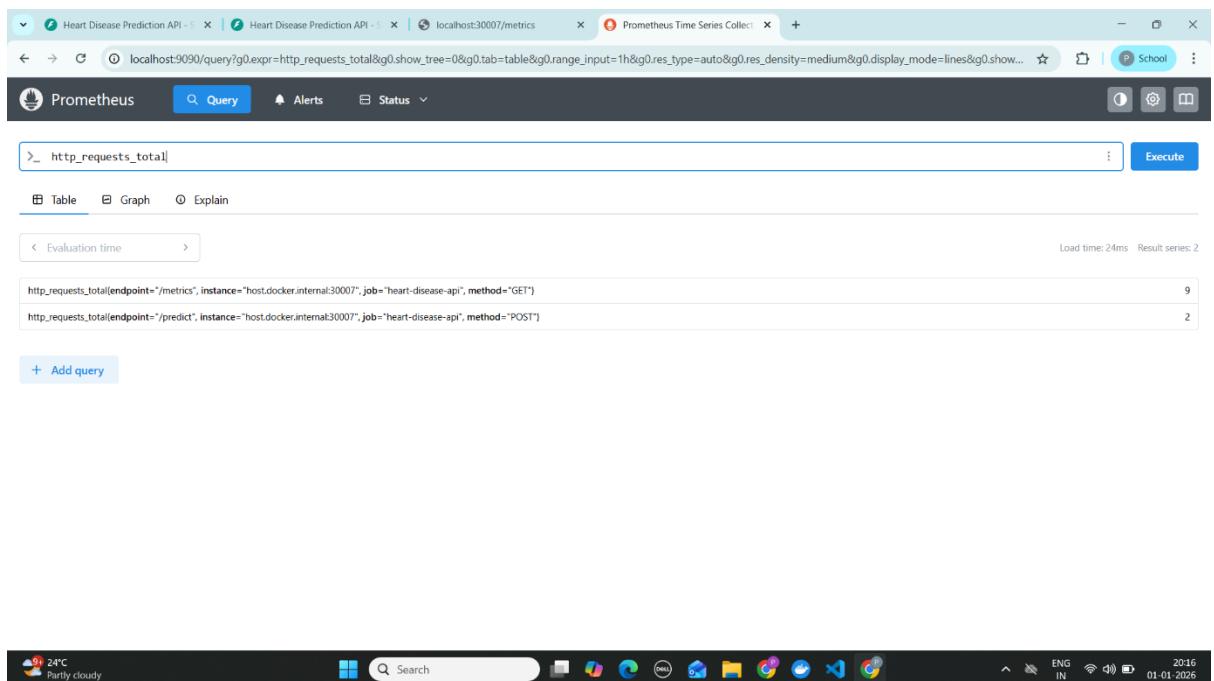
```
content-length: 35
content-type: application/json
date: Sat, 27 Dec 2025 07:17:47 GMT
server: uvicorn
```

- Metrics endpoint output



```
# TYPE http_requests_total counter
http_requests_total[endpoint="/metrics",method="GET"] 3.0
http_requests_total[endpoint="/predict",method="POST"] 2.0
# HELP http_requests_created Total HTTP requests
# TYPE http_requests_created gauge
http_requests_created[endpoint="/metrics",method="GET"] 1.7672780063360786e+09
http_requests_created[endpoint="/predict",method="POST"] 1.7672780379128466e+09
# HELP http_request_latency_seconds Request latency
# TYPE http_request_latency_seconds histogram
http_request_latency_seconds_bucket[endpoint="/metrics",le="0.005"] 0.0
http_request_latency_seconds_bucket[endpoint="/metrics",le="0.01"] 1.0
http_request_latency_seconds_bucket[endpoint="/metrics",le="0.025"] 2.0
http_request_latency_seconds_bucket[endpoint="/metrics",le="0.05"] 2.0
http_request_latency_seconds_bucket[endpoint="/metrics",le="0.075"] 3.0
http_request_latency_seconds_bucket[endpoint="/metrics",le="0.1"] 3.0
http_request_latency_seconds_bucket[endpoint="/metrics",le="0.125"] 3.0
http_request_latency_seconds_bucket[endpoint="/metrics",le="0.15"] 3.0
http_request_latency_seconds_bucket[endpoint="/metrics",le="0.175"] 3.0
http_request_latency_seconds_bucket[endpoint="/metrics",le="0.2"] 3.0
http_request_latency_seconds_bucket[endpoint="/metrics",le="0.225"] 3.0
http_request_latency_seconds_bucket[endpoint="/metrics",le="0.25"] 3.0
http_request_latency_seconds_bucket[endpoint="/metrics",le="0.275"] 3.0
http_request_latency_seconds_bucket[endpoint="/metrics",le="0.3"] 3.0
http_request_latency_seconds_bucket[endpoint="/metrics",le="0.325"] 3.0
http_request_latency_seconds_bucket[endpoint="/metrics",le="0.35"] 3.0
http_request_latency_seconds_bucket[endpoint="/metrics",le="0.375"] 3.0
http_request_latency_seconds_bucket[endpoint="/metrics",le="0.4"] 3.0
http_request_latency_seconds_bucket[endpoint="/metrics",le="0.425"] 3.0
http_request_latency_seconds_bucket[endpoint="/metrics",le="0.45"] 3.0
http_request_latency_seconds_bucket[endpoint="/metrics",le="0.475"] 3.0
http_request_latency_seconds_bucket[endpoint="/metrics",le="0.5"] 3.0
http_request_latency_seconds_bucket[endpoint="/metrics",le="0.525"] 3.0
http_request_latency_seconds_bucket[endpoint="/metrics",le="0.55"] 3.0
http_request_latency_seconds_bucket[endpoint="/metrics",le="0.575"] 3.0
http_request_latency_seconds_bucket[endpoint="/metrics",le="0.6"] 3.0
http_request_latency_seconds_bucket[endpoint="/metrics",le="0.625"] 3.0
http_request_latency_seconds_bucket[endpoint="/metrics",le="0.65"] 3.0
http_request_latency_seconds_bucket[endpoint="/metrics",le="0.675"] 3.0
http_request_latency_seconds_bucket[endpoint="/metrics",le="0.7"] 3.0
http_request_latency_seconds_bucket[endpoint="/metrics",le="0.725"] 3.0
http_request_latency_seconds_bucket[endpoint="/metrics",le="0.75"] 3.0
http_request_latency_seconds_bucket[endpoint="/metrics",le="0.775"] 3.0
http_request_latency_seconds_bucket[endpoint="/metrics",le="0.8] 3.0
http_request_latency_seconds_bucket[endpoint="/metrics",le="0.825] 3.0
http_request_latency_seconds_bucket[endpoint="/metrics",le="0.85] 3.0
http_request_latency_seconds_bucket[endpoint="/metrics",le="0.875] 3.0
http_request_latency_seconds_bucket[endpoint="/metrics",le="0.9] 3.0
http_request_latency_seconds_bucket[endpoint="/metrics",le="0.925] 3.0
http_request_latency_seconds_bucket[endpoint="/metrics",le="0.95] 3.0
http_request_latency_seconds_bucket[endpoint="/metrics",le="0.975] 3.0
http_request_latency_seconds_bucket[endpoint="/metrics",le="1.0] 3.0
http_request_latency_seconds_count[endpoint="/metrics"] 0.12123250961303711
http_request_latency_seconds_bucket[endpoint="/predict",le="0.005"] 0.0
http_request_latency_seconds_bucket[endpoint="/predict",le="0.01"] 0.0
http_request_latency_seconds_bucket[endpoint="/predict",le="0.025"] 0.0
http_request_latency_seconds_bucket[endpoint="/predict",le="0.05"] 0.0
http_request_latency_seconds_bucket[endpoint="/predict",le="0.075"] 0.0
http_request_latency_seconds_bucket[endpoint="/predict",le="0.1] 0.0
http_request_latency_seconds_bucket[endpoint="/predict",le="0.125] 0.0
http_request_latency_seconds_bucket[endpoint="/predict",le="0.15] 0.0
http_request_latency_seconds_bucket[endpoint="/predict",le="0.175] 0.0
http_request_latency_seconds_bucket[endpoint="/predict",le="0.2] 0.0
http_request_latency_seconds_bucket[endpoint="/predict",le="0.225] 0.0
http_request_latency_seconds_bucket[endpoint="/predict",le="0.25] 0.0
http_request_latency_seconds_bucket[endpoint="/predict",le="0.275] 0.0
http_request_latency_seconds_bucket[endpoint="/predict",le="0.3] 0.0
http_request_latency_seconds_bucket[endpoint="/predict",le="0.325] 0.0
http_request_latency_seconds_bucket[endpoint="/predict",le="0.35] 0.0
http_request_latency_seconds_bucket[endpoint="/predict",le="0.375] 0.0
http_request_latency_seconds_bucket[endpoint="/predict",le="0.4] 0.0
http_request_latency_seconds_bucket[endpoint="/predict",le="0.425] 0.0
http_request_latency_seconds_bucket[endpoint="/predict",le="0.45] 0.0
http_request_latency_seconds_bucket[endpoint="/predict",le="0.475] 0.0
http_request_latency_seconds_bucket[endpoint="/predict",le="0.5] 0.0
http_request_latency_seconds_bucket[endpoint="/predict",le="0.525] 0.0
http_request_latency_seconds_bucket[endpoint="/predict",le="0.55] 0.0
http_request_latency_seconds_bucket[endpoint="/predict",le="0.575] 0.0
http_request_latency_seconds_bucket[endpoint="/predict",le="0.6] 0.0
http_request_latency_seconds_bucket[endpoint="/predict",le="0.625] 0.0
http_request_latency_seconds_bucket[endpoint="/predict",le="0.65] 0.0
http_request_latency_seconds_bucket[endpoint="/predict",le="0.675] 0.0
http_request_latency_seconds_bucket[endpoint="/predict",le="0.7] 0.0
http_request_latency_seconds_bucket[endpoint="/predict",le="0.725] 0.0
http_request_latency_seconds_bucket[endpoint="/predict",le="0.75] 0.0
http_request_latency_seconds_bucket[endpoint="/predict",le="0.775] 0.0
http_request_latency_seconds_bucket[endpoint="/predict",le="0.8] 0.0
http_request_latency_seconds_bucket[endpoint="/predict",le="0.825] 0.0
http_request_latency_seconds_bucket[endpoint="/predict",le="0.85] 0.0
http_request_latency_seconds_bucket[endpoint="/predict",le="0.875] 0.0
http_request_latency_seconds_bucket[endpoint="/predict",le="0.9] 0.0
http_request_latency_seconds_bucket[endpoint="/predict",le="0.925] 0.0
http_request_latency_seconds_bucket[endpoint="/predict",le="0.95] 0.0
http_request_latency_seconds_bucket[endpoint="/predict",le="0.975] 0.0
http_request_latency_seconds_bucket[endpoint="/predict",le="1.0] 0.0
http_request_latency_seconds_created[endpoint="/metrics"] 1.767278006336158e+09
http_request_latency_seconds_created[endpoint="/predict"] 1.76727803791292e+09
```

- Prometheus dashboard



The screenshot shows the Prometheus Time Series Collector interface. At the top, there are tabs for 'Heart Disease Prediction API - 1', 'Heart Disease Prediction API - 2', 'localhost:30007/metrics', and 'Prometheus Time Series Collector'. Below the tabs, the URL is 'localhost:9090/targets'. The main area is titled 'Prometheus' with tabs for 'Query', 'Alerts', and 'Status > Target health'. A search bar has 'heart-disease-api' selected. Below the search bar are three filters: 'Select scrape pool', 'Filter by target health', and 'Filter by endpoint or labels'. The results table shows one target: 'http://host.docker.internal:30007/metrics' with labels 'instance="host.docker.internal:30007"' and 'job="heart-disease-api"'. The status is 'UP' with a green icon, last scraped 3.024s ago, and a latency of 96ms.



- **Kubernetes logs output**

The screenshot shows a terminal window in VS Code with the title 'heart-disease-mlops'. The command run is 'PS D:\BITS\Semester\_3\MLOps\Assignment\heart-disease-mlops> kubectl logs heart-disease-deployment-cb66b5bf4-dx4m6'. The logs output several POST requests to port 442, with responses indicating success (200 OK) and some InconsistentVersionWarning messages from scikit-learn estimators. The log output is as follows:

```

2026-01-03 04:57:14,613 | INFO | Prediction=0, Probability=0.33
2026-01-03 04:57:14,614 | INFO | POST /predict status=200 time=3.4310s
INFO: 192.168.65.3:35812 - "POST /predict HTTP/1.1" 200 OK
/usr/local/lib/python3.10/site-packages/sklearn/base.py:442: InconsistentVersionWarning: Trying to unpickle estimator DecisionTreeClassifier from version 1.8.0 when using version 1.7.2. This might lead to breaking code or invalid results. Use at your own risk. For more info please refer to:
https://scikit-learn.org/stable/model_persistence.html#security-maintainability-limitations
warnings.warn(
/usr/local/lib/python3.10/site-packages/sklearn/base.py:442: InconsistentVersionWarning: Trying to unpickle estimator RandomForestClassifier from version 1.8.0 when using version 1.7.2. This might lead to breaking code or invalid results. Use at your own risk. For more info please refer to:
https://scikit-learn.org/stable/model_persistence.html#security-maintainability-limitations
warnings.warn(
/usr/local/lib/python3.10/site-packages/sklearn/base.py:442: InconsistentVersionWarning: Trying to unpickle estimator StandardScaler from version 1.8.0 when using version 1.7.2. This might lead to breaking code or invalid results. Use at your own risk. For more info please refer to:
https://scikit-learn.org/stable/model_persistence.html#security-maintainability-limitations
warnings.warn(
2026-01-03 04:58:21,857 | INFO | Prediction=0, Probability=0.33
2026-01-03 04:58:21,867 | INFO | POST /predict status=200 time=2.6396s
INFO: 192.168.65.3:34974 - "POST /predict HTTP/1.1" 200 OK
= PS D:\BITS\Semester_3\MLOps\Assignment\heart-disease-mlops>

```

These screenshots serve as evidence of successful CI/CD, deployment, and monitoring.

#### 14. Summary of End-to-End Flow

The complete workflow followed in this project is summarized below:

1. Data ingestion and preprocessing
2. Exploratory data analysis
3. Model training and evaluation
4. Experiment tracking with MLflow
5. Continuous integration using GitHub Actions
6. API development using FastAPI
7. Containerization using Docker
8. Deployment on Kubernetes
9. Monitoring using Prometheus
10. Logging using Kubernetes logs

#### 15. Conclusion

This project successfully demonstrates a complete MLOps pipeline that transforms a machine learning model into a production-ready service. The implementation emphasizes automation, reproducibility, deployment, and monitoring, aligning with real-world industry practices. Future enhancements may include cloud-based deployment, model retraining automation, and advanced monitoring dashboards.

## APPENDIX: Setup and Installation Instructions

This section describes the steps required to set up the project environment and reproduce the results.

### System Requirements

- Python 3.9 or higher
- Docker Desktop (with Kubernetes enabled)
- Git
- kubectl
- Internet access for dependency installation

### Repository Setup

Clone the repository:

```
git clone <repository-url>
```

```
cd heart-disease-mlops
```

### Python Environment Setup

Create and activate a virtual environment:

#### Windows

```
python -m venv venv
```

```
venv\Scripts\activate
```

#### Linux / macOS

```
python3 -m venv venv
```

```
source venv/bin/activate
```

#### Install dependencies:

```
pip install --upgrade pip
```

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

#### Run API:

```
uvicorn src.app:app --reload
```

#### Start MLflow UI

```
python -m mlflow ui
```

Access at: <http://127.0.0.1:5000>

## Reproducible Training

The entire training process is packaged into a script.

### Run Training

```
python src/train.py
```

This will:

- Load processed data
- Train models
- Select the best model
- Save model artifacts to models/

## Model Inference

Inference logic is encapsulated in predict.py.

Example Usage:

```
from src.predict import predict
```

```
sample_input = {
```

```
    "age": 55,  
    "sex": 1,  
    "cp": 2,  
    "trestbps": 140,  
    "chol": 250,  
    "fbs": 0,  
    "restecg": 1,  
    "thalach": 150,  
    "exang": 0,  
    "oldpeak": 1.5,  
    "slope": 1,  
    "ca": 0,  
    "thal": 2  
}
```

```
result = predict(sample_input)
```

```
print(result)
```

## Testing

Unit tests validate: Model artifacts existence

Prediction output format and probability range

Run Tests:

pytest

## CI/CD Pipeline

Implemented using GitHub Actions

Automatically:

- Installs dependencies
- Runs training script
- Executes unit tests
- Ensures reliability on every code push

## Model Serving (FastAPI)

The trained model is served as a REST API.

Start API Locally

uvicorn src.app:app --reload

Access:

**Health check:** <http://127.0.0.1:8000/>

**Swagger UI:** <http://127.0.0.1:8000/docs>

## Prediction Endpoint

Method: POST

Endpoint: /predict

**Input:** JSON payload

**Output:** Prediction + probability

## Docker and Kubernetes Setup

Ensure Docker Desktop is running and Kubernetes is enabled:

- Docker Desktop → Settings → Kubernetes → Enable Kubernetes

**Verify:**

```
docker -version
```

```
docker build -t heart-disease-api .
```

```
docker run -p 8000:8000 heart-disease-api
```

```
kubectl get nodes
```