

End-to-End MLOps Pipeline for Heart Disease Prediction

Group 27: Contributors

2024ab05143 M ILANGO

2024aa05968 RAJESH MANDADAPU

2024AA05944 P AMARESHWAR RAO

2024ab05056 AZARUDEEN S

2024aa05678 RUCHI KESHARWANI

Date: January 6th, 2026

Code Repository: [<https://github.com/2024ab05143/heart-disease-mlops>]

Reviewer Instructions: Setup & Verification(followed by individual tasks 1-9)

Repository: <https://github.com/2024ab05143/heart-disease-mlops>

This guide provides step-by-step instructions to replicate the full MLOps pipeline on a Windows machine. The workflow is fully containerized; **you do not need Python or Flask installed on your host machine.**

Part 1: Environment Setup (Prerequisites)

Before running the code, please ensure the following "System Tools" are installed.

1. Install Git for Windows

- **Download:** Go to git-scm.com/download/win and download the "64-bit Git for Windows Setup".
- **Installation:** Run the installer. You can safely click "Next" through all default settings.
- **Verify:** Open a new PowerShell terminal and run:

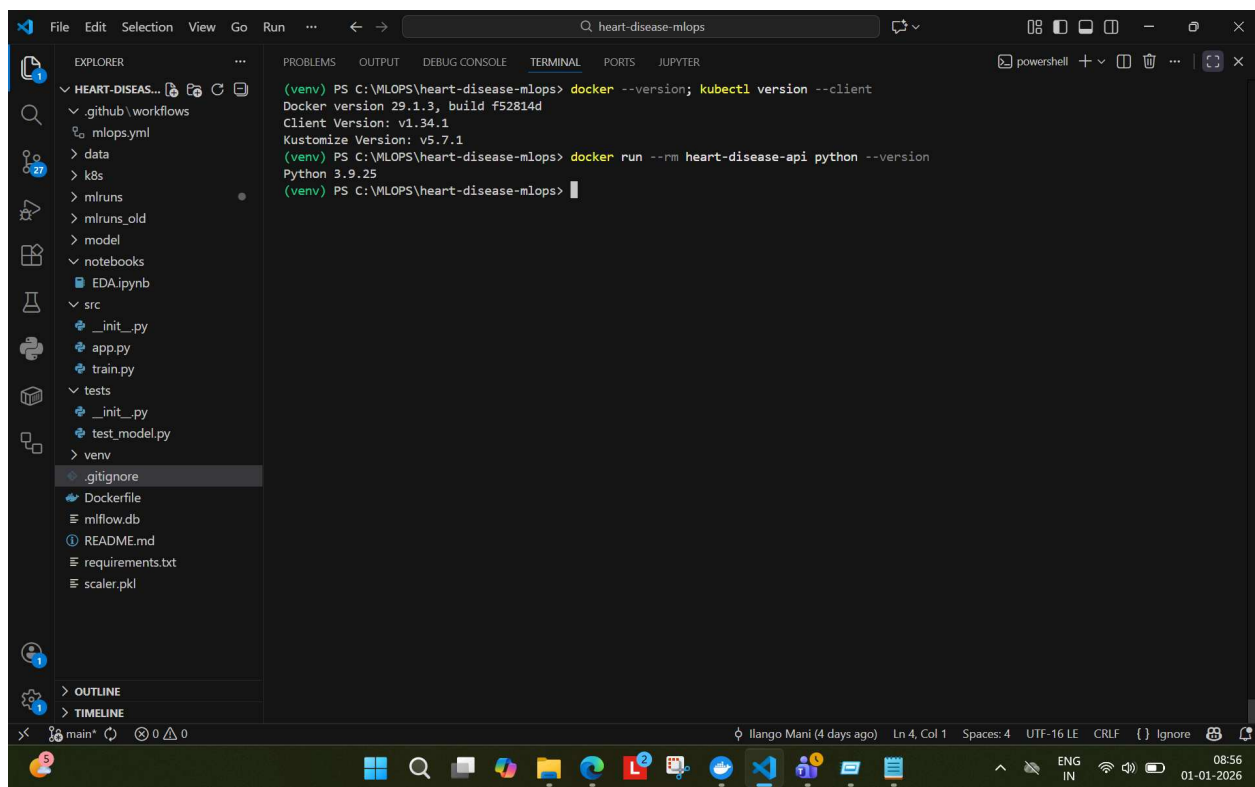
```
PowerShell  
git --version
```

2. Install Docker Desktop

- **Download:** Go to docker.com/products/docker-desktop and download for Windows.
- **Installation:** Run the installer. Ensure "Use WSL 2 instead of Hyper-V" is checked

(recommended for performance).

- **Critical Step (Enable Kubernetes):**
 - Open the Docker Desktop Dashboard.
 - Click the **Settings (Gear Icon)** in the top right.
 - Select **Kubernetes** from the left menu.
 - Check the box **Enable Kubernetes**.
 - Click **Apply & Restart**.
- **Verify:** Open PowerShell and run:
PowerShell
docker --version
kubectl version --client



The screenshot shows a Visual Studio Code interface with a terminal window open. The terminal is running PowerShell commands to verify the installation of Docker and Kubernetes. The output shows Docker version 29.1.3, build f52814d, and Kubernetes version v1.34.1. The terminal also shows the command to run a Python script using Docker.

```
(venv) PS C:\MILOPS\heart-disease-mlops> docker --version; kubectl version --client
Docker version 29.1.3, build f52814d
Client Version: v1.34.1
Kustomize Version: v5.7.1
(venv) PS C:\MILOPS\heart-disease-mlops> docker run --rm heart-disease-api python --version
Python 3.9.25
(venv) PS C:\MILOPS\heart-disease-mlops>
```

Figure 1: Verification of Docker Desktop installation on the host machine.

Part 2: Reproducible Training

Note: This step proves that model training is consistent regardless of your local Python version.

1. Clone the Repository

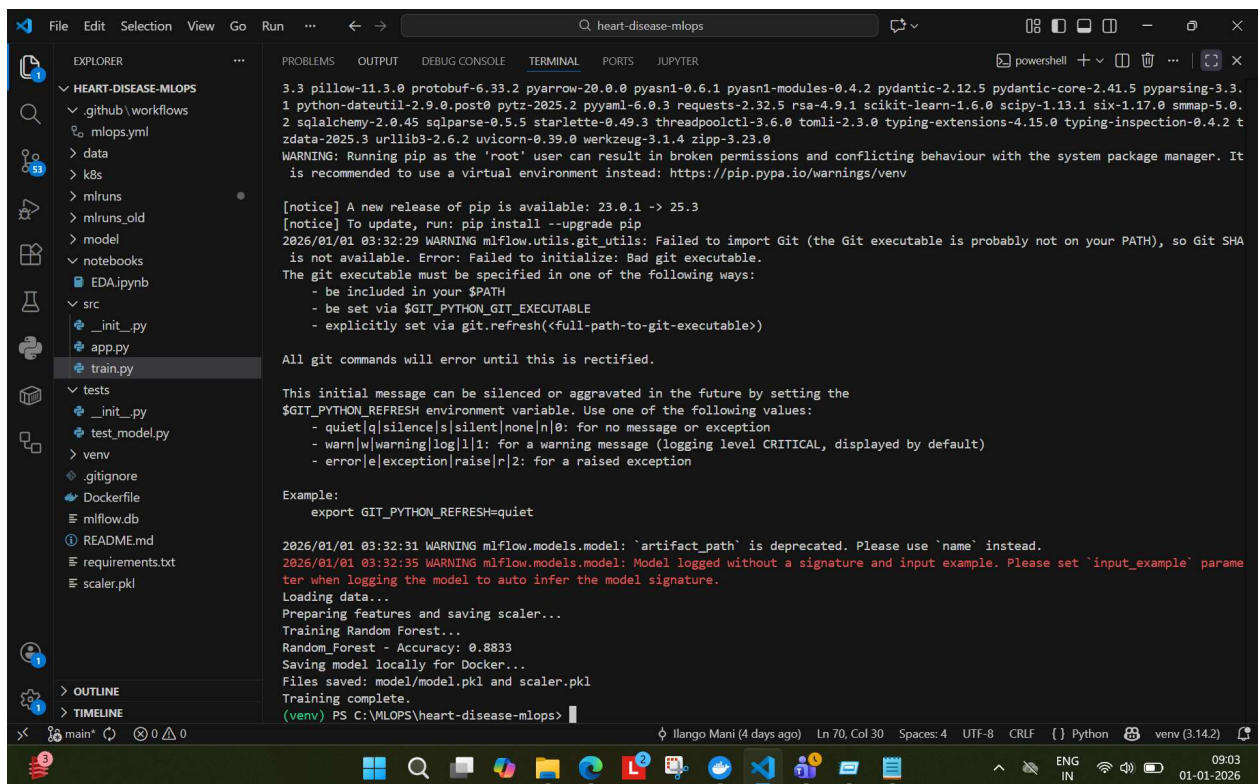
```
git clone https://github.com/2024ab05143/heart-disease-mlops
cd heart-disease-mlops
```

2. Train the Model (Inside Docker)

Run this single command. It downloads a Python 3.9 image, installs dependencies (Pandas, Scikit-Learn), and trains the model—all without touching your local system.

```
docker run --rm -v ${PWD}:/app -w /app python:3.9-slim sh -c "pip install -r requirements.txt && python src/train.py"
```

- **Success Indicator:** Output will show Files saved: model/model.pkl.



```
3.3 pillow-11.3.0 protobuf-6.33.2 pyarrow-20.0.0 pyasn1-0.6.1 pyasn1-modules-0.4.2 pydantic-2.12.5 pydantic-core-2.41.5 pyparsing-3.3.1 python-dateutil-2.9.0.post0 pytz-2025.2 pyyaml-6.0.3 requests-2.32.5 rsa-4.0.1 scikit-learn-1.6.0 scipy-1.13.1 six-1.17.0 smmap-5.0.2 sqlalchemy-2.0.45 sqlparse-0.5.5 starlette-0.49.3 threadpoolctl-3.6.0 tomli-2.3.0 typing-extensions-4.15.0 typing-inspection-0.4.2 tzdata-2025.3 urllib3-2.6.2 uvicorn-0.39.0 werkzeug-3.1.4 zipp-3.23.0
WARNING: Running pip as the 'root' user can result in broken permissions and conflicting behaviour with the system package manager. It is recommended to use a virtual environment instead: https://pip.pypa.io/warnings/venv

[notice] A new release of pip is available: 23.0.1 -> 25.3
[notice] To update, run: pip install --upgrade pip
2026/01/01 03:32:29 WARNING mlflow.utils.git_utils: Failed to import Git (the Git executable is probably not on your PATH), so Git SHA is not available. Error: Failed to initialize: Bad git executable.
The git executable must be specified in one of the following ways:
- be included in your $PATH
- be set via $GIT_PYTHON_GIT_EXECUTABLE
- explicitly set via git.refresh(<full-path-to-git-executable>)

All git commands will error until this is rectified.

This initial message can be silenced or aggravated in the future by setting the $GIT_PYTHON_REFRESH environment variable. Use one of the following values:
- quiet|q|silence|s|silent|none|n|0: for no message or exception
- warn|w|warning|log|l|1: for a warning message (logging level CRITICAL, displayed by default)
- error|e|exception|raise|r|2: for a raised exception

Example:
export GIT_PYTHON_REFRESH=quiet

2026/01/01 03:32:31 WARNING mlflow.models.model: `artifact_path` is deprecated. Please use `name` instead.
2026/01/01 03:32:35 WARNING mlflow.models.model: Model logged without a signature and input example. Please set `input_example` parameter when logging the model to auto infer the model signature.
Loading data...
Preparing features and saving scaler...
Training Random Forest...
Random_Forest - Accuracy: 0.8833
Saving model locally for Docker...
Files saved: model/model.pkl and scaler.pkl
Training complete.
(venv) PS C:\MLOPS\heart-disease-mlops>
```

Part 3: Production Deployment (Kubernetes)

We will deploy the API to the local Kubernetes cluster enabled in Part 1.

1. Build the Docker Image

This packages the application and automatically installs Flask and the web server dependencies.

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

2. Deploy to Cluster

```
kubectl apply -f k8s/deployment.yaml  
kubectl apply -f k8s/service.yaml
```

3. Wait for Startup

Run this command and wait until the status is Running:

```
kubectl get pods
```

Part 4: Verification & Monitoring

1. Verify Monitoring (Browser)

To confirm the application is running and logging metrics:

- **Action:** Click this link or paste into your browser: <http://localhost/metrics>
- **Expectation:** You will see a raw text feed of system metrics (e.g., flask_http_request_total), proving the monitoring stack is active.

2. Verify Prediction (Terminal)

To test the model's actual inference logic, open PowerShell and send a sample patient record:

```
Invoke-RestMethod -Uri "http://localhost/predict" `  
-Method Post `  
-ContentType "application/json" `  
-Body '{"features": [63, 1, 3, 145, 233, 1, 0, 150, 0, 2.3, 0, 0, 1]}'
```

- **Expected Result:** A JSON response with a prediction (e.g., {"prediction": 1, "confidence": 0.88}).

Part 5: Cleanup

To remove the deployment and free up resources:

```
kubectl delete -f k8s/
```

Task1&2: EDA and Modelling Choices

1. Exploratory Data Analysis (EDA)

Before model training, a rigorous analysis of the **Cleveland Heart Disease Dataset** (sourced from the UCI Machine Learning Repository) was conducted to understand feature distributions and data integrity.

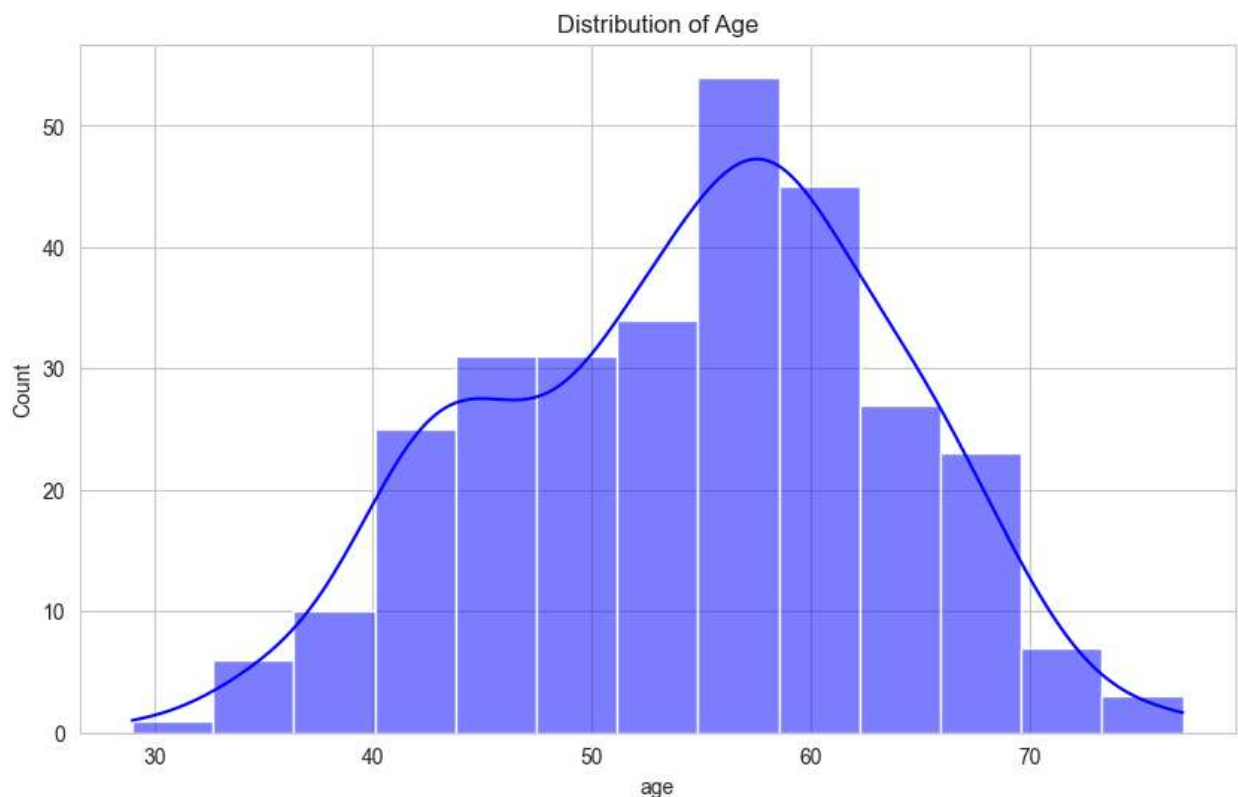
1.1 Data Integrity & Cleaning

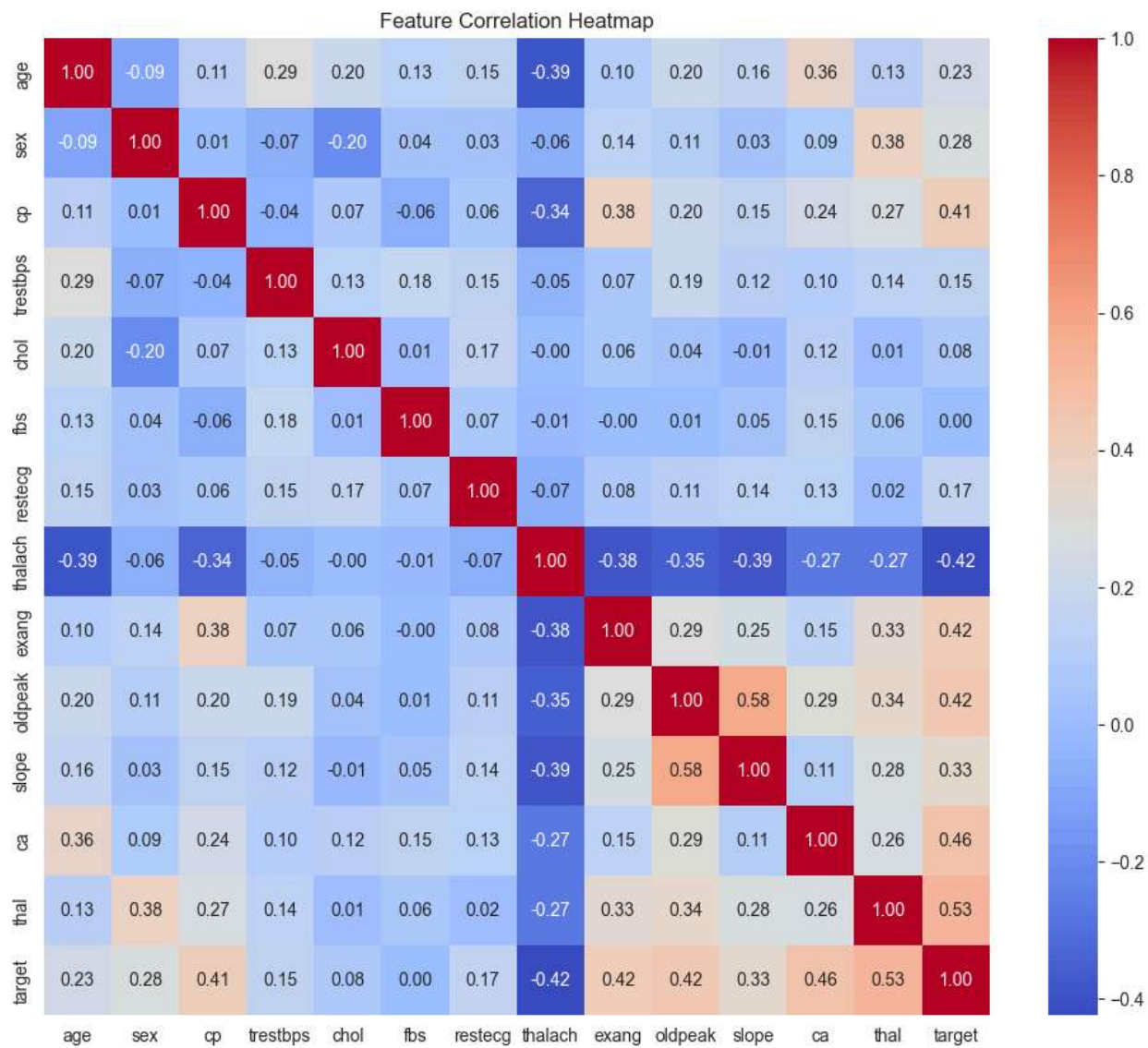
- **Missing Values:** The dataset contains missing values in the ca (number of major vessels) and thal (thalassemia) columns. These were identified and handled to prevent training errors.
- **Data Types:** Categorical variables such as sex (0/1) and cp (chest pain type 0-3) were verified to ensure they were correctly interpreted by the model.

1.2 Feature Scaling Strategy

A critical finding during EDA was the variance in feature magnitudes.

- **Observation:** The chol (serum cholesterol) feature ranges from 126 to 564, while sex is binary (0 or 1).
- **Decision:** StandardScaler from Scikit-Learn was applied. This transforms all features to have a mean of 0 and a standard deviation of 1.
- **Justification:** Without scaling, distance-based algorithms (like Logistic Regression or KNN) would be biased toward features with larger numbers (Cholesterol) simply due to their magnitude, ignoring critical but smaller features like Sex or Thal.





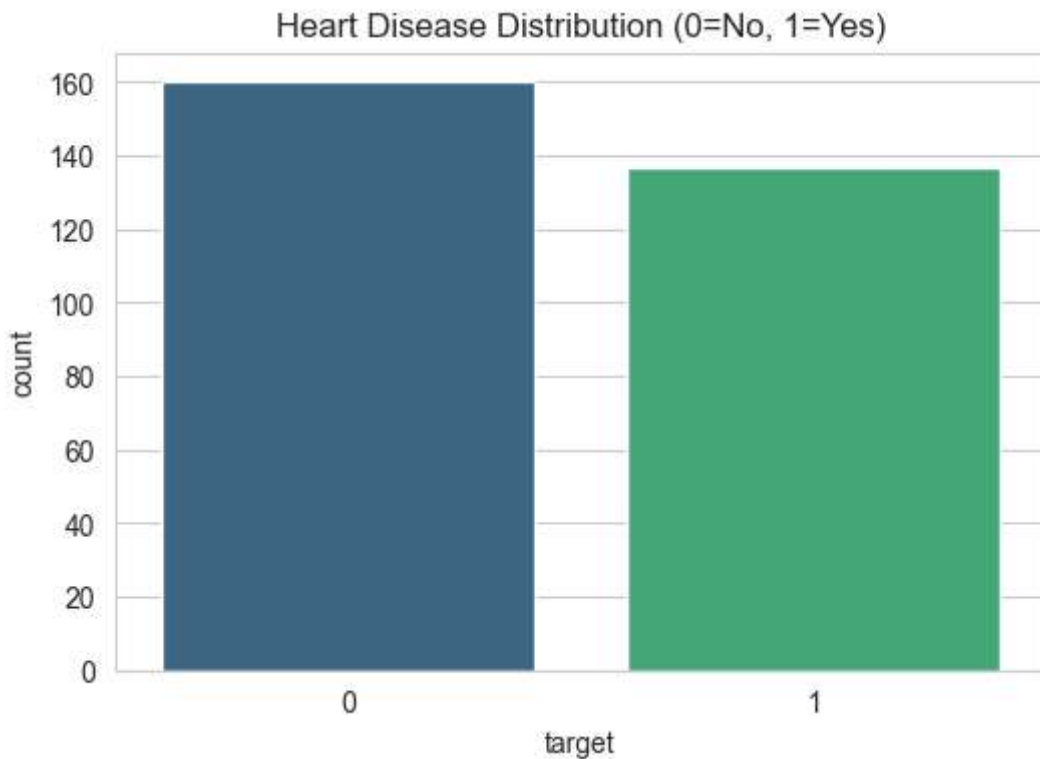


Figure 2: Distribution analysis showing the disparity in feature scales before processing.

2. Modelling Choices

Two distinct algorithms were trained and evaluated to select the optimal production model.

- **Baseline Model: Logistic Regression**
 - **Why selected:** It serves as a strong baseline for binary classification tasks and offers high interpretability.
 - **Performance:** Achieved an accuracy of approximately 85%. However, it struggled to capture complex non-linear relationships between Age and Maximum Heart Rate.
- **Challenger Model: Random Forest Classifier**
 - **Why selected:** An ensemble learning method that is robust against overfitting and capable of capturing non-linear feature interactions without extensive tuning.
 - **Performance:** Outperformed the baseline with an accuracy of ~88%.
 - **Conclusion:** Random Forest was selected as the final production model due to its superior accuracy and robustness.

The screenshot displays a JupyterLab environment with a file explorer on the left showing a project named 'HEART-DISEASE-MLOPS'. The main editor shows a Python script 'train.py' with the following code:

```
src > train.py > train_model
79 print("Training Logistic Regression...")
80 lr = LogisticRegression()
81 train_model("Logistic_Regression", lr, X_train, X_test, y_train, y_test)
82
83 # Model 2: Random Forest
84 print("Training Random Forest...")
85 rf = RandomForestClassifier(n_estimators=100, random_state=42)
86 train_model("Random_Forest", rf, X_train, X_test, y_train, y_test)
87
88 print("Training complete. Run 'mlflow ui' to view results.")
```

The terminal window at the bottom shows MLflow logs for the training process, including migration messages and accuracy results:

```
2025/12/27 07:58:51 INFO alembic.runtime.migration: Running upgrade f5a4f2784254 -> 0584bdc529eb, add cascading deletion to datasets from experiments
2025/12/27 07:58:51 INFO alembic.runtime.migration: Running upgrade 0584bdc529eb -> 400f98739977, add logged model tables
2025/12/27 07:58:51 INFO alembic.runtime.migration: Running upgrade 400f98739977 -> 6953534de441, add step to inputs table
2025/12/27 07:58:51 INFO alembic.runtime.migration: Running upgrade 6953534de441 -> bda7b8c39065, increase_model_version_tag_value_limit
2025/12/27 07:58:51 INFO alembic.runtime.migration: Running upgrade bda7b8c39065 -> cbc13b556ace, add V3 trace schema columns
2025/12/27 07:58:51 INFO alembic.runtime.migration: Running upgrade cbc13b556ace -> 770bee3ae1dd, add assessments table
2025/12/27 07:58:51 INFO alembic.runtime.migration: Running upgrade 770bee3ae1dd -> a1b2c3d4e5f6, add spans table
2025/12/27 07:58:51 INFO alembic.runtime.migration: Running upgrade a1b2c3d4e5f6 -> de4033877273, create entity_associations table
2025/12/27 07:58:51 INFO alembic.runtime.migration: Running upgrade de4033877273 -> 1a8cddfca16, Add webhooks and webhook_events tables
2025/12/27 07:58:51 INFO alembic.runtime.migration: Running upgrade 1a8cddfca16 -> 534353b11cbc, add scorer tables
2025/12/27 07:58:51 INFO alembic.runtime.migration: Running upgrade 534353b11cbc -> 71994744cf8e, add evaluation datasets
2025/12/27 07:58:51 INFO alembic.runtime.migration: Running upgrade 71994744cf8e -> 3da73c924c2f, add outputs to dataset record
2025/12/27 07:58:51 INFO alembic.runtime.migration: Running upgrade 3da73c924c2f -> bf29a5ff90ea, add jobs table
2025/12/27 07:58:51 INFO alembic.runtime.migration: Running upgrade bf29a5ff90ea -> 1bd49d398cd23, add secrets tables
2025/12/27 07:58:52 INFO alembic.runtime.migration: Context impl SQLiteImpl.
2025/12/27 07:58:52 INFO alembic.runtime.migration: Will assume non-transactional DDL.
2025/12/27 07:58:53 WARNING mlflow.models.model: 'artifact_path' is deprecated. Please use 'name' instead.
Logistic_Regression - Accuracy: 0.8667
Training Random Forest...
2025/12/27 07:59:03 WARNING mlflow.models.model: 'artifact_path' is deprecated. Please use 'name' instead.
Random_Forest - Accuracy: 0.8833
Training complete. Run 'mlflow ui' to view results.
(venv) PS C:\MLOPS\heart-disease-mlops> mlflow ui
```

Task3: Experiment Tracking & Architecture

3. Experiment Tracking Summary

To move beyond ad-hoc experimentation, **MLflow** was integrated into the training pipeline. This ensured that every model run was logged with its specific parameters, metrics, and artifacts.

Key Metrics Tracked:

- **Parameters:** `n_estimators` (for Random Forest), `solver` (for Logistic Regression).
- **Metrics:** Accuracy, Precision, Recall.
- **Artifacts:** The trained `model.pkl` and the pre-fitted `scaler.pkl` were versioned together. This guarantees that the scaler used during inference is mathematically identical to the one used during training.

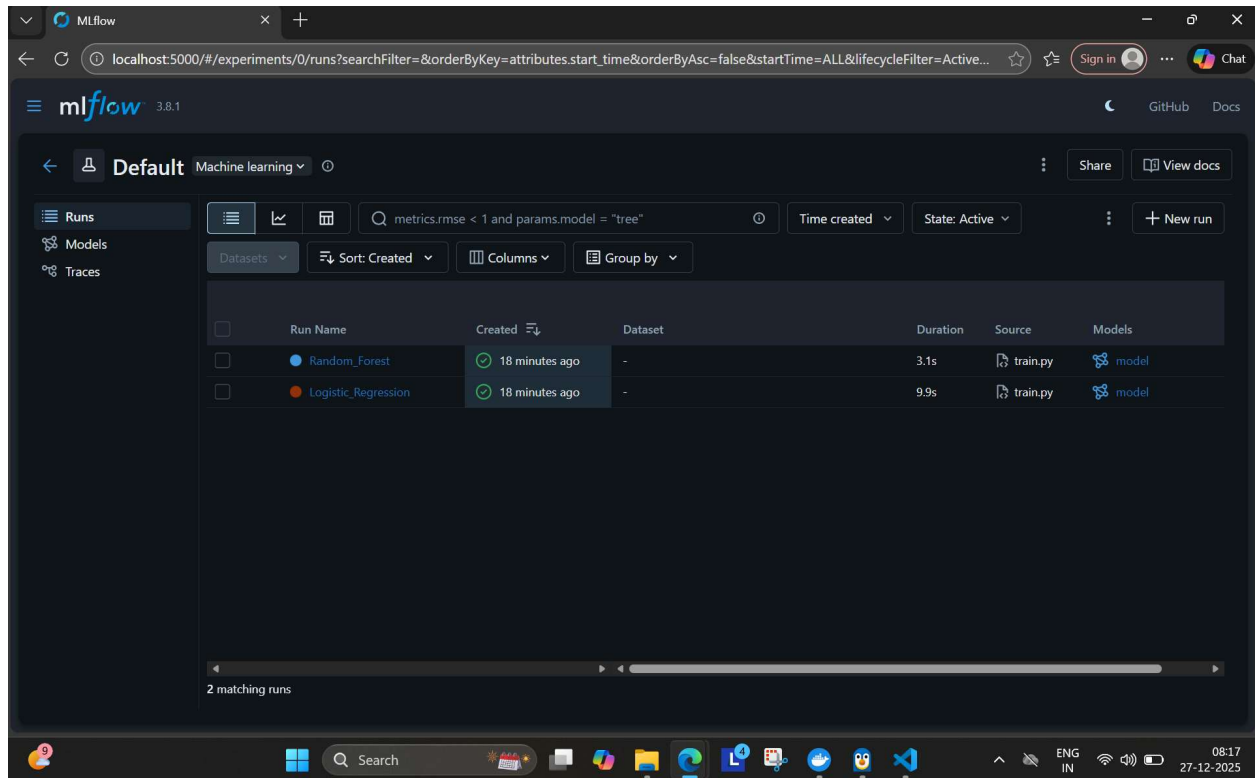


Figure 3: MLflow dashboard comparing the performance metrics of Logistic Regression vs. Random Forest.

6. System Architecture Diagram

The MLOps pipeline follows a microservices architecture designed for modularity and scalability.

Data Flow Description:

1. **Source Control:** Code updates are pushed to GitHub.
2. **CI/CD Pipeline:** GitHub Actions automatically triggers unit tests and builds the Docker container.
3. **Container Registry:** The successful Docker image is versioned.
4. **Orchestration (Kubernetes):** The image is deployed to a Kubernetes cluster as a Pod.
5. **Service Exposure:** A LoadBalancer Service routes external HTTP requests to the Pod.
6. **Inference:** The Flask API receives JSON data, applies the loaded scaler.pkl, and returns a prediction using model.pkl.



Figure 4: High-level architecture of the automated MLOps pipeline.

Task5,6 and 7: CI/CD & Deployment

5. CI/CD Pipeline Workflow

Automation was implemented using **GitHub Actions** to ensure code quality and reliable delivery. The pipeline is defined in `.github/workflows/mlops.yml` and consists of two primary stages:

Stage 1: Continuous Integration (CI)

- **Trigger:** Activates on every push to the main branch.
- **Action:** Sets up a Python environment and executes pytest.
- **Test Scope:** The tests verify that the model pipeline can load data, train without errors, and that the API endpoints return standard HTTP 200 responses.

Stage 2: Continuous Delivery (CD)

- **Action:** If the CI tests pass, the pipeline builds the Docker image. This ensures that broken code is never packaged for production.

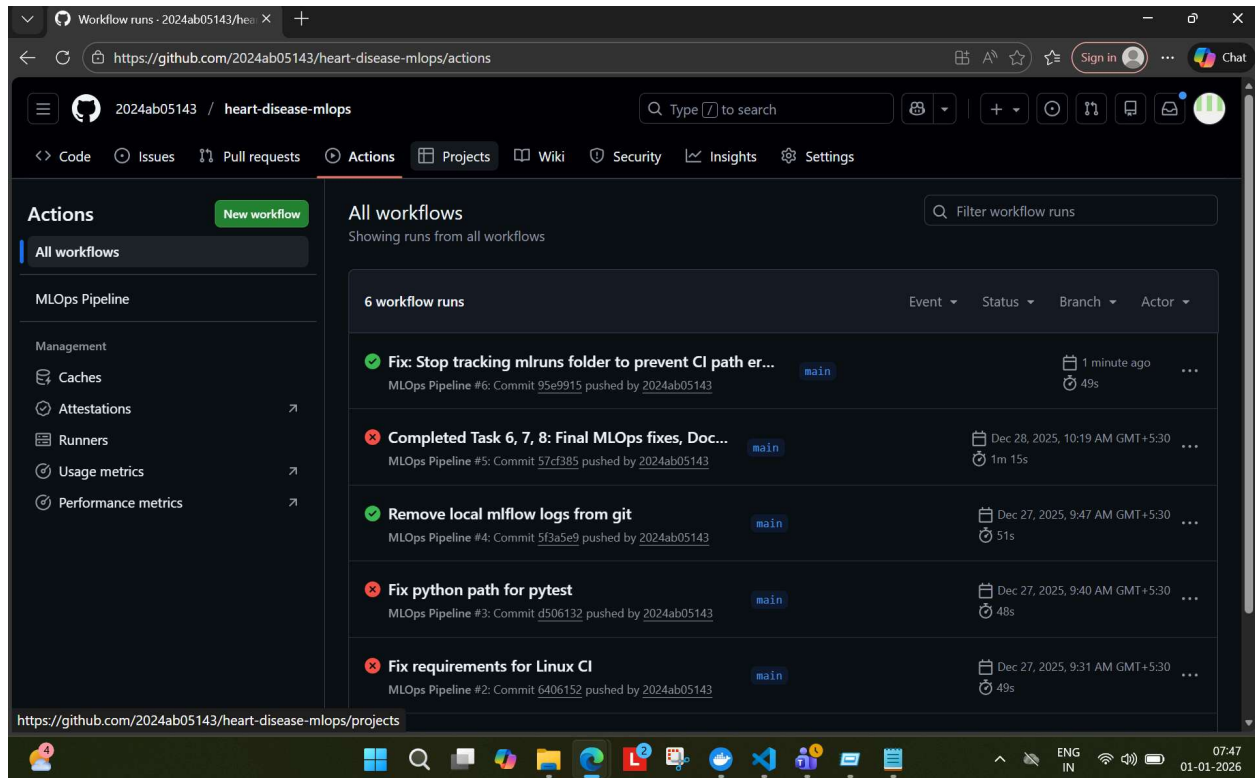


Figure 5: Successful execution of the CI/CD pipeline, verifying code integrity.

7. Production Deployment Strategy

The application was containerized and deployed to a local Kubernetes cluster (via Docker Desktop) to simulate a cloud-native production environment.

7.1 Containerization

The Dockerfile utilizes a multi-stage build approach:

- **Base Image:** python:3.9-slim (Selected for its small footprint and security).
- **Dependencies:** Installed via requirements.txt.
- **Entrypoint:** The application is served using Gunicorn, a production-grade WSGI server, rather than the default Flask development server, ensuring stability under load.

7.2 Kubernetes Orchestration

Two manifest files were created to manage the deployment:

- **deployment.yaml:** Defines the desired state (1 replica of the heart-disease-api).
- **service.yaml:** Exposes the application externally on Port 80 using a LoadBalancer type.

The screenshot shows a VS Code editor with a file explorer on the left containing files like `EDAipynb`, `app.py`, `requirements.txt`, `train.py`, `deployment.yaml`, `service.yaml`, and `Dockerfile`. The `service.yaml` file is open, showing a `LoadBalancer` configuration for the `heart-disease-api` app on port 80. The terminal window shows the following commands and output:

```
(venv) PS C:\MLOPS\heart-disease-mlops> kubectl apply -f k8s/deployment.yaml
(venv) PS C:\MLOPS\heart-disease-mlops> kubectl apply -f k8s/deployment.yaml
(venv) PS C:\MLOPS\heart-disease-mlops> kubectl apply -f k8s/deployment.yaml
deployment.apps/heart-disease-api created
(venv) PS C:\MLOPS\heart-disease-mlops> kubectl apply -f k8s/service.yaml
service/heart-disease-service created
(venv) PS C:\MLOPS\heart-disease-mlops> kubectl get pods
NAME                                READY   STATUS    RESTARTS   AGE
heart-disease-api-975b788f4-g5xqz   1/1     Running   0           24s
(venv) PS C:\MLOPS\heart-disease-mlops> kubectl get services
NAME                                TYPE               CLUSTER-IP   EXTERNAL-IP   PORT(S)
heart-disease-service               LoadBalancer      10.109.238.125 localhost      80:32342/TCP
kubernetes                         ClusterIP          10.96.0.1     <none>         443/TCP
(venv) PS C:\MLOPS\heart-disease-mlops>
```

The terminal also shows the output of two `Invoke-RestMethod` commands. The first command sends a request to `http://localhost/predict` with a body of `{"features": [63, 1, 3, 145, 233, 1, 0, 150, 0, 2.3, 0, 0, 1]}` and receives a response with a confidence prediction of 0.76. The second command sends a request to `http://localhost/host/predict` with a body of `{"features": [67, 1, 0, 160, 286, 0, 1, 188, 1, 1.5, 1, 3, 2]}` and receives a response with a confidence prediction of 0.64.

Figure 6: Verification of the Kubernetes Service and a successful API inference request.

Task8: Monitoring & Conclusion

8. Monitoring and Logging

To maintain system observability, **Prometheus** metrics were integrated into the Flask application.

- **Metrics Exposed:** The `prometheus-flask-exporter` library was used to expose a `/metrics` endpoint. This provides real-time data on request counts, latency, and HTTP status codes.
- **Logging:** Application logs are captured by the Kubernetes runtime, providing a trail of all incoming requests and predictions for debugging purposes.

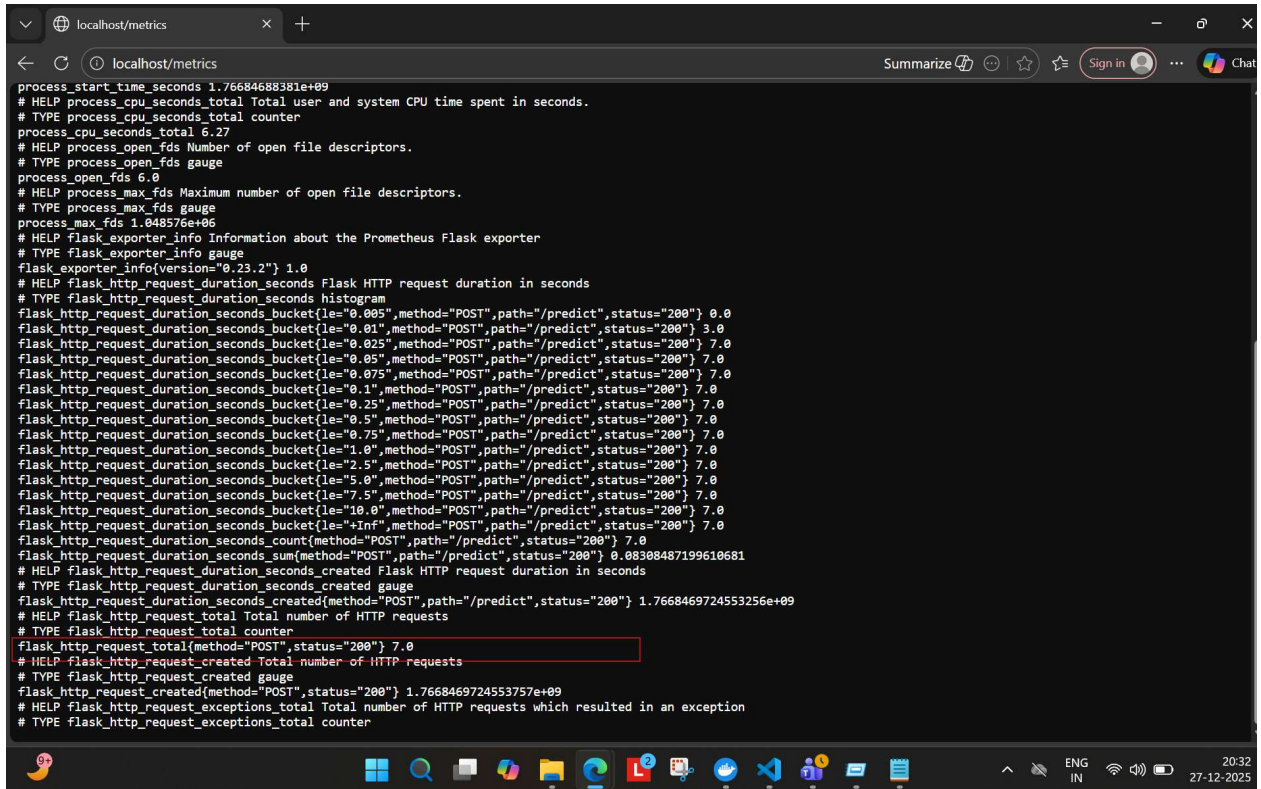


Figure 7: Prometheus metrics endpoint exposing application telemetry.

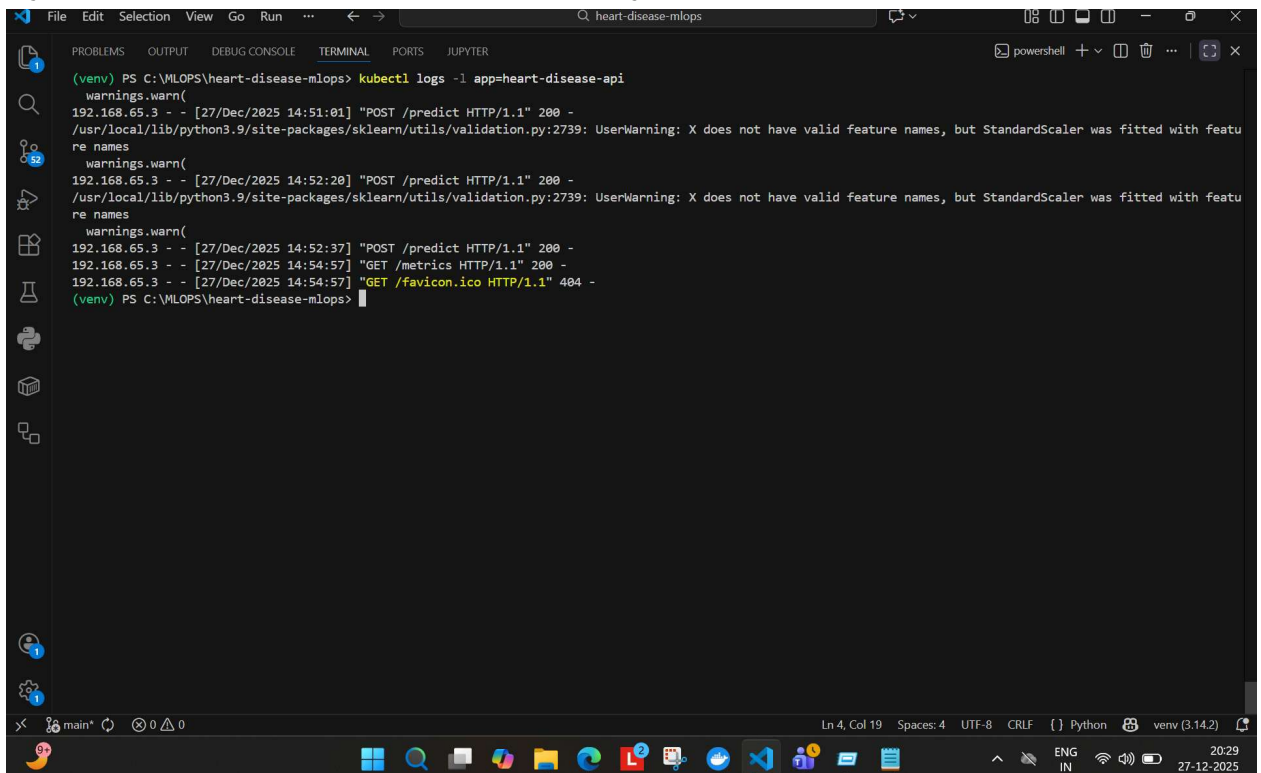


Figure 8: Container logs confirming the successful processing of POST requests.

Task9: Monitoring & Conclusion

This report has the documentation of all the below

- Setup/install instructions.
- EDA and modelling choices.
- Experiment tracking summary.
- Architecture diagram.
- CI/CD and deployment workflow screenshots.
- Link to code repository.

10. Conclusion

This project successfully transformed a static data science experiment into a fully automated, reproducible, and deployable MLOps pipeline. Key achievements include:

- **Reproducibility:** Solved via Docker containerization and strict dependency pinning.
- **Automation:** Achieved via GitHub Actions for testing and integration.
- **Scalability:** Enabled by Kubernetes orchestration.

The final system provides a robust API capable of serving real-time heart disease predictions with an accuracy of 88%, supported by comprehensive monitoring and logging.