

BITS ID	Name	Contribution
2024AA05041	Amit Kumar	100%
2024AA05042	Yogesh Agarwal	100%
2024AA05633	Yawar Abbas	100%

Heart Disease MLOps Solution: Step-by-Step Report with Code

Introduction & Project Overview

Purpose:

This project builds a full machine learning pipeline to predict heart disease, using the UCI Heart Disease dataset. The solution covers data download, cleaning, EDA, feature engineering, model training, experiment tracking, packaging, deployment, and monitoring.

Workflow Outline:

- Data download & cleaning
- Exploratory data analysis (EDA)
- Feature engineering
- Model training & evaluation
- Experiment tracking (MLflow)
- Model packaging & deployment (Docker, FastAPI, K8s)
- Monitoring & logging (Prometheus, Grafana)
- CI/CD automation (DVC, GitHub Actions)

Project Setup & Environment Checks

Cells 1–2: Project Initialization

Role:

- Sets up folders, checks Docker, and ensures a reproducible environment.

Key Code:

```
1 from pathlib import Path
2 import os
3
4 # Set up project root and folders
5 PROJECT_ROOT = Path("C:/.../mlops-assignment-final")
6 PROJECT_ROOT.mkdir(exist_ok=True)
7 os.chdir(PROJECT_ROOT)
8
9 # Print project structure
10 print(f"✓ PROJECT: {PROJECT_ROOT.absolute()}")
11 print("\n✓ FOLDERS:")
12 print(os.listdir("."))
13 ``
```

Why:

Ensures all files are organized and the environment is ready for reproducible work.

Data Acquisition & Cleaning

Cells 3–5: Downloading and Preparing Data

Role:

- Downloads the UCI Heart Disease dataset, assigns column names, handles missing values, and saves both raw and cleaned datasets.

Key Code

```

1 # Download dataset
2 url = "https://archive.ics.uci.edu/ml/machine-learning-databases/heart-
disease/processed.cleveland.data"
3 columns = ['age', 'sex', 'cp', 'trestbps', 'chol', 'fbs', 'restecg',
4             'thalach', 'exang', 'oldpeak', 'slope', 'ca', 'thal', 'target']
5 df_raw = pd.read_csv(url, names=columns, na_values='?')
6 df_raw.to_csv("data/heart_disease_raw.csv", index=False)
7
8 # clean data
9 df_clean = df_raw.copy()
10 df_clean['target'] = (df_clean['target'] > 0).astype(int)
11 # Impute missing values
12 imputer_num = SimpleImputer(strategy='median')
13 imputer_cat = SimpleImputer(strategy='most_frequent')
14 df_clean[numeric_features] = imputer_num.fit_transform(df_clean[numeric_features])
15 df_clean[categorical_features] = imputer_cat.fit_transform(df_clean[categorical_features])
16 df = df_clean.dropna()
17 df.to_csv("data/heart_disease_cleaned.csv", index=False)
18

```

Why:

Clean, well-structured data is essential for building reliable machine learning models.

Exploratory Data Analysis (EDA)

Cells 6–7: Data Exploration and Visualization

Role:

- Shows summary statistics and creates professional plots for class balance, age, cholesterol, heart rate, and a correlation heatmap.

Key Code:

```

1 # Class balance plot
2 sns.countplot(data=df, x='target')
3 plt.title('Heart Disease Distribution (0=No, 1=Yes)')
4 plt.savefig('data/task1_eda_plots.png')
5
6 # Correlation heatmap
7 corr_matrix = df.corr()
8 sns.heatmap(corr_matrix, annot=True, cmap='coolwarm')
9 plt.title('Feature Correlation Matrix')
10 plt.savefig('data/task1_correlation_heatmap.png')

```

Why:

EDA helps us understand the data, spot patterns, and decide which features might be important for prediction.

Feature Engineering & Preprocessing

Cells 8–10: Building the Feature Pipeline

Role:

- Splits data into train/test, builds a preprocessing pipeline (imputation, scaling, encoding), and transforms the data.

Key Code:

```

5  from sklearn.pipeline import Pipeline
6
7  # Split data
8  X = df.drop('target', axis=1)
9  y = df['target']
10 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, stratify=y)
11
12 # Preprocessing pipeline
13 preprocessor = ColumnTransformer([
14     ('num', Pipeline([
15         ('imputer', SimpleImputer(strategy='median')),
16         ('scaler', StandardScaler())
17     ])), numeric_features,
18     ('cat', Pipeline([
19         ('imputer', SimpleImputer(strategy='most_frequent')),
20         ('onehot', OneHotEncoder(drop='first', sparse_output=False,
21             handle_unknown='ignore'))
22     ])), categorical_features)
23 ])
24 X_train_processed = preprocessor.fit_transform(X_train)

```

Why:

Consistent preprocessing ensures the model sees data in the same format during training and prediction.

Model Training & Evaluation

Cells 11–13: Training and Comparing Models

Role:

- Trains Logistic Regression and Random Forest models, tunes hyperparameters, and compares their performance.

Key Code:

```
1 # Logistic Regression
2 lr_model = LogisticRegression(random_state=42, max_iter=1000)
3 lr_model.fit(X_train_processed, y_train)
4 y_pred_lr = lr_model.predict(X_test_processed)
5 print("Accuracy:", accuracy_score(y_test, y_pred_lr))
6
7 # Random Forest with GridSearchCV
8 param_grid = {
9     'n_estimators': [100, 200],
10    'max_depth': [5, 10, None],
11    'min_samples_split': [2, 5]
12 }
13 rf = RandomForestClassifier(random_state=42)
14 grid_search = GridSearchCV(rf, param_grid, cv=5, scoring='roc_auc')
```

Why:

Comparing models helps us pick the best one for deployment.

Experiment Tracking with MLflow

Cells 14–19: Logging Experiments

Role:

- Sets up MLflow, logs all model parameters, metrics, artifacts, and provides a summary table of experiments.

Key Code:

```
1 import mlflow
2 mlflow.set_tracking_uri("file:./mlruns")
3 mlflow.set_experiment("Heart_Disease_Experiments")
4
5 with mlflow.start_run(run_name="Logistic_Regression_Baseline"):
6     mlflow.log_params({'model_type': 'LogisticRegression', 'max_iter': 1000})
7     mlflow.log_metrics({'accuracy': accuracy_score(y_test, y_pred_lr)})
8     mlflow.sklearn.log_model(lr_model, "model")
```

Why:

MLflow makes it easy to track, compare, and reproduce experiments, which is essential for MLOps.

Model Packaging & Deployment

Cells 20–26: Production Pipeline and API

Role:

- Packages the best model and preprocessing pipeline, creates a FastAPI app, and generates Dockerfiles for deployment.

Key Code:

```
1 # Save full pipeline
2 from sklearn.pipeline import Pipeline
3 import joblib
4 full_pipeline = Pipeline([
5     ('preprocessor', preprocessor),
6     ('classifier', best_rf)
7 ])
8 joblib.dump(full_pipeline, "models/heart_disease_pipeline_prod.joblib")
9
10 # FastAPI app (main.py)
11 from fastapi import FastAPI
12 import joblib
13 import pandas as pd
14
15 app = FastAPI()
16 model = joblib.load("models/heart_disease_pipeline_prod.joblib")
17
18 @app.post("/predict")
19 def predict(data: dict):
```

Why:

Containerization and APIs make it easy to deploy the model anywhere and serve predictions reliably.

CI/CD Automation & Monitoring

Cells 27–49: Automation, Monitoring, and Reporting

Role:

- Sets up DVC for data/model versioning, GitHub Actions for CI/CD, Kubernetes manifests for deployment, and Prometheus/Grafana for monitoring.

Key Code:

```
1 # .github/workflows/mllops.yml (GitHub Actions)
2 name: MLOps CI/CD Pipeline
3 on:
4   push:
5     branches: [ main ]
6   pull_request:
7     branches: [ main ]
8 jobs:
9   test-train-deploy:
10     runs-on: ubuntu-latest
11     steps:
12       - uses: actions/checkout@v4
13       - name: Set up Python 3.11
14         uses: actions/setup-python@v4
15         with:
16           python-version: 3.11
```

Why:

Automation and monitoring are key to reliable, maintainable, and observable ML systems in production.

Docker: Containerizing the ML Model API

Role:

Docker packages the ML model and API (FastAPI/Flask) into a portable container. This ensures the solution runs identically on any system.

Key Steps:

- Write a Dockerfile to install dependencies and copy code/model files.
- Expose the /predict endpoint for inference.

Sample Dockerfile:

```
1 FROM python:3.11-slim
2 WORKDIR /app
3 COPY requirements.txt .
4 RUN pip install -r requirements.txt
5 COPY . .
6 CMD ["uvicorn", "main:app", "--host", "0.0.0.0", "--port", "8080"]
```

How to Build & Run:

```
1 docker build -t heart-disease-api .
2 docker run -p 8080:8080 heart-disease-api
```

Kubernetes: Production Deployment & Orchestration

Role:

Kubernetes (K8s) manages deployment, scaling, and reliability of the containerized API. It can run locally (Minikube, Docker Desktop) or on cloud (GKE, EKS, AKS).

Key Steps:

- Write a deployment manifest or Helm chart.
- Expose the API via LoadBalancer or Ingress.

Sample Deployment Manifest:

```
apiVersion: apps/v1
kind: Deployment
metadata:
  name: heart-disease-api
spec:
  replicas: 2
  selector:
```

```
matchLabels:  
  app: heart-disease-api  
  
template:  
  metadata:  
    labels:  
      app: heart-disease-api  
  
  spec:  
    containers:  
      - name: api  
        image: heart-disease-api:latest  
    ports:  
      - containerPort: 8080  
  
---  
  
apiVersion: v1  
kind: Service  
  
metadata:  
  name: heart-disease-service  
  
spec:  
  type: LoadBalancer  
  selector:  
    app: heart-disease-api  
  ports:  
    - protocol: TCP  
      port: 80  
      targetPort: 8080
```

Grafana: Monitoring & Visualization

Role:

Grafana visualizes metrics collected from the API, such as request rates, errors, and latency. It connects to Prometheus as a data source.

Key Steps:

- Install Grafana and Prometheus (often via Helm charts).
- Import a dashboard JSON to visualize API metrics.

Sample Grafana Dashboard Setup:

- Add Prometheus as a data source.
- Import dashboard JSON (provided in your repo).
- View real-time metrics for API health.

Loki: Centralized Logging

Role:

Loki collects and indexes logs from the API containers, making it easy to search and analyze logs in Grafana.

Key Steps:

- Install Loki and configure it as a Grafana data source.
- Ensure API logs are written to stdout/stderr for collection.

Sample Logging Code (Python/FastAPI):

```
1 import logging
2 logging.basicConfig(level=logging.INFO)
3 logger = logging.getLogger("api_logger")
4
5 @app.post("/predict")
6 def predict(data: dict):
7     logger.info(f"Received prediction request: {data}")
8     # ...rest of code...
```

Prometheus: Metrics Collection

Role:

Prometheus scrapes metrics from the API (using /metrics endpoint) and stores them for visualization in Grafana.

Key Steps:

- Add Prometheus client code to expose metrics.
- Configure Prometheus to scrape the API.

Sample Prometheus Metrics Code:

```
1  from prometheus_client import Counter, generate_latest
2  REQUEST_COUNT = Counter('api_requests_total', 'Total API requests', ['method', 'endpoint',
   'status'])
3
4  @app.middleware("http")
5  async def prometheus_middleware(request, call_next):
6      response = await call_next(request)
7      REQUEST_COUNT.labels(request.method, request.url.path, response.status_code).inc()
8      return response
9
10 @app.get("/metrics")
11 def metrics():
12     return Response(generate_latest(), media_type="text/plain")
```

Prometheus Config Example:

```
1  scrape_configs:
2      - job_name: 'heart-disease-api'
3          static_configs:
4              - targets: ['heart-disease-service:8080']
```

Integration Workflow

1. **Build Docker image** for the API.
2. **Deploy on Kubernetes** using manifests/Helm charts.
3. **Expose metrics** and logs from the API.
4. **Prometheus** scrapes metrics; **Loki** collects logs.
5. **Grafana** visualizes both metrics and logs for monitoring.

