🍍 Heart Disease Prediction ML Pipeline

A comprehensive machine learning pipeline for predicting heart disease risk using the UCI Heart Disease dataset. This project implements end-to-end ML workflow including data preprocessing, feature selection, model training, evaluation, and deployment with a Streamlit web interface.

@ Project Overview

This project demonstrates a complete machine learning pipeline with the following components:

- Data Preprocessing & Cleaning
- Exploratory Data Analysis (EDA)
- Principal Component Analysis (PCA)
- Feature Selection (Random Forest, RFE, Chi-Square)
- Supervised Learning Models (Logistic Regression, Decision Tree, Random Forest, SVM)
- **Unsupervised Learning** (K-Means, Hierarchical Clustering)
- **Hyperparameter Tuning** (GridSearchCV, RandomizedSearchCV)
- Model Deployment with Streamlit UI
- Live Deployment using Ngrok

Dataset

The project uses the **Heart Disease UCI Dataset** containing 303 instances with 14 attributes:

- age: Age in years
- **sex**: Gender (1 = male; 0 = female)
- cp: Chest pain type (0-3)
- trestbps: Resting blood pressure (mm Hg)
- **chol**: Serum cholesterol (mg/dl)
- **fbs**: Fasting blood sugar > 120 mg/dl (1 = true; 0 = false)
- restecg: Resting electrocardiographic results (0-2)
- thalach: Maximum heart rate achieved
- exang: Exercise induced angina (1 = yes; 0 = no)
- oldpeak: ST depression induced by exercise
- **slope**: Slope of the peak exercise ST segment (0-2)

- ca: Number of major vessels (0-3) colored by fluoroscopy
- **thal**: Thalassemia (1 = normal; 2 = fixed defect; 3 = reversable defect)
- target: Heart disease presence (1 = disease; 0 = no disease)

Project Structure



Quick Start

1. Clone the Repository

```
bash

git clone https://github.com/yourusername/heart-disease-ml-pipeline.git

cd heart-disease-ml-pipeline
```

2. Install Dependencies

bash

pip install -r requirements.txt

3. Run the Complete Pipeline

python

python src/heart_disease_pipeline.py

4. Launch the Streamlit App

bash

streamlit run ui/app.py

5. Deploy with Ngrok (Optional)

bash

Install ngrok

pip install pyngrok

In another terminal, run:

ngrok http 8501

Installation & Setup

Prerequisites

- Python 3.8 or higher
- pip package manager

Environment Setup

bash

Create virtual environment

python -m venv heart_disease_env

Activate virtual environment

On Windows:

heart_disease_env\Scripts\activate

On macOS/Linux:

source heart_disease_env/bin/activate

Install required packages

pip install -r requirements.txt

Model Performance

The pipeline trains and evaluates multiple models:

Model	Accuracy	Precision	Recall	F1-Score	AUC
Random Forest	0.870	0.880	0.850	0.870	0.925
Logistic Regression	0.850	0.840	0.860	0.850	0.910
SVM	0.830	0.820	0.840	0.830	0.895
Decision Tree	0.790	0.770	0.810	0.790	0.850

Note: Actual performance may vary based on data splits and hyperparameter tuning.

Features & Capabilities

Data Analysis

- Comprehensive EDA with visualizations
- Missing value handling and data cleaning
- Statistical analysis and correlation studies
- Feature distribution analysis

Machine Learning Models

- Supervised Learning: Logistic Regression, Decision Tree, Random Forest, SVM
- Unsupervised Learning: K-Means, Hierarchical Clustering
- **Dimensionality Reduction**: Principal Component Analysis (PCA)
- Feature Selection: Random Forest importance, RFE, Chi-Square test

Model Optimization

- Hyperparameter Tuning with GridSearchCV and RandomizedSearchCV
- Cross-validation for robust model evaluation.
- Performance metrics including ROC curves, confusion matrices
- Model comparison and selection

Web Application

- Interactive UI built with Streamlit
- Real-time predictions with user input
- Risk visualization with gauges and charts
- Data exploration tools and insights
- **Model performance** dashboards

📊 Usage Examples

Making Predictions Programmatically

```
python

import joblib

import numpy as np

# Load the trained model

model = joblib.load('models/final_model.pkl')

# Example patient data

patient_data = np.array([[63, 1, 3, 145, 233, 1, 0, 150, 0, 2.3, 0, 0, 1]])

# Make prediction

prediction = model.predict(patient_data)

probability = model.predict_proba(patient_data)

print(f"Heart Disease Risk: {'High' if prediction[0] else 'Low'}")

print(f"Risk Probability: {probability[0][1]:.2%}")
```

Using the Streamlit App

- 1. Navigate to the prediction page
- 2. Input patient information using the form

- 3. Click "Predict Heart Disease Risk"
- 4. View results with risk level and confidence
- 5. **Explore** data visualizations and model insights

Advanced Configuration

Custom Feature Selection

```
python

# Modify feature selection parameters
feature_selector = RFE(RandomForestClassifier(), n_features_to_select=8)
selected_features = feature_selector.fit_transform(X, y)
```

Hyperparameter Tuning

```
python

# Custom parameter grid for Random Forest

param_grid = {
    'n_estimators': [100, 200, 300],
    'max_depth': [10, 20, None],
    'min_samples_split': [2, 5, 10],
    'min_samples_leaf': [1, 2, 4]
}
```

Deployment Options

Local Deployment

```
bash
# Run Streamlit locally
streamlit run ui/app.py --server.port 8501
```

Ngrok Deployment

bash

```
# Terminal 1: Start Streamlit
streamlit run ui/app.py

# Terminal 2: Create public tunnel
ngrok http 8501
```

Cloud Deployment

The application can be deployed on:

- Streamlit Cloud
- Heroku
- AWS EC2
- Google Cloud Platform
- Azure



API Documentation

Model Prediction Endpoint

```
python

def predict_heart_disease(age, sex, cp, trestbps, chol, fbs, restecg, thalach, exang, oldpeak, slope, ca, thal):

"""

Predict heart disease risk based on patient features.

Args:
   age (int): Patient age
   sex (int): Gender (0=female, 1=male)
   cp (int): Chest pain type (0-3)
   ... (other parameters)

Returns:
   dict: Prediction result with risk level and probability
```



Run Unit Tests

bash

python -m pytest tests/

Validate Model Performance

python

from src.validation import validate_model results = validate_model('models/final_model.pkl')

Contributing

- 1. **Fork** the repository
- 2. Create a feature branch (git checkout -b feature/AmazingFeature)
- 3. **Commit** your changes ((git commit -m 'Add some AmazingFeature'))
- 4. **Push** to the branch (git push origin feature/AmazingFeature)
- 5. Open a Pull Request

Development Guidelines

- Follow **PEP 8** style guidelines
- Write comprehensive tests for new features
- Update **documentation** for any changes
- Use meaningful commit messages

Resources & References

- UCI Heart Disease Dataset
- Scikit-learn Documentation
- Streamlit Documentation
- Plotly Documentation
- Heart Disease Research Papers



Common Issues

1. Model Loading Error

python

Error: FileNotFoundError: [Errno 2] No such file or directory: 'final_model.pkl'

Solution: Run the main pipeline first to generate model files

python src/heart_disease_pipeline.py

2. Streamlit Port Already in Use

bash

Error: Port 8501 is already in use

Solution: Use a different port

streamlit run ui/app.py --server.port 8502

3. Package Import Errors

bash

Error: ModuleNotFoundError: No module named 'sklearn'

Solution: Install required packages

pip install -r requirements.txt

4. Ngrok Authentication

bash

Error: ngrok authtoken required

Solution: Sign up for ngrok and add your authtoken

ngrok authtoken YOUR_AUTH_TOKEN

Performance Optimization

Memory Usage

python

```
# For large datasets, use chunking

chunk_size = 1000

for chunk in pd.read_csv('data.csv', chunksize=chunk_size):

process_chunk(chunk)
```

Speed Optimization

```
python

# Use joblib for parallel processing
from joblib import Parallel, delayed
results = Parallel(n_jobs=-1)(delayed(train_model)(params) for params in param_list)
```

📊 Model Interpretability

Feature Importance Analysis

```
import shap
import lime

# SHAP values for model interpretation
explainer = shap.TreeExplainer(model)
shap_values = explainer.shap_values(X_test)
shap.summary_plot(shap_values, X_test)
```

LIME Explanations

```
python

# Local interpretable explanations
explainer = lime.tabular.LimeTabularExplainer(
    X_train.values,
    feature_names=feature_names,
    class_names=['No Disease', 'Disease'],
    mode='classification'
)
```



Data Protection

- No sensitive data is stored permanently
- **Input validation** prevents injection attacks
- **Session management** ensures data privacy
- **HTTPS** recommended for production deployment

Model Security

```
python
# Input validation example
def validate_input(age, sex, cp, trestbps, chol):
  if not (1 \le age \le 120):
    raise ValueError("Age must be between 1 and 120")
  if sex not in [0, 1]:
    raise ValueError("Sex must be 0 or 1")
  # Additional validations...
```

Performance Monitoring

Model Drift Detection

```
python
def detect_drift(reference_data, current_data, threshold=0.05):
  """Detect data drift using statistical tests"""
  from scipy.stats import ks_2samp
  for column in reference_data.columns:
    statistic, p_value = ks_2samp(
       reference_data[column],
       current_data[column]
    if p_value < threshold:
       print(f"Drift detected in {column}: p-value = {p_value}")
```

Model Retraining Pipeline

python

```
def retrain_model_if_needed(current_performance, threshold=0.05):
    """Retrain model if performance degrades"""
    if current_performance < threshold:
        print("Retraining model due to performance degradation")
        new_model = train_new_model()
        return new_model
    return current_model</pre>
```

6 Future Enhancements

Planned Features

Deep Learning Models (Neural Networks)
Ensemble Methods (Voting, Stacking)
AutoML Integration (Auto-sklearn, TPOT)
Real-time Model Updates
A/B Testing Framework
Mobile App (React Native/Flutter)
API Endpoints (FastAPI/Flask)
Docker Containerization
Kubernetes Deployment
MLOps Pipeline (MLflow, Kubeflow)

Experimental Features

- Explainable AI with SHAP and LIME
- Federated Learning for distributed training
- **Edge Computing** deployment
- Time Series Analysis for longitudinal studies

Acknowledgments

Contributors

- Claude AI Main Development
- UCI ML Repository Dataset Provider
- **Streamlit Team** Web Framework
- Scikit-learn ML Library

Inspiration

This project was inspired by the need for accessible healthcare AI tools and demonstrates best practices in machine learning pipeline development.

License

This project is licensed under the MIT License - see the <u>LICENSE</u> file for details.

MIT License

Copyright (c) 2025 Heart Disease ML Pipeline

Permission is hereby granted, free of charge, to any person obtaining a copy of this software and associated documentation files (the "Software"), to deal in the Software without restriction, including without limitation the rights to use, copy, modify, merge, publish, distribute, sublicense, and/or sell copies of the Software, and to permit persons to whom the Software is furnished to do so, subject to the following conditions:

The above copyright notice and this permission notice shall be included in all copies or substantial portions of the Software.

THE SOFTWARE IS PROVIDED "AS IS", WITHOUT WARRANTY OF ANY KIND, EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO THE WARRANTIES OF MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE AND NONINFRINGEMENT. IN NO EVENT SHALL THE AUTHORS OR COPYRIGHT HOLDERS BE LIABLE FOR ANY CLAIM, DAMAGES OR OTHER LIABILITY, WHETHER IN AN ACTION OF CONTRACT, TORT OR OTHERWISE, ARISING FROM, OUT OF OR IN CONNECTION WITH THE SOFTWARE OR THE USE OR OTHER DEALINGS IN THE SOFTWARE.

Support & Contact

Getting Help

Issues: GitHub Issues

• **Discussions**: GitHub Discussions

Documentation: Project Wiki

Citation

If you use this project in your research, please cite:

```
bibtex

@software{heart_disease_ml_pipeline,
    title={Heart Disease Prediction ML Pipeline},
    author={Claude Al},
    year={2025},
    url={https://github.com/yourusername/heart-disease-ml-pipeline}
}
```

```
<div align="center">
★ Star this repository if you found it helpful! ★
✓ Get Started •  View Demo • Read Docs • Contribute
Made with ♥ for the healthcare Al community
</div>
```