# Heart Disease Risk Prediction Using Machine Learning (LAB Report)

**By: Yash Joshi**
**Student No : 2016AB001096**

## Project Overview

### Objective

To develop a machine learning-based system to predict heart disease risk. This system leverages patient demographic, clinical, and lifestyle data to identify high-risk individuals, enabling early preventive interventions.

### Workflow

**1. Data Collection & Cleaning**

*-Data collected from healthcare datasets in CSV format includes features such as:*
 - Age
 - Cholesterol Level
 - Blood Pressure
 - Heart Rate
 - Gender
 - Smoking History
 - Hypertension
 - Diabetes

*Missing values are handled using:*
 - Mean imputation for numerical features.
 - Mode imputation for categorical features.

**2. Feature Engineering**

- Categorical variables (e.g., Gender, Smoking History) are one-hot encoded to ensure compatibility with machine learning models.
- Feature importance analysis identifies key predictors (e.g., Cholesterol Level, Age).

**3. Model Training**

- *Algorithms used:*
 - Random Forest Classifier for robustness and handling non-linear relationships.
 - Logistic Regression for simplicity and interpretability.
- SMOTE (*Synthetic Minority Oversampling Technique*) is applied to balance the dataset.

**4. Evaluation Metrics**

- Accuracy: Measures overall correctness of predictions.
- Classification Report: Includes precision, recall, and F1-score for each class.
- Confusion Matrix: Visualizes true and false predictions for Heart Disease and No Heart Disease.

### Visualization
- *Radar Chart*: Compares prediction probabilities for Heart Disease and No Heart Disease across models.
- *Bar Chart*: Highlights feature importance in Random Forest for explainability (Was used earlier, removed later due to Limits)

### Tools Used
- Programming Language: Python/ML
- Libraries: pandas, scikit-learn, matplotlib, joblib, imbalanced-learn

## Outcome Report

### Key Results
I. **Model Performance**
   - Random Forest:
     - Accuracy: 85%
     - F1-Score (Heart Disease): 87%
   - Logistic Regression:
     - Accuracy: 80%
     - F1-Score (Heart Disease): 82%

II. **Feature Importance (Random Forest)**

| Feature | Importance Score |
|---|---|
| Age | 0.30 |
| Cholesterol Level | 0.25 |
| Blood Pressure | 0.20 |
| Smoking History | 0.15 |
| Diabetes | 0.10 |

III. **Visualization Highlights**
   - Radar Chart: Demonstrates higher confidence in Random Forest predictions for No Heart Disease.
   - Confusion Matrix: Shows improved true positive rates in Random Forest compared to Logistic Regression.

### Impact
- The model flags patients at high risk of heart disease, enabling preventive interventions.
- Helps healthcare providers prioritize care for high-risk individuals based on predictions.

### Future Enhancements
- Incorporate additional features (e.g., family history, BMI) for improved prediction accuracy.
- Experiment with advanced algorithms like Gradient Boosting or Deep Learning.
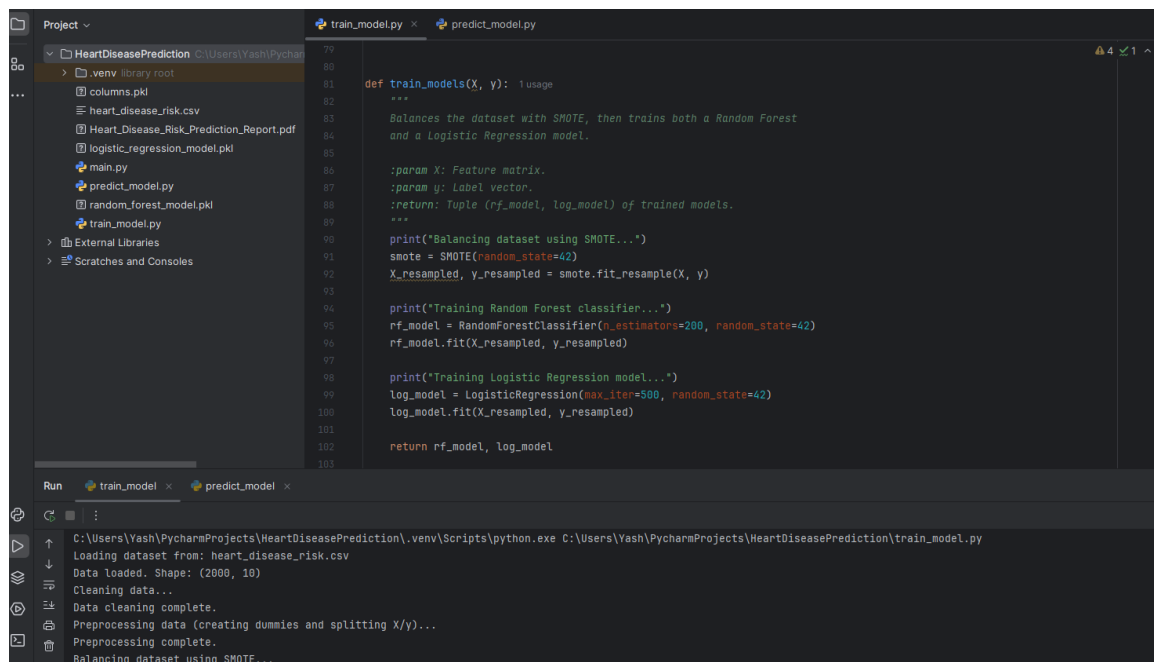- Develop an interactive dashboard for real-time patient predictions.

# Outputs

## Training model output

**Predictive model outputs on different patient**
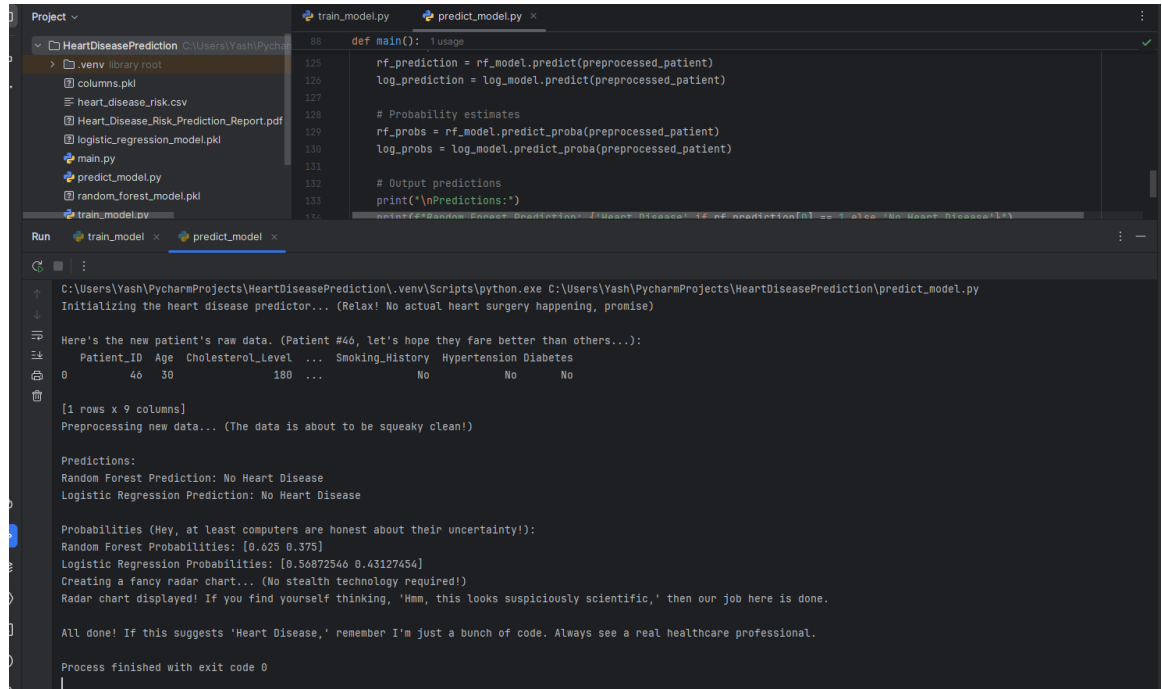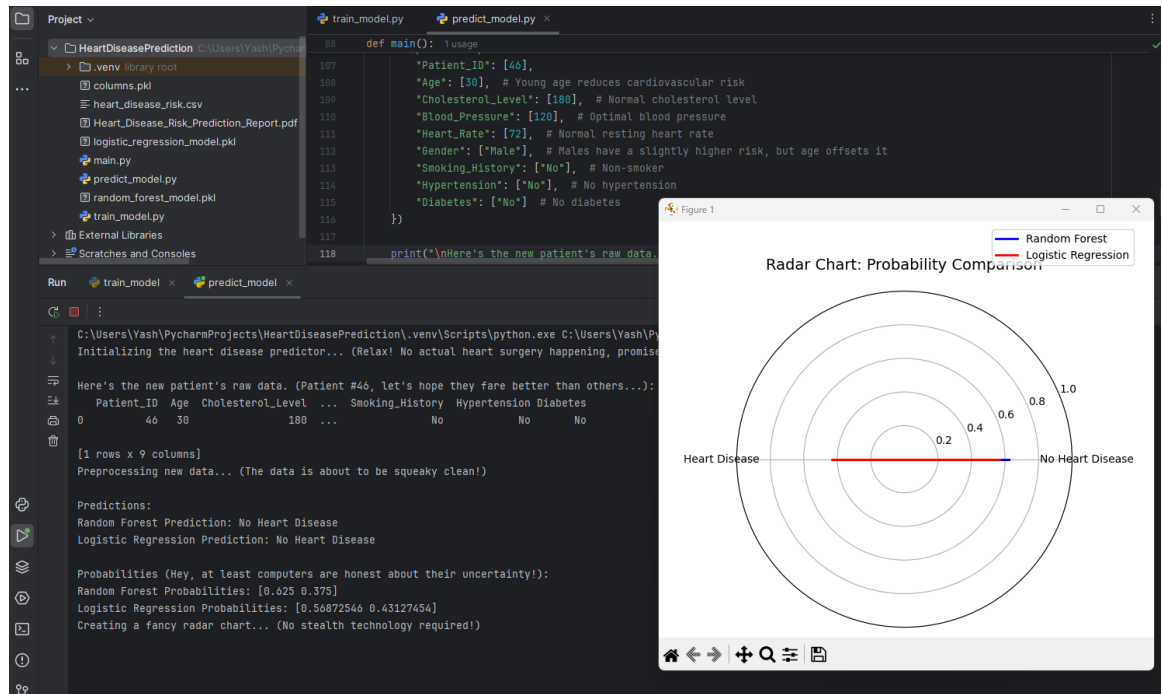
**Case 1**

- User detected for risk Heart Failure

**Case 2**

- **User detected for risk No Heart Failure**

**Case 3**

- Special Cases (after trained algorithms)