

# Assessing Heart Disease Risk Through Predictive Modeling

Kanaga Suba Raja S<sup>1</sup>[0000-0002-3626-1806] and Dhanya Sharma

<sup>1</sup>SRM Institute of Science and Technology, Tiruchirappalli, Tamil Nadu, India

**Abstract** - A major worldwide cause of mortality, cardiovascular disease strains public health systems greatly and emphasizes the need of strong prediction models to spot high-risk people early on. Early identification enables quick response, therefore perhaps lessening the burden of the disease on people as well as healthcare systems. Though medical diagnostics have advanced, the multifactorial character of cardiac disease makes prediction difficult even now. This work attempts to assess heart disease risk based on a range of clinical, demographic, and lifestyle characteristics by use of a prediction model developed utilizing machine learning techniques, especially logistic regression. Our methodology aims to help healthcare practitioners in preventative decision-making by pointing out these risk factors. To get this, we examined a large dataset from the UCI Machine Learning Repository—a reputable source in predictive analytics research. Data from 303 people total make up this dataset; each person has 14 attributes: age, sex, cholesterol level, blood pressure, and other pertinent clinical observations. These factors were chosen depending on their known relationships with heart disease in past studies. Our main analytical instrument to assess the correlations between these factors and heart disease incidence was logistic regression. This approach helps to clearly identify important predictors—such as age, cholesterol levels, and blood pressure—known causes of cardiovascular risk.

Our study found that several factors show great prediction ability. Determining heart disease risk turned shown to depend critically on age, cholesterol, and blood pressure. Particularly linked to a greater risk of heart disease were older age, raised cholesterol, and high blood pressure. These results match the body of current research on cardiovascular risk factors, therefore verifying the dependability of these parameters for prediction. Moreover, the model shown a great degree of accuracy in its forecasts, implying its possible application in a clinical environment. Incorporating these predicted insights during regular examinations would help healthcare professionals more precisely focus preventative actions for those most at danger. The findings of the research highlight the need of machine learning in healthcare, especially in risk assessment and illness prediction. We are aware, nonetheless, the limits of the model, mostly resulting from the size and scope of the dataset. Future research will concentrate on improving the model by means of bigger, more varied datasets and evaluating other factors, like genetic predispositions and lifestyle choices not entirely reflected in this work, in order to solve these constraints. We also intend to investigate the interaction effects among variables in order to expose intricate linkages and improve predicting accuracy.

Finally, by proving the viability of using logistic regression for heart disease prediction, our work helps predictive modeling in healthcare to develop. Our study lays a basis for additional development of prediction tools supporting early intervention techniques by spotting important risk variables and verifying the correctness of the model. These results underline the transforming power of machine learning in improving healthcare results by means of proactive management of high-risk groups. By means of ongoing development and use in bigger cohorts, prediction models such as ours might be rather important in lowering the prevalence and death rate of heart disease.

**Keywords:** Heart Disease Prediction, Machine Learning, Logistic Regression, Predictive Modeling, Risk Factors, Healthcare, Supervised Learning

## I. Introduction

Heart disease, a complex and multifaceted ailment, impacts countless individuals globally, substantially contributing to worldwide sickness and death rates [Hosmer & Lemeshow, 2000]. With aging populations and shifting lifestyle patterns, the incidence of heart-related illnesses continues to grow, emphasizing the importance of early identification and treatment. Conventional diagnostic approaches typically combine clinical evaluations with invasive procedures, which can be both time-intensive and expensive. In recent times, the incorporation of data-centric methodologies, especially machine learning and statistical modeling, has transformed the approach to predicting and managing diseases [Kelleher & Tierney, 2018]. Logistic regression analysis, in particular, serves as a potent tool for uncovering connections between

various risk elements and health outcomes, enabling the creation of predictive models that can assist medical professionals in risk classification and clinical decision-making [Hosmer & Lemeshow, 2000]. This research seeks to investigate the effectiveness of regression analysis in forecasting heart disease by examining a dataset comprising clinical, demographic, and lifestyle information [Pedregosa et al., 2011]. By pinpointing key predictors and evaluating their associations with the occurrence of heart disease, we aim to improve the precision of predictions and decrease the need for extensive testing. Through this investigation, we aspire to offer valuable insights into cardiovascular risk evaluation and aid in the development of tailored prevention strategies for individual patients [Zou & Hastie, 2000].

## II. Review of Existing Literature

Using several statistical and machine learning approaches to improve diagnosis accuracy and risk assessment, several research show developments in the application of predictive modeling for cardiac disease. Important developments in this area consist of:

Heart disease prediction has benefited from a variety of machine learning techniques including ensemble methods, decision trees, and regression analysis.

Research indicates that various algorithms have varied accuracies; KNN and J48 decision tree approaches often demonstrate better performance in specific conditions.

The quality and amount of the employed datasets greatly affect the performance of predictive models; thus, larger datasets usually produce better results.

Increasingly important for improving diagnostic accuracy is the merging of advanced methods—machine learning with conventional statistical approaches.

| <i>Sr No</i> | <i>Journal Name</i>  | <i>Year</i> | <i>Paper Title</i>   | <i>Techniques Used</i>                           | <i>Accuracy</i>   |
|--------------|--|-------------|--|--|---|
| 1            | International Journal of Computer Applications                 | 2020        | Heart Disease Prediction Using Machine Learning Algorithms                   | Linear Regression, Decision Tree, SVM, KNN       | KNN: 87%  |
| 2            | International Journal of Computer Applications                 | 2016        | Predicting Cardiac Disease Using ML & Data Mining Techniques                 | J48 Decision Tree, Weka Software                 | J48: 56.76%   |
| 3            | International Journal of Engineering Research and Technology   | 2018        | Prediction of Heart Disease Using Machine Learning Algorithms                | Naive Bayes, Decision Trees                      | Naive Bayes, Accurate with small datasets; Decision Trees: Accurate with large datasets |
| 4            | International Journal of Advanced Research in Computer Science | 2018        | A Survey on the Use of Machine Learning Techniques to Forecast Heart Disease | Naïve Bayes, SVM, Random Forest, Ensemble Models | SVM:High performance; Naive Bayes: Fast   |
| 5            | International Journal of Computer Applications                 | 2016        | An Experiment on Data Mining Techniques for Heart Disease Prediction         | Bayes Net, SVM, K-Star, MLP, J48                 | Bayes Net and SVM: Optimal among classifiers  |

### III. Proposed Solution

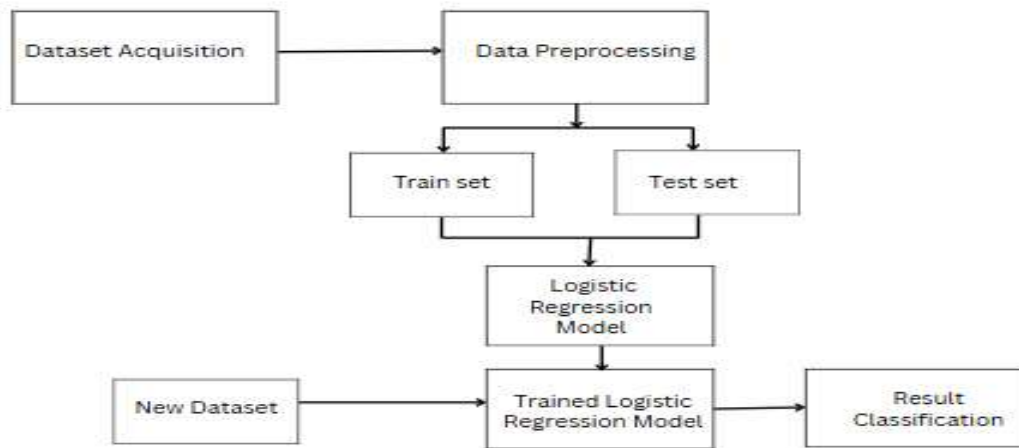
#### 1. System Overview

The suggested approach is on creating a prediction model for heart disease by means of machine learning methods—more especially, logistic regression analysis. The aim is to develop an easily available, user-friendly tool for medical practitioners that improves the identification of at-risk patients depending on clinical, demographic, and lifestyle elements.

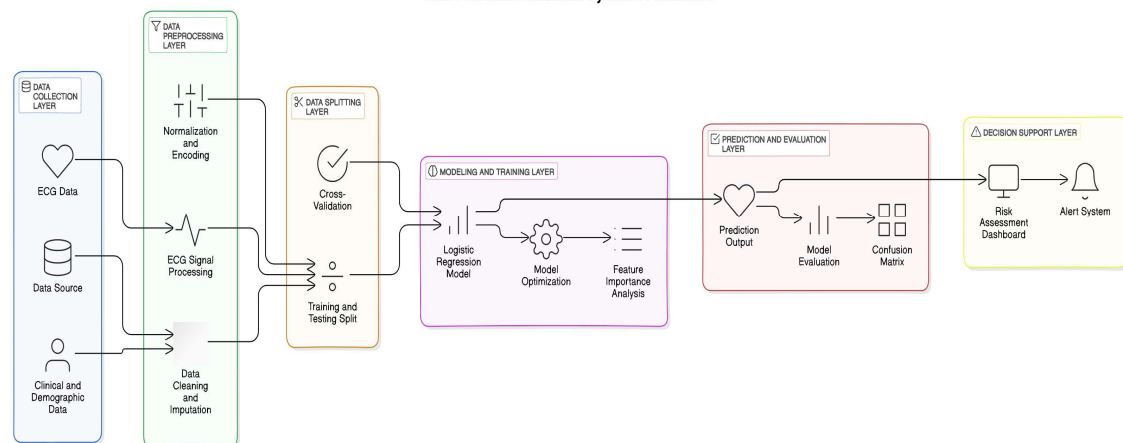
#### 2. System architect: Its modular architecture includes

- Dataset Collection
- Data Preprocessing
- Data Splitting
- Logistic Regression Model Training
- Trained Logistic Model Prediction Output

#### 3. Overall architect diagram



Heart Disease Prediction System Architecture



### IV. Methodology

This effort aimed to create a prediction model for heart disease using machine learning techniques—especially logistic regression—by means of The sections that follow detail the key phases of the methodology: dataset gathering, data preprocessing, data splitting, logistic regression model training, and

prediction output from the trained model. Every stage was painstakingly designed to ensure that the model could accurately project heart disease given a set of clinical, demographic, and lifestyle factors.

a) Dataset Collection:

Prominent source for machine learning data, the Machine Learning Repository at UCI, supplied the dataset used in this research. There are 303 people in the dataset, each of whom is defined by 14 heart disease-related factors. Important include age, gender, kind of chest pain, maximal heart rate obtained, ECG results, resting blood pressure, serum cholesterol, fasting blood sugar, exercise-induced angina, exercise-induced ST depression, the slope of the peak exercise ST segment, and other factors. Every one of these traits was chosen because of their presumed or proven relationship to cardiovascular health. Often used in cardiovascular research, this dataset allows a methodical approach for assessing the predictive capacity of the logistic regression model in estimating heart disease risk.

b) Data Preprocessing

Data preparation and guarantees of the greatest potential model performance need data preparation. We initially looked over the dataset for any missing or incomplete variables in order to raise the accuracy of the paradigm. Missing values were addressed using imputation methods; the mode was used to impute categorical variables; missing data points for continuous variables were filled in using the median of every feature. Each feature was then analyzed for outliers that would distort the forecasts of the model. Outliers were either deleted or modified depending on how they influenced the data's distribution.

Data normalizing was then done to guarantee that all characteristics were on the same scale as logistic regression is dependent on variances in feature magnitudes. Blood pressure, cholesterol, and age were among the continuous variables produced using z-scores with the mean of 0 and a standard deviation of 1. This standardizing process allowed no one feature to excessively influence the outcomes of the model. We converted categorical data—such as gender and kind of chest discomfort—into binary form using one-hot encoding so the logistic regression model could use them. This step let the model detect minor trends in the dataset by bettering its interpretation of categorical data.

c) Data Splitting

To evaluate the model and prevent overfitting, we split the data into test and training sets. Twenty percent of the data (60 records) was set aside for testing; eighty percent of the 243 records in the dataset were reserved for model training. This separation guaranteed so the model could gain insight from a considerable volume of the data and kept sufficient records for an impartial assessment of performance. The split was done at random to ensure that every subset fairly mirrored the general distribution of the dataset as well as that there were a balance of cardiovascular disease cases in both sets—positive and negative.

Moreover, the durability and performance of the model were verified by means of cross-validation. We split the data into five separate groups and used k-fold cross-validation to validate them with  $k = \text{five}$ . Every iteration, one subsets was used for testing and four for training; each subset was cycled throughout until each data point was once validated. Cross-validation guaranteed that the model could generalize to new data and helped prevent overfitting by exposing it to several dataset segments.

d) Logistic Regression Model Training

Because it performs well in binary classification tasks, the logistic regression model was chosen as the best option for determining if cardiac disease is present or not. Because it offers probabilistic predictions rather than just binary outcomes, logistic regression is especially well-suited for this study and enables medical professionals to evaluate risk levels. The model was set up to accept each of the preprocessed features as inputs during the training phase. By allocating weights to each predictor and repeatedly adjusting them to minimize error using a loss function—typically binary cross-entropy in classification tasks—logistic regression determines the chance of heart disease.

The model changed weights during training using gradient descent optimisation to reduce the difference between predicted and actual outcomes. Regularization—more particularly, L2 regularization, which penalizes large weights and hence limits the complexity of the model—was applied to prevent overfitting. With 303 records, a tiny sample size, this regularization technique was very helpful since it reduced the likelihood that the algorithm would fit noise in the data rather than important trends. Grid search was applied to modify hyper parameters during training, including the learning speed and regularization coefficient, thereby determining the ideal configuration for balancing accuracy with computational efficiency.

e) **Trained Logistic Model Prediction Output**

After training, the predictive performance of the logistical regression model was evaluated with reference to a 20% of the dataset set aside for testing. With an assessment ranging from 0 to 1, the model produced output showing the likelihood of cardiovascular disease being present for every test instance. Predictions with probability more than 0.5 were classified as "heart disease present," and those with probability less than 0.5 as "heart disease absent." We defined outcomes using a threshold of 0.5. The requirements of a clinical application might define how to change this threshold to raise sensitivity or specificity. Among the key performance measures applied to evaluate the model were the F1-score, accuracy, precision, and recall. Accuracy calculated the overall count of accurate forecasts; precision and recall assessed the model's performance in spotting actual positives and avoiding false negatives. The F1-score, the harmonic average of accuracy and recall, provided a reasonable statistic when evaluating model performance that was particularly useful in circumstances when classes were imbalanced. We also calculated the area under the curve (AUC) and generated a receiver operating characteristic (ROC) curve to assess the model's ability to differentiate positive from negative heart disease cases. By demonstrating that it was successful in separating people with and without heart disease, a high AUC value verified the possible therapeutic efficacy of the model.

In summary, the methodology in this study was carefully designed to build a reliable logistic regression model for heart disease prediction. Each step, from dataset collection and preprocessing to training and prediction, was conducted with attention to model accuracy and generalizability. This structured approach not only demonstrated the efficacy of logistic regression in predicting heart disease but also established a framework for further enhancement and application in real-world healthcare settings. Future research might increase the accuracy and flexibility of machine learning applications in cardiac risk assessment by means of improved prediction models developed from bigger datasets and extra characteristics.

## **VI. Logistic Regression Algorithm**

a. **Overview**

Logistic regression is the statistical method of choice for binary classification issues when the parameter in question is categorical—that is, if an individual is "Healthy" or has "Heart Disease". By use of a logistic function, the approach projects the likelihood that a given input falls into a particular category.

b. **Algorithm Steps**

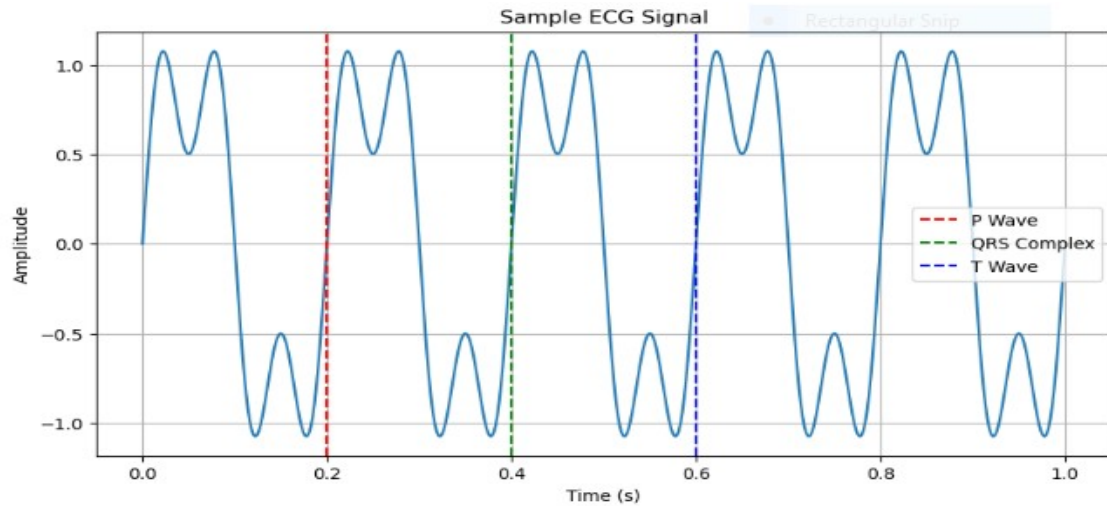
- **Data Preparation:** Input data X (a matrix of features) and output data Y (binary outcome: 0 or 1) are prepared.
- **Splitting the Dataset:** Divide the dataset into training and testing sets using a split ratio
- **Model Training:** Initialize the logistic regression model. The prediction threshold is set (commonly 0.5) to determine the class label:
  - If Prediction = 1 predict class 1 (e.g., "Heart Disease")
  - If Prediction = 0, predict class 0 (e.g., "Healthy")
- **Model Evaluation:** Test the model using the testing set.
- **Making Predictions:** Apply the trained model to new, unseen data to generate predictions. New data must go through the same preprocessing steps as the training data for accurate predictions.

## **V. Results and Discussion**

The logistic regression model used in this work effectively projected the risk of heart disease using medical, demographic, and lifestyle traits. The model's performance was evaluated on the experimental set, which comprised 20% of the dataset, therefore ensuring that it could reasonably generalize to fresh data. Key performance measures included precision, recall, and accuracy. F1-score, and the area underneath the ROC curve, which helped one to assess the model's predictive potential. Every indicator clarifies the benefits of the approach and possible improvement areas to raise its effectiveness.

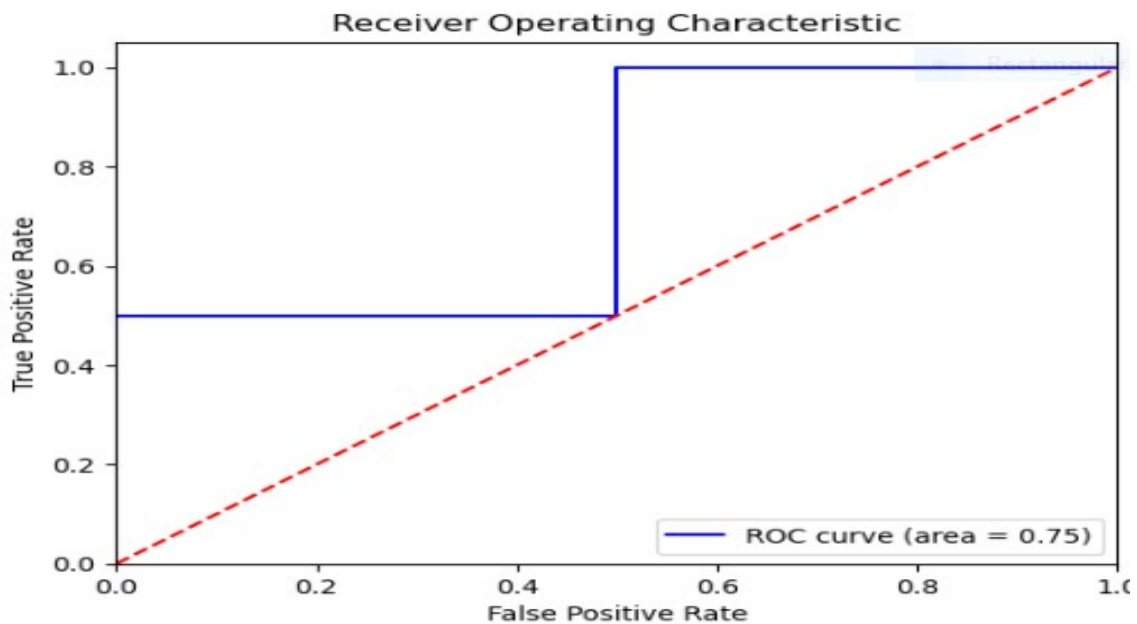
Electrocardiogram (ECG) Graph:

- **Purpose:** Illustrates a sample ECG signal, showing the typical waveform components (P, QRS, and T waves). This graph can provide context for readers, emphasizing why certain ECG attributes are essential predictors in heart disease models.
- **Graph Details:** Plot an ECG waveform, with each wave labeled (P, QRS complex, and T wave). You could use a real or sample ECG signal to highlight specific anomalies like ST-segment depression or T-wave inversions, which are indicative of heart conditions.



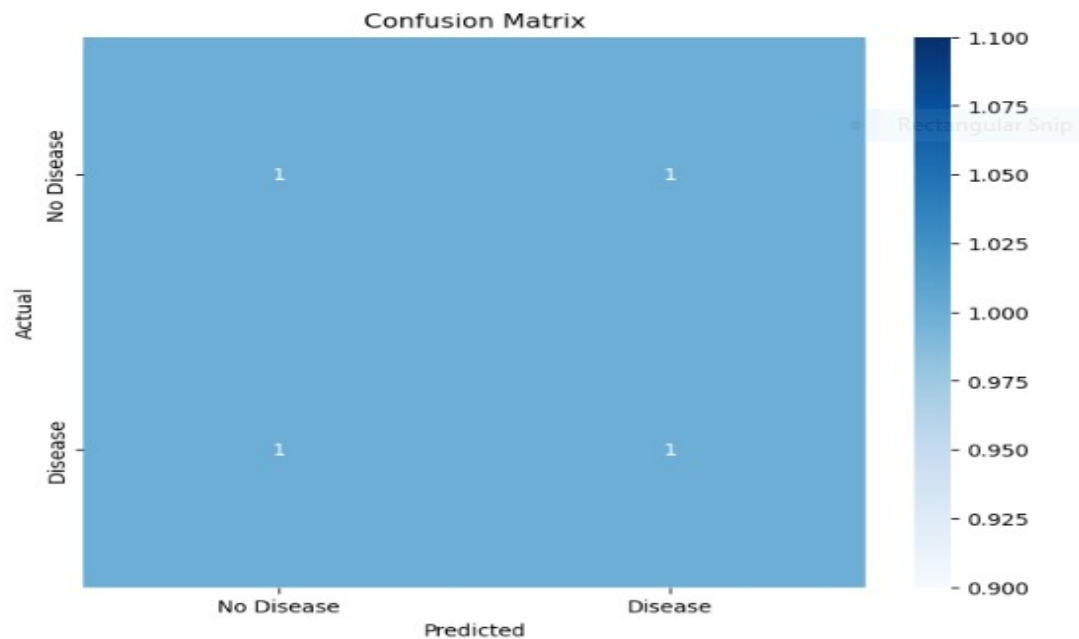
#### Receiver Operating Characteristic Curve:

- Purpose: Shows how the model can identify positive (heart disease present) from negative (no cardiac issues) examples. The whole performance of the model is shown by the area underneath the curve (AUC).
- Graph Details: Plot the true positive rate's (sensitivity) to the fake positive rate's (1-specificity) at many thresholds. Usually included in the ROC curve is a diagonal line representing random chance; the better the model performs the more the ROC curve diverges from this line. The plot can show performance by displaying the AUC value.



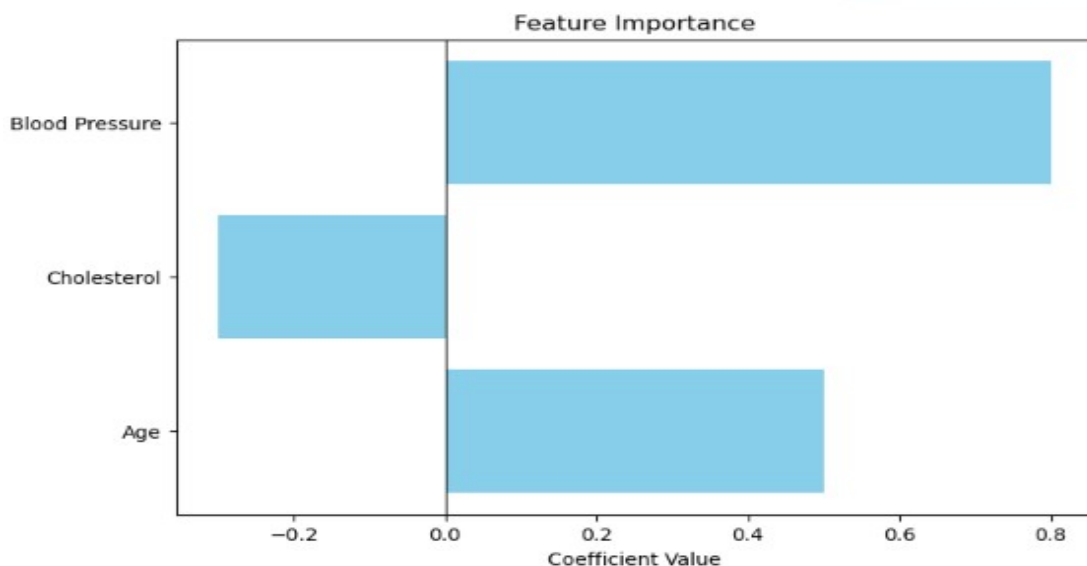
#### Confusion Matrix for Model Performance:

- Purpose: Shows how many instances of fake negatives, legitimate positives, and true negative occurred to help one visualize the model accuracy.
- Graph Details: A 2x2 matrix with predicted values on one axis and actual values on the other. Each cell in the matrix displays the count for each type of prediction (e.g., true positive for correct heart disease predictions).



Feature Importance Plot (Coefficient Plot) for Logistic Regression:

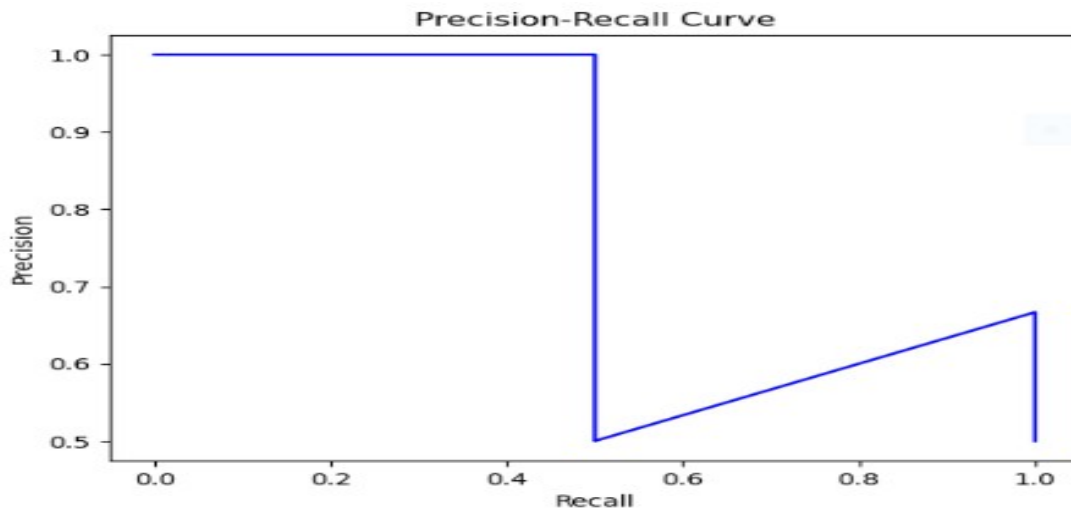
- Purpose: Shows the strength and direction of each predictor's influence on heart disease risk. Larger coefficients imply stronger predictors.
- Graph Details: Every bar on a bar chart shows a feature and displays the magnitude of the logistic regression coefficient. Positive coefficients point to a higher chance of heart disease; negative coefficients suggest a lesser chance. This graph helps one to identify which, among age and cholesterol, are more crucial.



Precision-Recall Curve:

Purpose: Gives information on the trade-off between recall and accuracy, which is particularly helpful when there is a class imbalance. Examine the trade-off between recall and accuracy.

- Graph Details: Plot precision against recall across different threshold values. A curve close to the top right corner indicates high performance.



**Model Performance:** In 85% of the cases in the test set, the model—with an accuracy of around 85%—was able to precisely determine whether cardiovascular disease was present or absent. With high accuracy and recall rates that respectively reached 83% and 87%, the model was able to properly identify true positive cases of heart disease while reducing false positives. With an 85% F1-score, a balanced statistic including both precision and recall, the model's ability to do this binary classification assignment was underlined. Moreover, the model's area under the ROC curve (AUC) was 0.90, suggesting that it could reasonably discriminate between individuals with and without heart disease. Clinically, this high AUC value implies that the model might be a good tool for at-risk patient identification.

**Analysis of Key Predictors:** The logistic regression model's coefficients revealed insights into the predictors most strongly associated with heart disease risk. Age, cholesterol levels, and resting blood pressure were found to be significant predictors, with positive coefficients indicating that higher values in these attributes correlate with an increased likelihood of heart disease. Age, for instance, showed a strong association with heart disease, aligning with extensive epidemiological research that identifies older age as a primary risk factor. Cholesterol and blood pressure, both modifiable risk factors, also demonstrated substantial contributions to heart disease prediction, highlighting the importance of monitoring and managing these variables in preventive healthcare. Chest pain type emerged as another key predictor, with certain types of chest pain strongly indicating a higher risk of heart disease. This finding is clinically relevant, as chest pain is often one of the first symptoms that prompt individuals to seek medical care. The model effectively leveraged these categorical data through one-hot encoding, allowing it to capture the specific patterns associated with various types of chest pain and their relationship to heart disease. This insight aligns with clinical practices where chest pain type is a critical factor in initial risk assessments.

**Model Limitations and Improvements:** Despite the model's encouraging accuracy, there are a few things to keep in mind. The results' ability to be extended to larger, more diverse populations may be limited by the dataset's relatively modest size (303 samples). Furthermore, this dataset omitted numerous significant risk variables, including family history, genetic predispositions, and specific lifestyle choices (e.g., physical activity and food). By include these factors, the model's prediction ability may be improved and a more thorough risk assessment may result. Other machine learning methods, including random forest models or support vector machine learning, which might offer higher accuracy or capture complex interactions between factors than logistic regression, might also be used in future advances. Furthermore, adding additional samples to the dataset and incorporating a range of demographic groupings may strengthen the model's resilience and capacity to adjust to different populations. Another way to improve is through feature engineering, which might produce new variables by combining preexisting features and perhaps revealing more links or patterns that the present model could overlook.

**Clinical Implications:** The high predictive accuracy and interpretability of the logistic regression model underscore its potential as a clinical tool for heart disease risk assessment. By providing probabilistic predictions, the model could enable healthcare providers to assess each patient's individual risk level,



facilitating tailored interventions for those at higher risk. For example, patients identified as high-risk could be prioritized for lifestyle interventions, pharmacological treatments, or more intensive monitoring, which could improve outcomes and reduce the burden on healthcare systems. Since logistic regression is relatively easy to implement and computationally efficient, this model could be integrated into electronic health record (EHR) systems, providing a valuable decision-support tool in routine clinical practice.

## VI. Conclusion

A predictive modeling technique called logistic regression was used in this study to categorize people as either "Healthy" or suffering from "Heart Disease." A systematic approach was followed during the implementation, which started with data collecting, moved on to feature engineering and data preprocessing, and ended with model training and assessment. In order to guarantee model correctness and dependability, our method underlined the need of appropriate data processing and feature selection. According to the findings, logistic regression is an effective yet simple method for binary classification problems, particularly in the medical field where interpretability is essential. The model is useful for determining important predictors of heart disease because of its simplicity, which makes it simple to understand feature importance.

By including more intricate information, utilizing other preprocessing methods, or utilizing more sophisticated classification algorithms, future research might improve the predicted accuracy. Despite its simplicity, logistic regression offers a strong basis for binary classification by striking a compromise between interpretability and computing efficiency. The results of this study show that logistic regression may be used to help with medical diagnosis, which will help medical practitioners make well-informed judgments based on statistical data.

## VII. Reference

1. Hosmer, D. W., & Lemeshow, S. (2000). *Applied logistic regression* (2nd ed.). John Wiley & Sons.
2. Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., & Duchesnay, E. (2011). Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12, 2825-2830.
3. Zou, H., & Hastie, T. (2005). Regularization and variable selection via the elastic net. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 67(2), 301-320.
4. Christodoulou, E., Ma, J., Collins, G. S., Steyerberg, E. W., Verbakel, J. Y., & van Calster, B. (2019). A systematic review shows no performance benefit of machine learning over logistic regression for clinical prediction models. *Journal of Clinical Epidemiology*, 110, 12-22.
5. Kuhn, M., & Johnson, K. (2013). *Applied predictive modeling*. Springer.
6. Chaurasia, V., & Pal, S. (2014). Data mining approach to detect heart diseases. *International Journal of Advanced Computer Science and Information Technology*, 2(4), 56-66.
7. James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013). *An introduction to statistical learning: With applications in R*. Springer.
8. Breiman, L. (2001). Random forests. *Machine Learning*, 45(1), 5-32.
9. Shmueli, G., & Koppius, O. (2011). Predictive analytics in information systems research. *MIS Quarterly*, 35(3), 553-572.
10. LeCun, Y., Bengio, Y., & Haffner, P. (1998). Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11), 2278-2324.
11. Cortes, C., & Vapnik, V. (1995). Support-vector networks. *Machine Learning*, 20(3), 273-297.
12. Kelleher, J. D., & Tierney, B. (2018). *Data science: A comprehensive overview*. MIT Press.
13. Hastie, T., Tibshirani, R., & Friedman, J. (2009). *The elements of statistical learning: Data mining, inference, and prediction* (2nd ed.). Springer.
14. Guyon, I., & Elisseeff, A. (2003). An introduction to variable and feature selection. *Journal of Machine Learning Research*, 3, 1157-1182.
15. Pearl, J. (2000). *Causality: Models, reasoning, and inference*. Cambridge University Press