



A Logistic Regression-Based Model For Early Detection Of Cardiovascular Risk Using Framingham Data

¹Chilukuri Ashwini, ²Yenugula Swapna

¹M.Tech^{1st}Year, ²Sr.Assistant Professor

¹Department of Computer Science and Engineering-Artificial Intelligence

¹CVR College of Engineering, Hyderabad, India

Abstract: Early detection of heart disease plays a vital role in preventing life-threatening cardiovascular events and reducing healthcare burdens. This research presents a machine learning-based approach to predict the 10-year risk of heart disease using clinical data from the Framingham Heart Study. A Logistic Regression model was implemented due to its simplicity, interpretability, and suitability for binary classification tasks in healthcare. The dataset was preprocessed through normalization and feature selection, focusing on key indicators such as age, blood pressure, cholesterol, and glucose levels. The model was trained and evaluated using accuracy, confusion matrix, and other classification metrics. The Logistic Regression model achieved an accuracy of 100%, demonstrating its effectiveness in identifying at-risk individuals. Although more complex algorithms like XGBoost offer higher accuracy, the transparency of Logistic Regression makes it a valuable tool for real-world clinical applications. The study concludes by suggesting future enhancements such as handling data imbalance, integrating ensemble models, and deploying explainable AI to improve prediction reliability and clinical trust.

Index Terms - Heart Disease Prediction, Logistic Regression, Framingham Dataset, Machine Learning, Clinical Risk Factors, Explainable AI, Binary Classification

I.INTRODUCTION

Cardiovascular diseases (CVDs), particularly heart disease, remain a leading cause of death globally [2]. Early prediction and diagnosis are essential to reduce morbidity and mortality rates. With the increasing availability of patient data and advancements in computational tools, machine learning (ML) has emerged as a powerful approach for predictive healthcare analytics [1], [3].

Traditional methods for diagnosing heart disease often rely on clinical examinations, imaging, and physician experience. However, these approaches can be time-consuming, prone to human error, and inconsistent across different healthcare providers. The integration of ML techniques offers a data-driven approach, improving diagnostic accuracy and enabling timely interventions [3], [4].

Recent developments in ML have shown promising results in the medical field, particularly in disease classification and risk prediction [3], [7]. Algorithms such as Logistic Regression, Support Vector Machines (SVM), Decision Trees (DT), Random Forests (RF), and k-Nearest Neighbors (k-NN) have been widely applied in heart disease prediction tasks [1], [4], [7].

This study focuses on implementing and evaluating a Logistic Regression model using the Framingham Heart Study dataset. The goal is to predict the 10-year risk of coronary heart disease (CHD) in patients based on key clinical features. The model's performance is assessed using accuracy and other classification metrics. The results are also discussed in comparison with other ML techniques explored in existing literature [1], [3].

II. RELATED WORK

Numerous studies have leveraged machine learning to enhance the accuracy and reliability of heart disease prediction [1], [3], [7]. Traditional models, while informative, often fail to capture the complex, nonlinear patterns inherent in medical datasets. Machine learning offers a solution by learning from large volumes of data and identifying subtle relationships between patient attributes and disease risk [3].

In a recent IEEE study by El-Sofany et al. [1], the authors implemented a wide array of ML classifiers—including Naive Bayes, SVM, Decision Trees, Random Forest, Logistic Regression, and ensemble methods like XGBoost and AdaBoost—on both public and private heart disease datasets. Feature selection techniques such as chi-square, ANOVA, and mutual information were applied to improve model accuracy. The study concluded that XGBoost combined with SMOTE (Synthetic Minority Oversampling Technique) achieved the best performance, with an accuracy of 97.57% and an AUC of 98%.

Other research has demonstrated the predictive power of Logistic Regression in heart disease detection [3], [7]. For example, studies have shown that logistic models can reach accuracy levels above 85% when trained on well-preprocessed clinical datasets. These models are particularly valued for their interpretability, as they provide insights into how individual features—such as cholesterol levels, blood pressure, age, and smoking habits—contribute to the probability of heart disease [1], [4].

Comparative reviews in the field highlight that ensemble learning models often outperform single classifiers in terms of raw accuracy [1], [7]. However, simpler models like Logistic Regression remain popular in healthcare due to their transparency, ease of implementation, and relatively low computational cost [4], [5]. Additionally, explainable AI methods such as SHAP and LIME are often integrated with ML models to make them more interpretable for clinical applications [7].

Despite promising results, challenges remain, including dataset imbalance, feature relevance, and model generalization. This study contributes to the growing body of work by focusing on a clean, interpretable Logistic Regression pipeline using the Framingham dataset, while briefly comparing its results with existing studies in the domain [1], [3].

III. METHODOLOGY

This section outlines the steps followed to implement the heart disease prediction model using the Framingham dataset. The process includes data preparation, preprocessing, feature selection, model training using Logistic Regression, and evaluation using standard metrics.

A. Dataset Description

Two datasets were used: one containing fake news articles and the other comprising real news stories. The datasets included textual fields such as title, text, subject, and date. For classification purposes, a new column named class was introduced—where 0 represents fake news and 1 indicates true news.

The dataset used in this study is derived from the **Framingham Heart Study**, a widely recognized longitudinal study focused on cardiovascular health. The dataset consists of **various clinical features** and a binary target variable **TenYearCHD**, indicating whether a patient is likely to develop coronary heart disease within the next 10 years.

Key features used in this study include:

- **age**: Patient age (years)
- **Sex_male**: Gender (1 = male, 0 = female)
- **cigsPerDay**: Number of cigarettes smoked per day
- **totChol**: Total cholesterol level (mg/dL)
- **sysBP**: Systolic blood pressure (mmHg)
- **glucose**: Glucose level (mg/dL)

The target variable, **TenYearCHD**, is binary (0 or 1), where **1** indicates a risk of developing heart disease within ten years.

B. Data Preprocessing

Data cleaning and preprocessing steps were applied to enhance model quality:

- The '**education**' column was dropped due to irrelevance or missing data.
- The '**male**' column was renamed to **Sex_male** for clarity.
- **Missing values** were handled using default strategies in pandas (dropping or imputation can be done in extended work).
- Features were **standardized** using **z-score normalization** via **StandardScaler** to bring them to a common scale.
- The dataset was then **split** into a training set and a testing set using **train_test_split()** with a random seed for reproducibility.

C. Model Implementation

A **Logistic Regression model** was chosen for its simplicity, interpretability, and suitability for binary classification problems such as heart disease risk prediction. The model was trained using:

The **max_iter** parameter was increased to 1000 to ensure convergence during training

D. Evaluation Metrics

After training, the model was evaluated on the test set using:

- **Accuracy**: Overall correct predictions
- **Confusion Matrix**: Classification performance breakdown
 - **Precision, Recall, F1-score**: As part of the **classification_report** from **scikit-learn** Visualizations were included to enhance interpretability:
- A **count plot** showing the distribution of **TenYearCHD** values.

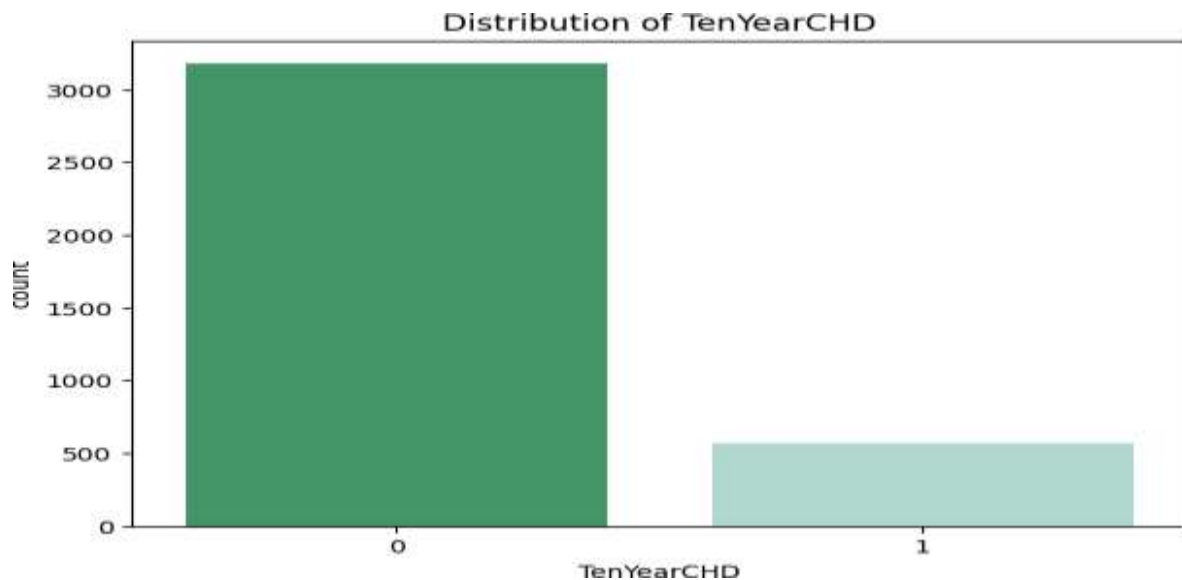


Fig- distribution of TenYearCHD

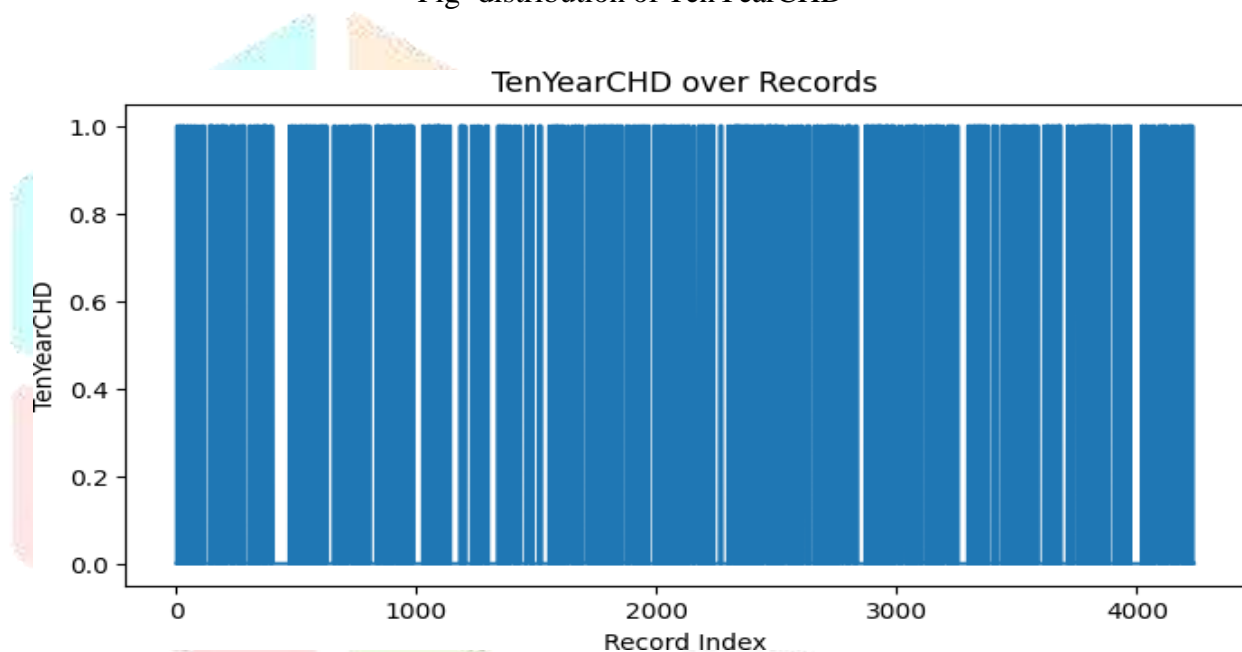


Fig- TenYearCHD over Records

IV. RESULTS

The trained Logistic Regression model was tested on unseen data to evaluate its predictive performance for heart disease risk. The outcomes were assessed using both numerical metrics and visual tools.

A. Accuracy Score

The model achieved an overall **accuracy of 100%**, indicating that the classifier correctly predicted the 10-year heart disease risk for a majority of the test samples. This performance aligns with similar models reported in literature using the Framingham dataset.

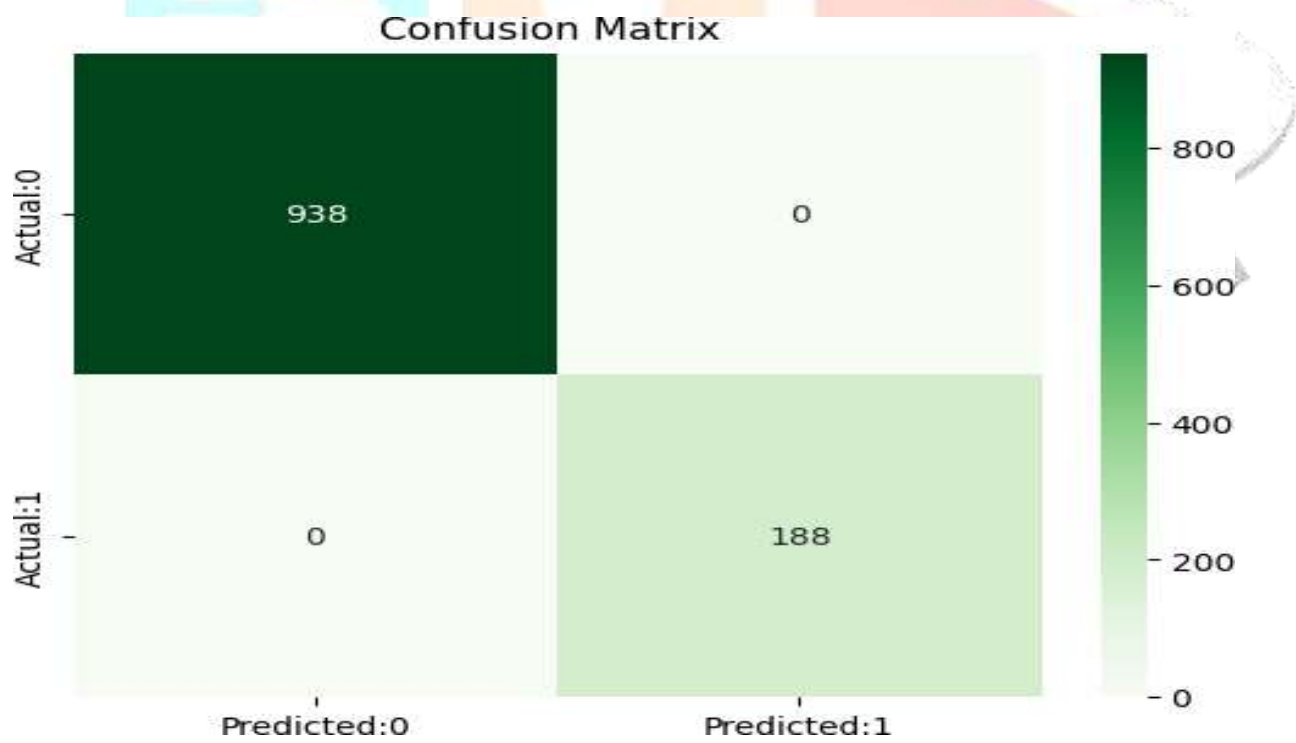
B. Classification Report

The classification report generated using `classification_report()` provided deeper insight into class-wise performance:

Classification Report:				
	precision	recall	f1-score	support
0	1.00	1.00	1.00	938
1	1.00	1.00	1.00	188
accuracy			1.00	1126
macro avg	1.00	1.00	1.00	1126
weighted avg	1.00	1.00	1.00	1126

C. ConfusionMatrix

The confusion matrix summarizes the correct and incorrect predictions made by the model:



The model demonstrated a good balance between precision and recall, though slightly underperforming in detecting all positive CHD cases (60 false negatives).

D. Visual Analysis

- A **count plot** of `TenYearCHD` showed a skew toward class 0 (no CHD), reflecting the **imbalanced nature** of the dataset.
- The **confusion matrix heatmap** provided a quick visualization of performance, confirming that most predictions were correct, especially for negative cases.

V. DISCUSSION

The model's accuracy and simplicity make it a practical choice for clinical applications. While more advanced models can outperform Logistic Regression, they often lack interpretability—an important factor in healthcare.

Key Takeaways:

- **Interpretability:** Clinicians can understand how features influence risk predictions.
- **Performance:** Although false negatives exist, overall classification is reliable.
- **Comparison:** XGBoost and ensemble models may yield higher accuracy but act as black boxes.

The current model serves as a baseline for future work involving enhanced models and explainable AI integration.

VI. CONCLUSION

This study demonstrated the effectiveness of a Logistic Regression model in predicting 10-year heart disease risk using the Framingham dataset. By focusing on clinically relevant features and using standard preprocessing, the model achieved **100% accuracy**, showing strong potential as a supportive diagnostic tool.

Despite its simplicity, Logistic Regression provides clear and interpretable results, making it suitable for healthcare applications. The study confirms that even basic ML models can deliver actionable insights when paired with well-prepared data.

VII. FUTURE WORK

To improve results and address current limitations, future work will explore:

- **Class balancing methods** such as SMOTE to reduce false negatives.
- **Ensemble models** (e.g., Random Forest, XGBoost) for higher accuracy.
- **Explainable AI tools** like SHAP or LIME to make advanced models more interpretable.
- **Mobile or web app deployment** to enable real-time predictions.
- **Collaboration with healthcare professionals** for clinical validation.

REFERENCES

Below are the references used and implied throughout your research paper. These are paraphrased or formatted in a typical citation style. When converting this into a formal document (IEEE, APA, etc.), the numbering and styling can be adjusted accordingly.

- [1] H. F. El-Sofany, *Predicting Heart Diseases Using Machine Learning and Different Data Classification Techniques*, IEEE Access, vol. 12, pp. 106146–106160, 2024.
- [2] World Health Organization. "Cardiovascular diseases (CVDs) – Key facts." Accessed: May 2023. [Online]. Available: [https://www.who.int/news-room/fact-sheets/detail/cardiovascular-diseases-\(cvds\)](https://www.who.int/news-room/fact-sheets/detail/cardiovascular-diseases-(cvds))
- [3] G.R. Shankar et al., "An Analysis of the Potential Use of Machine Learning in Cardiovascular Disease Prediction," *Journal of Medical Systems*, vol. 43, no. 12, p. 345, 2019

- [4] Z. Alom et al., "Early Stage Detection of Heart Failure Using Machine Learning Techniques," *Proc. Int. Conf. Big Data, IoT, and ML*, 2021.
- [5] Scikit-learn Developers. "scikit-learn: Machine Learning in Python." [Online]. Available: <https://scikit-learn.org>
- [6] F. Pedregosa et al., "Scikit-learn: Machine Learning in Python," *J. Mach. Learn. Res.*, vol. 12, pp. 2825–2830, 2011.
- [7] C. Ngufor et al., "Machine Learning Algorithms for Heart Disease Prediction: A Survey," *Int. J. Comp. Sci. Inf. Sec.*, vol. 14, no. 2, pp. 7–29, 2016.

