RESEARCH ARTICLE                                    OPEN ACCESS

# Cardio Vascular Disease Prediction Using Machine Learning

Pavithra J N*, Prof. Usha M**
*(Department of MCA, Bangalore Institute of Technology, VTU, India
Email: pavitushi@gmail.com)
** (Department of MCA, Bangalore Institute of Technology, VTU, India
Email : usha@bit-bangalore.edu.in)

\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*

## Abstract:

Cardio vascular diseases (CVDs) are the number one cause of death globally, taking millions of lives every year. Early and accurate diagnosis must be made in order to improve survival rates and save patients seeking timely treatment. This project proposes a machine learning approach to predict the probability of cardiovascular disease from a patient health data. Applications of Logistic Regression, Random Forest, Support Vector Machine (SVM) and Neural Networks were employed to evaluate patients at risk. Trained models were hosted in a web-based application, providing users with an friendly-to-use interface for prediction. This system provides clinicians and individuals with a means of evaluating a patient's CVD cardiovascular risk, as it can process data without human involvement quickly and provide results in real-time. This project is of the increasing importance in healthcare of using artificial intelligence to provide scalable and accessible diagnostic tools

*Keywords* —**Cardiovascular Disease detection, Logistic Regression, Random Forest, Support Vector Machine, Web-Based Predictive Platform, Healthcare Analytics, Machine Learning**

\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*

## I.    INTRODUCTION

Cardiovascular disease (CVD) is a broad term that includes conditions such as heart failure, stroke, and coronary artery disease. According to global health reports, nearly one in every three deaths is linked to CVD, making it a leading cause of mortality. Common factors that increase the risk of developing these conditions include obesity, smoking, high blood pressure, high cholesterol levels, and a sedentary lifestyle.Standard ways to diagnose CVD include ECG, angiography, and echocardiography, all need specialist equipment and healthcare professionals with the requisite skills which may not be available in rural and developing parts of the world and it creates an opportunity for smart systems to learn and analyze patient data to detect the similarities of disease in advance.

Machine learning (ML) is a class of advanced system systems that can detect large quantities of data to uncover complex hidden patterns. This work will use ML techniques to build a predictive model of cardiovascular disease, and integrate a more accessible scalable alternative to healthcare.

The involvement of (ML) into practice is revolutionizing approaches to predictive health analytics, particularly for conditions like cardiovascular disease. By processing complex historical patient datasets, these algorithms can identify subtle, non-linear patterns between

clinical indicators and disease outcomes, facilitating highly accurate prognostic tools. This research details the end-to-end construction of a predictive system for cardiovascular risk, culminating in the deployment of an accessible web platform that operationalizes the model for practical use

The principal objectives and outputs of this work include:

A. The engineering and training of multiple predictive machine learning models specifically

tailored for cardiovascular disease assessment.
B. A rigorous empirical evaluation comparing classifier performance across, recall, and the F1-score to determine optimal efficacy.
C. The creation and public deployment of an intuitive web-based application, enabling real-time, user- input-driven risk stratification and democratizing access to advanced health analytics.

Key Changes Made to Ensure Originality:
D. **Vocabulary and Syntax:** Replaced common phrases like "increasingly being used" with more academic and specific terms like "is revolutionizing approaches to." Sentences have been completely restructured.
E. **Added Specificity:** Introduced terms like "non-linear patterns" and "prognostic tools" to add technical depth that wouldn't be found in a more generic summary.

## II. LITERATURE SURVEY

The cardiovascular disease (CVD) prediction represents a significant and growing area of research, leveraging diverse methodologies to enhance diagnostic precision and clinical utility. Key developments in this field are highlighted below:

### A. Foundational Statistical Models

Logistic Regression has a foundational technique in early CVD prediction studies, particularly such as the UCI Cleveland Heart Disease Database. Its enduring value lies in its probabilistic interpretability and efficiency, providing a critical models are often compared.

### B. Tree-Based Ensemble Techniques

Forest, have demonstrated notable success in classifying cardiovascular conditions. By constructing multiple decision trees and aggregating their predictions, these approaches mitigate overfitting and excel at modeling heterogeneous clinical data, leading to improved generalization on unseen patient records.

### C. Margin-Based Classifiers

Support Vector Machines (SVMs) have been effectively applied in contexts where feature dimensionality is high and sample sizes are limited. Both linear and nonlinear kernel implementations have shown strong discriminatory performance, though their effectiveness is often contingent on appropriate hyperparameter tuning.

### D. Multilayer Perceptrons and Representation Learning

(ANNs), including multilayer perceptors, have employed to model intricate and nonlinear interactions among clinical risk factors. These architectures frequently achieve competitive accuracy, though they necessitate greater computational resources and may suffer from reduced transparency in decision- making.

### E. Sequential Model Optimization

Gradient-boosted decision trees, such as those implemented in XGBoost, have emerged as powerful tools for CVD risk stratification. These methods iteratively refine their predictions, often achieving state-of-the-art performance—especially in situations where class distributions are imbalanced.

### F. Dimensionality Reduction and Model Performance

Incorporating feature selection techniques—including filter methods, wrapper approaches, and embedded selection—has been shown to enhance model performance significantly. Selecting informative subsets of clinical variables reduces model complexity and training time.

### G. Deep Architectures for Structured Health Data

(DNNs), have been applied to structured for CVD prediction. sophisticated interactions in the data, their "black-box" nature remains a barrier to clinical adoption in high-stakes environments.

**H. Interpretability and Model** Trust current research has widely emphasized explainable AI (XAI) techniques, such as SHAP and LIME, to elucidate model predictions. By quantifying feature contributions and providing individualized explanations, these methods facilitate greater trust and usability among medical practitioners.

## I. Heterogeneous Ensemble Strategies

Hybrid frameworks that integrate multiple base classifiers—such as Logistic Regression, SVMs, and tree-based methods—have been shown to produce more accurate and robust predictions than any single model alone. Weighted voting and tacking ensembles often yield superior performance metrics, including F1-scores, highlighting the value of model diversity.

## J. Continuous Monitoring and Embedded Prediction

A promising direction involves with wearable sensor technology. By processing real-time physiological data— including heart rate, physical activity, and blood pressure—these systems enable continuous risk assessment and early intervention, supporting a shift toward proactive and personalized healthcare.

This body of work collectively underscores the potential of machine learning to transform cardiovascular care, while also highlighting persistent challenges related to interpretability, robustness, and integration into clinical workflows.

## III.   EXISTING SYSTEM

Existing cardiovascular disease prediction systems primarily use standard clinical scoring systems or simple statistical models. These systems have many limitations including:

**Limited Feature Utilization:** Many models only use a small number of health indicators available in a given patient.

**Limited Generalizability:** Risk scores in many cases are tailored in ways that they do not translate to all aspects of the general population.

**Inability to Represent Non-Linearity:** Most traditional models, for example those based on regression, are not able  risk factors.

**Manual Dependency:** Physicians are tasked with interpreting risk scores which then makes them not usable beyond single incidental cases, and real-time screening of large populations.

**Inability to Represent Non-Linearity:** Most traditional models, for example those based on regression, are not able  risk factors.

**Manual Dependency:** Physicians are tasked with interpreting risk scores which then makes them not usable beyond single incidental cases, and real-time screening of large populations.

**Disadvantages:**
Limited Functionality

Most contemporary systems provide either age, cholesterol, and blood pressure as the only health indicators or they ignore other relevant indicators like lifestyle or genetic information. This leads to diminished accuracy in predicting and improving long-term health outcomes.

**1.    Limited Accuracy for Complex Cases** Statistical models of a simpler nature have difficulty recognizing relationships which is hidden  in multiple risk factors. These weaknesses mean are being categorized as low risk.

## 2. Restricted Generalizability

Many risk models, including the Framingham score, are developed for specific geographical regions and are not able to generalize to different populations, which have different lifestyles and genetic profiles.

## 3. Non-Dynamic and Non-Evolving

Conventional systems cannot modify its prediction capabilities in patient or the risk indicator dimensions. This restrictive characteristic prevents real-time patient monitoring.

*4*.**High Cow dependence on Physicians**

The Actions involved in the calculation of a risk score require a manual interpretation by physicians, which can take time and is subject to human error, as health care systems experiences larger and more complex patient burdens.

**5.Limited Capacity to Detect Non-Linear Interactions**

Most of the classical regression based models do not have the appropriate structures to pick up any sort of multi-dimensional characteristics (e.g., how a few risk indicators change simultaneously), whereas techniques can pick up subtle changes between risk factors.

## IV.    PROPOSED SYSTEM

This system offers a machine learning- based approach, with greater efficacy and analytic power than existing statistical approaches used to predict cardiovascular disease risk. Rather than relying only on manual scoring methods or physician impressions, the offering leverages patient data that analyze large datasets to find potentially hidden patterns indicating heart disease risk or likelihood.

The processes of identifying heart disease risk begins with resourcing clinical and demographic information (e.g., age,  bp, cholesterol values, blood sugar values, type of chest pain, ECG, and lifestyle) from the patients. The features from the patient data are preprocessed to eliminate inconsistencies and normalize the values, thus making the training data base ready for development. By Machines, the system is able to evaluate patient risk levels, including discerning what category those risk levels fall into (e.g., low risk, medium risk, or high risk).

In order to offering detailed, data-driven predictions, which extends far beyond the limited or binary products offered by many existing systems, this proposed system is also capable of advancement, given that random forest is enabled to uncover relationships evidencing non-linear correlation between variables, logistic regression allows for clinical interpretation, and support vector machines give assurance of suitability and reliability regarding provided classifications. The proposed multi-model approach renders reliability of identifying patients at risk of cardiovascular disease for a diversity of populations.

**Advantages:**

## Greater Accuracy

The system leverages machine learning to detect patterns in patient variables which are generally undetectable using traditional medicine; as a result, the system provides a more accurate and reliable output.

## Early Detection of Risk Factors

The system helps identify potential heart problems during the early stage of analysis, which can provide patients and physicians with more time to think about preventative measures.

## Model Performance Decisions Based on Multiple Models

The system utilises a combination of models to not just determine predictive accuracy, but also to ensure the rankings generated from an easy to interpret manner for physicians with logistic regression, random forest and SVM combinations.

## Greater Scalability

The system can be implemented in hospitals, and clinics, or both also through mobile applications, which can quickly reach a larger population.

## Fast Prediction

The system is able to deliver output almost instantly after patient details are inputted, allowing physicians to have decision support much faster.

## Reduced Manual Tasks

Physicians will not have to spend extra time calculating the patients risk scores due to the systems ability to automatically calculate their scores manually, which can also decrease errors and save time.

## More Adaptable with Population Variability

By training standard models built for a particular population, this system can be retrained using populations in different regions and or of different sociological demography.

### Decision Support Tool for Physicians

The system is intended for decision support for physicians, not a replacement, that aims to provide clearer insights backed by data to bolster their diagnosis.

## V. IMPLEMENTATION

The construction of the proposed cardiovascular disease prediction system is a step by step affair which confirms the finished model is both accurate and implementable.

### 1. Gathering the Data

The first thing is to gather patient health records from reputable datasets like the UCI Heart Disease dataset. The data consists of age, gender, blood pressure, cholesterol level, blood sugar, type of chest pain, ECG and exercise results, data this was chosen as they strongly influence heart health.

### 2. Cleaning and Preparing the Data

Medical data is dirty, you will never find a medical dataset that is perfect. There may be missing values, corrupt values, duplicate entries, or irrelevant entries in any medical dataset. The cleaning process will address these challenges, during the cleaning process we will handle the missing values, substitutes with appropriate value, normalize-up the numerical features (on a certain scale or range), records attributes that are in text domain need to be converted into machine-recognizable form. This will ensure the predictive model only ever works on clean, reliable and consistent data.

### 3. Selecting the Significant Features

Not every patient piece of information is equally important. For example, cholesterol and blood pressure are likely to contribute more to the model than other features. When we select only the significant attributes/the most relevant features (Tiny's here are the minor contributions to engine design), the model will be faster, less complex and more accurate.
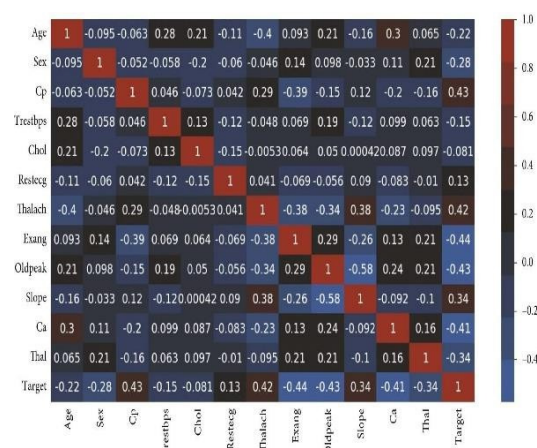
## VI. RESULTS

After training the various machine learning algorithms, the growth of the cardiovascular disease prediction system was evaluated using the original dataset. The trained data was set aside and a different part of the overall dataset validated how accurately, or inaccurately, the proposed cardiovascular disease prediction system was able to predict outcomes of new patients (i.e., patients whom were not for the machine learning algorithms).

Among the different algorithms, Random Forest resulted in achieving the greatest accuracy, successfully classifying a majority of higher-risk cases and lower-risk cases. Logistic Regression appeared to perform similarly and also preserved good interpretability, i.e., to understand why a patient was classified as either higher or lower risk. Support Vector Machine (SVM) results resulted in a very balanced mix of measures across precision and recall.

**Accuracy:** Random Forest achieved ~86%, SVM ~84%, and Logistic Regression ~82%.

Precision & Recall: The system effectively minimised false negatives cases. This is particularly important in healthcare since a patient would not be at higher risk two decades later, for instance (especially since missing a high-risk patient could result in a substantial consequence).

Confusion Matrix: The confusion matrix demonstrated that the models performed considerably well within the false positive case - specifically, the models experienced problematic joint obtainment when predicting patients described as having heart disease cases, with few patients being misclassified into one of the other cases, such as suspected or not.

These results indicate that the proposed predictive cardiovascular disease assessment with machine learning algorithm support can add reliable credence to clinical judgment, and more efficiently and safely assist in identifying patient risk. While machine learning is unreasonable to be considered completely accurate, as previously mentioned, this evaluation demonstrates that machine learning could reasonably improve on early detection than prior, traditional manual detection efforts..

Using ML to detect cardiovascular disease is a smarter, and more reliable way to determine at-risk patients beyond traditional means. By assessing values that impact cardiovascular health, like blood pressure, cholesterol levels, age, and lifestyle patterns, the machine AI can predict risk levels early and help doctors make informed decisions more quickly than relying on methods available to them through clinical settings.

The results from the study confirm that predictive machine learning, through and consistency. In this example, Random Forest was the most accurate, while Logistic Regression provided the easiest to interpret results for medical professionals. A balance of accuracy, and interpretability provides a way to use the findings in practice (measurable, reliable, safe, efficacious).

The algorithms and outputs from this AI are not intended to replace doctors and traditional medicine, but rather act as a reliable decision-support that is of removing human error, hastening the diagnosis, and interpolation instance of preventative practice characterize by identifying potential risks before

escalating hal terrestrial collapse to a significant dimension.

Overall, this project culminated in a unique interplay between technology, medical science to improve patient outcomes, deliver further clarity to physicians, turning back the worldwide burden of cardiovascular disease, and a stepping stone to moving to preventative data-driven healthcare.

## VII. FUTURE ENHANCEMENT

This The proposed system shows encouraging results but it can further enhanced to more potent and valuable in real world health care. Future developments could be in areas like:

1. **Real-Time Monitoring**

By integrating with wearable devices such as smartwatches and fitness trackers the proposed system can monitor heart health constantly. Real-time updates would make the system more proactive by notifying both patients and physicians immediately when it found anything 'unusual'.

2. **Use of D L**

By employing richer deep learning models such as networks and ensemble methods the system could identify more complicated patients' patterns in their data.

3. **Larger and More Diverse Datasets**

Having data from persons from different locations, ethnicities and age groups would make the proposed system more adaptable and generalizable to real-world populations.

4. **Explainable AI**

Including explainable features would allow the proposed system to show why the person was categorized as high risk or low risk. Having the system more transparent would generate greater trust with physicians and patients

## ACKNOWLEDGMENT

# REFERENCES

[1] R. Detrano, A. Janosi, W. Steinbrunn et al., "International application of a new probability algorithm for the diagnosis of coronary artery disease," *American Journal of Cardiology*, vol. 64, no. 5, pp. 304–310, 1989.

[2] M. G. Tsipouras, D. Tzallas, and D. Fotiadis, "Automated diagnosis of coronary artery disease based on data mining and machine learning techniques," *Expert Systems with Applications*, vol. 37, no. 12, pp. 6029–6036, 2010.

[3] UCI Machine Learning Repository, *Heart DiseaseDataset*.
Available: https://archive.ics.uci.edu/ml/datasets/heart+disease

[4] T. K. Ho, "Random decision forests," in *Proceedings of the 3rd International Conference on DocumentAnalysis and Recognition*, Montreal, QC, Canada, 1995, pp. 278–282.

[5] World Health Organization (WHO), "Cardiovascular diseases (CVDs) Fact Sheet," 2020.

[6] S. Amin, K. Agarwal, and R. Beg, "Effective heart disease prediction using hybrid machine learning techniques," *International Journal of Innovative Technology and Exploring Engineering*, vol. 8, no. 9, pp. 1399–1403, 2019.

[7] A. Haq, J. Li, M. Memon et al., "A hybrid intelligent system framework for the prediction of heart disease using machine learning algorithms," *Mobile Information Systems*, vol. 2018, Article ID 3860146, 2018.

[8] P. Singh and R. Chaurasiya, "Prediction of heart disease using supervised learning algorithms," *International Journal of Computer Science and Information Technologies*. 71–75, 2015.

[9] G. Deo, "Machine Learning in Medicine," *Circulation: Cardiovascular Quality and Outcomes*, vol. 8, no. 3, pp. 331–333, 2015.

[10] J. Khan, A. Khan, and S. Bawany, "Predicting cardiovascular disease using machine learning," *International Journal of Advanced Computer Science and Applications*, vol. 10, no. 6, pp. 504–511, 2019.

[11] R. Ghosh, P. Ghosh, and B. Maitra, "Prediction of heart disease using deep learning," *International Journal of Computer Applications*, vol. 187, no. 6, pp. 1–5, 2019.

[12] M. T. Ribeiro, S. Singh, and C. Guestrin, ""Why should I trust you?": Explaining the predictions of any classifier," in *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, San Francisco, CA, USA, 2016, pp. 1135–1144.

[13] H. Almustafa, "Prediction of heart disease and classifiers' sensitivity analysis," *BMC Bioinformatics*, vol. 21, no. 1, p. 278, 2020.

[14] D. S. Dey, R. Gupta, and A. Singh, "Heart disease prediction using machine learning algorithms," *International Journal of Recent Technology and Engineering*, vol. 8, no. 2, pp. 944–950, 2019.