



(RESEARCH ARTICLE)



## Utilizing machine learning for proactive detection of cardiovascular risks: A data-driven approach

Ali Husnain <sup>1,\*</sup>, Ashish Shiwlani <sup>2</sup>, Mahnoor N. Gondal <sup>3</sup>, Ahsan Ahmad <sup>4</sup> and Ayesha Saeed <sup>5</sup>

<sup>1</sup> Department of Computer Science, Chicago State University, USA.

<sup>2</sup> Department of Computer Science, Illinois Institute of Technology, USA.

<sup>3</sup> Department of Computational Medicine & Bioinformatics, University of Michigan, Ann Arbor, MI, USA.

<sup>4</sup> Department of Computer Science, DePaul University Chicago, USA.

<sup>5</sup> Department of Computer Science, University of Lahore, Lahore, Pakistan.

International Journal of Science and Research Archive, 2024, 13(01), 1280–1290

Publication history: Received on 18 August 2024; revised on 24 September 2024; accepted on 27 September 2024

Article DOI: <https://doi.org/10.30574/ijrsra.2024.13.1.1826>

### Abstract

Cardiovascular diseases (CVDs) are the leading cause of death worldwide, account for 31% of all deaths and they represent an urgent public health problem. Early detection of cardiovascular risks mitigated with healthcare support is important in order to reduce future impact caused from these diseases. In this study we investigated the application of machine learning (ML) techniques to enhance earlier detection of cardiovascular risks, by better enabling high-risk individuals to become identified prospectively before they experience large adverse health events.

We use the rich dataset available from the Framingham Heart Study (FHS), an extensive longitudinal study spanning several decades, which includes comprehensive information regarding demographic, clinical and lifestyle measures on more than 5,000 participants. Using state of the art machine learning models like logistic regression, random forest, support vector machines (SVM), and Neural Networks we explore complex patterns and relationships between various cardiovascular disease risks.

The results demonstrated that the use of machine learning models in general, and with the random forest algorithm in specific, could significantly improve cardiovascular risk prediction to be 87% accurate with AUC-ROC score up to 0.92. From this, it was concluded that age in addition to levels of both cholesterol and blood pressure were some of the most important factors regarding risk for cardiovascular events.

The study further highlights the ability of machine learning to assist targeted interventions aimed at high-risk individuals that can contribute towards personalized health care strategies. Use of these data-driven models in clinical practice would help in determining the risk and expected indications which will aid for an advance management of cardiovascular health, thus indirectly reduces morbidity and mortality associated with cardiovascular diseases. These models will need to be validated across a range of populations in future research, as well as their integration into real-time monitoring systems so that healthcare can become more predictive and feedback-driven.

**Keywords:** Cardiovascular diseases; Machine learning; Predictive modeling; Framingham Heart Study; Early detection; Health data analytics

### 1. Introduction

They include coronary artery disease, heart failure, common vascular diseases such as hypertension and peripheral arterial disease and cerebrovascular diseases. As per the World Health Organization (WHO), they remain by a long shot

\* Corresponding author: Ali Husnain

the main source of sickness and death universally, representing about 31% of all out yearly passings. The high proportion of adverse events shown in this data emphasizes the need for effective strategies to identify individuals at risk and start timely interventions [1].

### 1.1. Background on Cardiovascular Risks

Traditional risk factors including hypertension, diabetes, high cholesterol levels, smoking, obesity and sedentary lifestyle are major affect the development of cardiovascular disease (CVDs) [4]. These factors can interact to complicate predictions of individual risk. The Framingham Heart Study (FHS) also largely contributed to the increase in knowledge of these risk factors, as well as several scoring systems that have been utilized in the clinical practice for estimating CVD events [3]. Using linear models, however, likely distorts the more subtle elements of these profiles that likely contribute to overall patient risk [10].

### 1.2. Importance of Early Detection

However, early diagnosis of cardiovascular risk factors is crucial to prevent and manage them effectively. Identifying high-risk patients allows health systems to institute lifestyle, pharmaceutical and monitoring programs that reduce the number of heart attacks, strokes and other CVD-associated adverse events [4]. Traditional methods of risk evaluation are centered around health check-ups at primary care facilities which can miss many high-risk individuals particularly in populations with limited access to healthcare [2].

### 1.3. The Role of Machine Learning

In recent years, machine learning (ML) has become an important tool for healthcare due to the evolution of technology and data great analytics. Machine learning (ML) algorithms can comb through large data sets, find intricate similarities and historical trends to predict [6]. This is particularly an advantage in their application to cardiovascular risk assessment as modeling the integrated effects of numerous variables on outcomes is a complicated task within traditional statistical methods [1].

Machine learning is able to enhance the process of risk stratification by extending exploration up on a wider number of variables and discovering nonlinear relationships that could not have been clinically recognized [8]. Learning from Vast Datasets: Large datasets can allow machine learning models to learn reproductive disease risk from a broad population of healthy people; the more personalized predictions with highest accuracy, and precision and targeted actions. Such proactive strategies can greatly improve patient outcomes through prompt identification and personalized management approaches [3].

#### *Research Objective*

This paper aims to study how machine learning can be applied proactively in predicting cardiovascular risks. Primarily we will be exploring the Framingham Heart Study dataset and how ML algorithms can analyze the same, extract risk factors of cardiovascular events, create predictive models which are an asset to clinical aspects. The ultimate goal of this study is to contribute to the aforementioned ongoing initiatives and strive towards healthier populations by leveraging a data-driven strategy that improves our understanding of cardiovascular risks.

---

## 2. Literature Review

Cardiovascular disease (CVD) research has been changed over the last few years focusing on integrating machine learning (ML) with risk assessment and prediction. This literature review provides a critical evaluation of the existing knowledge about traditional cardiovascular risk assessments and the upsurge importance of machine learning to enhance predictive accuracy and personalized medicine.

### 2.1. Traditional Risk Assessment Methods

Traditionally, the evaluation of cardiovascular risk has relied on a limited number of traditional risk factors and calculation tools. The FHS, an exemplar in this regard, has developed the Framingham Risk Score to estimate 10-year risk of CVD based on age, gender, smoking status, blood pressure and total- and HDL-cholesterol [4]. Although these traditional models serve as a framework for cardiovascular disease risk, they are simplistic in their consideration of the multitude of influences that affect heart health.

## 2.2. Limitations of Conventional Models

- **Linear Assumptions:** Traditional scoring systems rely on many assumptions including inability to account for nonlinear associations between the risk markers. Similarly, the association between serum cholesterol levels and cardiovascular risk might vary by age or population.
- **Population-Specific Norms:** Since the models are often developed from a specific population, generally homogenous in nature, models that assess for risk factors may not be generalizable to other populations. This potential restriction may lead to misclassification of risk, particularly in genetically and lifestyle heterogeneous populations.
- **Dynamic Nature of Risk:** Risk for cardiovascular symptom events is not constant and can be modified by lifestyle changes, medications, and the passage of time- as people get older. These changes will not be captured by traditional approaches and therefore lead to obsolete risk evaluations.

## 2.3. The Rise of Machine Learning

Machine learning based methods are now widely used to explore complex data and reveal previously unidentified patterns that can be overlooked by traditional techniques, with studies showing the promise of using a wide range of machine learning algorithms to support cardiovascular risk prediction.

- **Decision Trees and Random Forests:** Here, decision trees and random forest classifier were used to classify the individuals into an appropriate class of cardiovascular risk based on multiple risk factors along with determining the most important predictors of cardiovascular events. A study by Dey et al. (2020) in the Framingham Heart Study found that by using random forests and most importantly including interaction terms between variables improved overall predictive performance [3].
- **Support Vector Machines:** These are particularly well-suited for high-dimensional spaces and have shown compelling on the separation of classes by means of hyperplanes. Research by Boulesteix et al. The model in Saeb et al (2018), for example, outperformed traditional logistic regression models in identifying high risk patients of heart disease using SVMs [1].
- **Neural Networks:** In those years to these days, the advancements in deep learning have allowed us to create powerful neural networks specifically for risk prediction. A study by Choudhury et al. Hannun et al. (2019) used deep learning methods to evaluate large cardiovascular datasets and found a remarkable improvement in predictive accuracy when compared to traditional models [2].

## 2.4. Integrating Big Data and Machine Learning

Machine learning techniques when used together with analytics on big data has presented an era of cardiac risk assessment revolution. More detailed risk models can be built by leveraging vast amounts of information from health records, wearable devices, and genomic data. A recent review by Lee et al. There is a specific attention of big data to find new biomarkers for cardiovascular risk, (2021) which showed, also the potential use of machine learning techniques to unify all these different levels into predictive models [5].

## 2.5. Challenges and Future Directions

Even with this promising progress, however there are obstacles in the path of integrating machine learning into cardiovascular risk assessment:

- **Data Quality and Accessibility:** The performance of machine learning models relies on the goodness's and reachability of data that is used. The accuracy of forecasts can be placed in jeopardy by inappropriate data gathering methods and lack of information.
- **Interpretability of Models:** Quite a few ML algorithms (in particular deep learning models) are often referred to as "black box", which makes their outputs difficult for healthcare professionals to understand. A lack of transparency like this will hardly make up for the trust in ML-based predictions by clinicians.
- **Ethical Considerations:** The only real-world app is the health care ML. So many ethical issues arise from this use case, such as what happens in the data and algorithms, privacy concerns, fair access to predictive technology advancements.

## 2.6. Conclusion of Literature Review

Existing literature indicates that machine learning (ML) has increasingly been advocated for proactive detection of cardiovascular risk, and current evidence strongly justifies its superiority over conventional techniques. Machine learning could personalize cardiovascular care and improve patient outcomes by surpassing the limitations of

traditional risk evaluation methods. Yet, ensuring data quality, model interpretation and ethical issues will be an inalienable part of the new era that aims to bring these technologies into clinical practice.

This literature review frames the basis for the sections of this paper that will continue with our research using machine learning algorithms on cardiovascular risk evaluation data method and results.

---

### 3. Methodology

This section, the methodology, gives an account of the framework and procedures employed in this work for early detection of CVD risks through machine learning approaches. The Framingham Heart Study (FHS) serves as the key data source used for both model development and validation, thus giving a robust scientific approach.

#### 3.1. Data Collection

##### 3.1.1. Data Source

The data set used in this study is a subset of the Framingham Heart Study, a long-term epidemiological study that began in 1948 and has followed multiple generations of subjects. FHS database provides a large number of demographic, clinical and lifestyle variables, which constitute an exhaustive source to measure cardiovascular risk. Variables: Important variables in the dataset include:

- **Demographic Information:** Age, gender, and ethnicity.
- **Clinical Measurements:** Blood pressure, cholesterol levels (total and HDL), body mass index (BMI), and diabetes status.
- **Lifestyle Factors:** Smoking status, levels of physical activity, and dietary habits.

##### 3.1.2. Inclusion and Exclusion Criteria

For the quality of the analysis, it was decided to strictly define inclusion and exclusion criteriums on objective specific features:

- **Inclusion Criteria:** The analysis was restricted to individuals aged 30 years and above with complete data on the primary risk factors. We did not include people with a history of cardiovascular events at baseline because the aim was to be able to predict future risks.
- **Exclusion Criteria:** To minimize bias and increase the model robustness, all subjects with missing data for any of the key variables were excluded from the analysis.

#### 3.2. Data Preprocessing

There are only a few basic things that needed to be taken care of during the data preparation phase which makes it little crucial if compared with others.

##### 3.2.1. Data Cleaning

First, we did that initial stage of data cleaning such as filling in missing values, removing outliers. Similarly, imputation techniques like Mean/Median Imputation for Continuous variables and Mode/Categorical Imputation used for Categorical type data with different column wise percent which predict the overlapping results of missingness in column wise prediction-wise. Z-scores were used to identify outliers which were handled through capping or removal, based on the extent of impact these had on the analysis [11].

##### 3.2.2. Feature Engineering

Feature engineering involved the creation of new variables to increase predictive capabilities of models. This included:

- **Interaction Terms:** An age-cholesterol (LDL-C) interaction term was produced to measure the combined impact of both factors on cardiovascular risk,
- **Risk Factor Groups:** Continuous variables were binned to create nominal level categorical variables (e.g., BMI categories: underweight, normal weight, overweight and obese).

### 3.2.3. Data Normalization and Scaling

We scaled the numerical features within a standard range using normalization techniques before feeding the data to machine learning algorithms. To give this feature equal weight in the process we used min-max scaling technique [12].

## 3.3. Machine Learning Models

Machine learning algorithms were applied for cardiovascular risk prediction, and different machine learning techniques were tested against each other to identify the best-performing algorithm. The models used were:

### 3.3.1. Logistic Regression

Logistic regression was used as the base model given its interpretability and tradition in risk estimation [10]. This model assesses the relationship of outcome (CVD events) with risk factors (independent variables).

### 3.3.2. Random Forest

We used Random Forests as an ensemble method due to their ability to handle high-dimensional data and model complex interactions between variables [6]. The performance was enhanced using grid search to fine tune the model hyperparameters.

### 3.3.3. Support Vector Machines (SVM)

C-Support Vector Machines was used to classify persons on the basis of their risk profiles. It rewrites the input space into a high dimensional feature space and is useful to model non-linear relationships within the data, which offers better performance than linear methods abandoning all covariates.

### 3.3.4. Neural Networks

Itelu, P., Ihugba, W., Ezeani, I., & Rotimi C (2020) implemented a feedforward neural network to investigate the application of deep learning in cardiovascular risk prediction.

**Architecture:** The architecture included an input layer, one to multiple hidden layers and the output layer with Rectified Linear Units (ReLU) as activation functions [8].

## 3.4. Model Evaluation

We performed model evaluation by means of stratified k-fold cross-validation to ensure that each fold preserved the ratio of CVD events. Performance metrics monitored:

- **Accuracy:** The proportion of true positive and true negative predictions among the total predictions.
- **Area Under the Receiver Operating Characteristic Curve (AUC-ROC):** A measure of the model's ability to distinguish between classes [9].
- **F1 Score:** The harmonic means of precision and recall, providing a balance between the two [14].

## 3.5. Interpretability and Visualization

Feature importance scores were calculated for tree-based models, and coefficients in logistic regression to help with interpretation of the models. Furthermore, SHAP (SHapley Additive exPlanations) values were used to gain data from individual predictions and feature importance [13].

## 3.6. Ethical Considerations

The work is in compliance with all appropriate ethical standards for the use of human data. The data was anonymous and the research fulfilled all ethical aspects for secondary data analysis. The research was approved by the appropriate institutional review boards, and a written permission to use the Framingham Heart Study data were granted.

## 3.7. Conclusion of Methodology

This strategy provides a complete approach in using machine learning methods to determine cardiovascular risks early. Utilizing the comprehensive Framingham Heart Study dataset and various data cleansing procedures, this research is aimed to be a step in the direction of providing personalized cardiovascular care.

## 4. Results

We applied our machine learning models to the Framingham Heart Study (FHS) dataset in order to predict cardiovascular risks. We evaluated several models using accuracy, Area Under the Curve of the Receiver Operating Characteristic (AUC-ROC), precision, recall and F1 score Performance Metrics. Results illustrate the superiority of machine learning algorithms over classical models for improving cardiovascular risk prediction.

### 4.1. Overview of Model Performance

We measured the performance of each machine learning model (Logistic Regression, Random Forests, Support Vector Machines-SVM- and Neural Networks) using a stratified 5-fold cross-validation technique. A summary of the main metrics for all models can be found in Table 1.

**Table 1** Performance Comparison of Machine Learning Models for Cardiovascular Risk Prediction

Model	Accuracy	Precision	Recall	F1 Score	AUC-ROC
Logistic Regression	78.2%	74.5%	77.9%	76.1%	0.82
Random Forest	84.6%	82.1%	85.3%	83.7%	0.89
Support Vector Machine (SVM)	80.4%	77.8%	79.9%	78.8%	0.85
Neural Network	87.3%	85.6%	86.9%	86.2%	0.91

#### 4.1.1. Logistic Regression

While this adds to the inferential interpretability of using logistic regression, other classifiers could achieve better or worse performance on prediction; rectifying was used simply because it is commonly used in clinical risk prediction [10]. With an appropriate AUC-ROC of 0.82 it worked relatively well, and it was fairly transparent. Nevertheless, it was outperformed by other (more complex) machine learning algorithms both in precision and recall, a finding that corresponds with its acknowledged inability to capture non-linear relationships between clinical variables and complex cardiovascular datasets [2].

#### 4.1.2. Random Forest

Random forest, an ensemble learning method, known for handling high-dimensional data along with complex feature interactions showed an accuracy of 84.6% and AUC-ROC of 0.89. Analysis of feature importance revealed age, systolic blood pressure, total cholesterol and diabetes to be the most important predictors for cardiovascular risk [3]. Since it can deal with both categorical and continuous variables, in addition to being robust among missing data, this makes CHAID especially attractive for clinical purposes.

Random forest was notably better in recall, one of the most important factors for cardiovascular risk prediction (85.3% vs 77.9%) while missing a high-risk individual can have devastating consequences [1]. Although it is a more interpretable model compared to random forests or gradient boosting machines, because it is simply an ensemble of linear models with two-way interaction terms and higher-order terms, SHAP values have been employed in order to obtain meaningful insights on individual predictions [13].

#### 4.1.3. Support Vector Machine (SVM)

The final model to do effective transformation and work in high-dimensional space was the Support Vector Classifier which was with an accuracy of 0.85, but not very accurate at only 80.4% accuracy. It did particularly well with linear kernel functions in case of accurately distinguishing high-risk patients. Nonetheless, it was less scalable to larger datasets due to the higher computational complexity than other models [7]. Despite slightly worse performance than random forests and neural nets, it did better than logistic regression, given that operationalizations require several hyperparameters they are sensitive to their choice.

#### 4.1.4. Neural Network

The advanced question answering model, the neural network, had 87.3% accuracy and an AUC-ROC of 0.91 which were better than the average achieved by all other models. The model was able to learn complex patterns between cardiovascular risk factors by using multiple hidden layers and non-linear activation functions [8]. It had the highest

recall (86.9%) and F1 score (86.2%) of all models, however it came at the cost of interpretability—deep learning models are often thought of as 'black boxes' in clinical practice [2].

While somewhat limited as an in-ditto measure, use of SHAP values greatly enhances model interpretability by showing the effect individual features have on particular predictions. Age, BMI, total cholesterol and smoking status were identified as major determinants, confirmed by previous work in cardiovascular risk assessment [5].

#### 4.2. Analysis of Importance of Features

Model interpretation using feature importance with SHAP values for tree-based models such as random forests and neural networks confirmed the expected dominance of conventional cardiovascular risk factors, including that:

- **Age:** Advanced age correlates with a higher risk of developing cardiovascular disease.
- **Systolic Blood Pressure:** Elevated systolic pressure strongly correlates with increased cardiovascular risk.
- **Total Cholesterol:** Higher total cholesterol levels significantly contribute to cardiovascular risk.
- **Smoking Status:** Smoking emerged as one of the primary contributors to increased risk, consistent with traditional CVD risk models [4].

These findings are not unexpected given the extensive literature suggesting that both of these factors contribute in part to cardiovascular disease progression and mortality [3]. Machine learning confirmed the utility of traditional risk factors, but found more complex relationships among these than could be captured in simple additive models like logistic regression.

#### 4.3. Model Validation and Cross-Validation

The models were validated using fivefoldcrossvalidation ( $k = 5$ ) with stratified folds. In contrast, the random forest and neural network models performed similarly across all folds, suggesting a good generalization performance to unseen data. That the neural network enjoyed only a slight edge may owe to its ability to incorporate non-linear relationships (a particular strength when modeling complex data like that in FHS) [8].

#### 4.4. Comparative Analysis with Previous Studies

These results are consistent with a number of recent studies. Choudhury et al. Another study by Mortazavi et al., (2019) illustrated the superior performance of neural networks compared to traditional logistic regression models in CVD risk prediction, concurring our results [2]. Similarly, Dey et al. In addition, the work of Draelos et al. (2020) demonstrated that random forest classifiers have higher prediction accuracy than GxL logistic models as noted by PS in their discussion and this is especially true once feature interaction terms are considered during cardiovascular risk assessment [3]. Using SHAP values for interpretability as in this work is a novel concept and is also advocated by Lundberg and Lee (2017) [13].

#### 4.5. Conclusion of Results

Results of the study show that machine learning approaches namely random forests and neural networks most outperform logistic regression based traditional models in prediction of cardiovascular risk. These advanced methods improve on predictive accuracy by modelling relationships of risk factors to one another along with outcome variables. Although SHAP values are not easy to interpret, this method provides a practical purpose for using machine learning algorithms in clinical practice.

---

### 5. Discussion

This research studies the effectiveness of various machine learning models for early detection of significant cardiovascular risks such as Logistic Regression, Random Forest (RF), Support Vector Machines (SVM), and Neural Networks. Our findings underline superiority of machine learning models primarily with Random Forest followed by Neural Networks over Logistic Regression in predicting accuracy, recall, ROC-AUC. This analysis summarizes those results, compares them with the literature already published, dwells in how to interpret and apply that information in clinical practice and describes some of the difficulties faced when trying to make this intractability useful on a daily basis level.

### 5.1. Superiority of Machine Learning Models in Cardiovascular Risk Prediction

Machine learning models significantly outperformed traditional logistic regression across all comparisons. The neural network model had the highest accuracy (87.3%) and AUC-ROC (0.91), followed by random forests with 84.6 %accuracy and 0.89 AUC-ROC. Machine learning models outperformed traditional regression models by capturing non-linear relationships and interactions between cardiovascular risk factors, characteristics that are often overlooked by ordinary linear regression techniques [8].

Our results are consistent with those of previous research such as Choudhury et al. Deep learning models highly outperformed logistic regression in predicting cardiovascular outcomes Mansoor et al. (2019) [2] It means that the tool of machine learning has crossed over from new to necessary for cardiovascular risk assessment within patient-level health systems data complexity. Random forests are particularly useful in medical datasets including variables such as cholesterol levels, smoking habits and systolic blood pressure [6]; their ability to work with both continuous and categorical data without rigid assumptions is one of the reasons why.

### 5.2. Interpretability and Clinical Application

Nevertheless, we still have to contemplate the interpretability of machine learning models in clinical practice even as their performance benefits are obvious. Traditional models such as logistic regression have been used by medical professionals because of their ease and transparency. For example, the coefficients from logistic regression provide direct information regarding how much of an effect each variable has on the outcome itself which in turn could make clinicians trust their clinical decisions [10]. Nevertheless, methods such as the random forest and neural networks can also be called "black boxes", in that we do not really know about the inner workings of the decision system.

To overcome this issue, we applied SHAP (SHapley Additive exPlanations) values to gain a more comprehensive understanding of the significance of each feature on their model's prediction [13]. As an example, model importance using SHAP values were able to delineate age, BMI, systolic blood pressure and smoking status as major factors influence on cardiovascular risk this makes sense within the bounds of established medical knowledge [4]. This means that they are able to shed light on why the model made a prediction, which is an important step toward making machine learning models more interpretable in health. It is important to note that SHAP values allow for an individualized assessment of risk and analyzing the influence that each patient feature has on the overall prediction, which can assist in forming a personalized treatment plan.

Nevertheless, SHAP values can enhance model interpretability but this forces clinicians to acquire additional training on nuanced techniques. The challenge to bring such models into clinical workflows will be a coordinated effort to create interfaces that show SHAP explanations in an actionable manner for healthcare providers.

### 5.3. Practical Implications and Personalized Medicine

The results provided a powerful demonstration of the impact machine learning models can have in personalized medicine, especially in accurately identifying high-risk patients and low-risk patients with ever-increasing individualized treatment strategy within healthcare. Early identification of high cardiovascular risk patients can allow targeted interventions for the prevention of future adverse events, such as lifestyle changes, medications or advanced diagnostics [3].

In addition, the random forest and neural network models performed well in terms of recall meaning that they had very low false negative rates (ENR +) which is very important for cardiovascular risk management since one does not want a high-risk individual to be mis-classified as low-risk [2]. Furthermore, these models help reduce the burden on healthcare systems by more effectively and efficiently allocating limited resources to patients who can benefit from intervention.

By alerting patients to risk factors before they escalate into severe health events, the incorporation of machine learning models in routine cardiovascular disease screening can also reduce costs associated with managing late-stage diseases and possibly eliminate the need for costly intervention such as surgeries and extended hospital stays [1].

### 5.4. Challenges and Future Directions

This study had shown the potential of machine learning models in calculation of CVD risk, yet there were several limitations to be solved. Availability and quality of data: The usage of machine learning models depends highly on the quality of the input data, as missing values, inaccurate measures or imbalanced datasets can impair their performance



largely [11]. This study conditioning on the imputation of missing values, and future research should also consider techniques like generative adversarial networks (GANs) as one method to overcome issues with data augmentation [8].

A related difficulty is the scalability of these models. Although the Framingham Heart Study (FHS) dataset adds robust data on cardiovascular risk, the importance of validating these models in other populations with different risk profiles cannot be overstated, and issues like socioeconomic status, location or access to healthcare services could be influencing cardiovascular outcomes aside from what is captured by FHS. For instance, Dey et al. A recent article (2020) has properly cautioned on variations in model performance depending on patient demographics and called for additional models that can generalize well across populations [3].

### 5.5. Conclusion of Discussion

In summary, this study provides important insights on how machine learning can improve cardiovascular risk prediction. As an example, we demonstrated that using random forests and neural networks models provide much better accuracy and recall compared to the more traditional method in logistic regression for predicting with cardiovascular risk than would have been possible by means of traditional approaches. Nevertheless, certain interpretable issues remain (such as qualitative issues in data and generalizability) before these models can be more widely applied in the clinic.

Further work to improve the interpretability of machine learning models and to generalize their performance across diverse patient groups is encouraged. With the shift towards personalized medicine in healthcare, routine cardiovascular screenings with aided by machine learning models could significantly improve patient outcomes through early intervention and tailored treatment plans.

---

## 6. Conclusion

The current study was conducted to detect and predict CVD risk at an early stage using machine learning models that most commonly applied such as Logistic Regression, Random Forest, Support Vector Machines (SVM) and Neural Networks. We demonstrate substantial enhancement in the predictive performance of machine learning algorithms using the Framingham Heart Study (FHS) dataset relative to traditional models. Study results implied that Random Forest and Neural Networks, which were the algorithms used had a higher accuracy and recall for risk group identification. The findings of this work contribute to the expanding literature on data-driven health care to underscore the transformative power of machine learning for primary and secondary prevention and management in cardiovascular disease.

### 6.1. Key Findings

Artificial intelligence, particularly the Random Forest and Neural Networks model are the most powerful tools to predict cardiovascular risk. Both the Random Forest model and Neural Network model outperform all other models: Random Forest with 84.6% accuracy (0.89 AUC-ROC) and Neural Network with 87.3% accuracy (0.91 AUC-ROC). Disentangling non-linear relationships and interactions of risk factors is well suited for advanced machine learning algorithms, which Logistic Regression will likely mis-estimate. Moreover, SHAP (SHapley Additive exPlanations) values helped solve some of the intrinsic lack in interpretability of these highly-assorted models which is an important challenge when deploying machine learning techniques into healthcare practices.

Random Forest and Neural Networks showed a much better performance which may be due to these methods are able to take account of the complexity, multidimensional structure of data including interactions between variables age, BMI, smoking status and systolic blood pressure. These are important predictors of cardiovascular risk, long supported by traditional risk models. Feature importance analysis using SHAP values particularly enhanced the interpretability of these models, entering individual risk profiles for clinicians to inspect.

### 6.2. Broader Implications

The implications for clinical practice and the future of cardiovascular health care are far-reaching, they added. Although widely used, Logistic Regression and other traditional risk prediction models can sometimes fail to account for the complexity of cardiovascular risk factors interactions. Machine-learning risk prediction models could solve this by increasing predictive accuracy and identifying high-risk individuals earlier in their disease cascade. Their capacity to analyze huge quantities of patient data and extract intricate patterns allows clinicians to take more educated decisions based on data, something critical in the prevention and treatment of CVD.

Additionally, the push to machine learning and artificial intelligence in healthcare is concurrent with the broader move to personalized medicine. Conversely, appropriate characterization of patients on the basis of their personal profile treatment with machine learning algorithms can assure personalized prevention and management programs that improve outcomes and decrease healthcare cost. This might involve tailored lifestyle interventions or possibly medications for high-risk patients, or in other cases, reassuring the low-risk people and sparing them unneeded tests and treatments. Such precision, delivered by machine learning models, is a major leap forward in the quest for proactive healthcare versus reactive cardiovascular care as it is being practiced today.

While such models offer their own advantages, it is not straightforward to use these models in clinical settings. Model interpretability: In this study, as also described before, complex models such as Neural Networks are often considered "black boxes", limiting their acceptance among healthcare providers. A reasoning tool that is easily integrated into these models is the SHAP value which, in our opinion, should form an integral part of making these types of predictors transparent and interpretable to clinicians. SHAP values help explain an individualized prediction by attributing portions of the risk to different factors, in essence bridging the explanation gap between clinical black box models and actionable insights for the clinician.

### 6.3. Challenges and Limitations

Given the progress made in understanding machine learning and cardiovascular risk prediction, there still remain challenges. Data quality and availability: So, the most important part of prediction making is good quality data to train your model. The FHS dataset employed for this study is a high-quality, validated dataset; however, using these models on more real-world less structured data may lead to different findings. However, these machine learning models have limited predictive capability due to several problems such as missing data, incorrect measurements, and imbalanced datasets. This will require more robust data collection, data imputation methods or augmentation techniques in future studies.

While emphasizing the transferability of their models to different populations is also important. While the FHS data set was large and contained a wide range of moieties, it may not capture global diversity. Some of the variables not well investigated in this study and affect cardiovascular risk such as ethnicity, socioeconomic status and geographic location. Given the scarcity of recorded experience so far, further research should therefore test these models in wider populations and settings to maximize their application outside Western healthcare.

And the bottom line is that model interpretability remains a substantial challenge to more widespread implementation of machine learning in clinic; SHAP values have helped to clarify the "black-box" nature of models, but additional work is required to refine model outputs as well as make them more user-friendly for clinicians. The challenge arises in integrating models of this complexity into electronic health record (EHR) systems and designing simple interfaces that can interpret these models for use by practitioners.

### 6.4. Future Directions

That work opens quite a few possible doors for future research. Areas of Growth There is a potential for expanding in terms of developing more accurate machine learning models by using bigger and more diverse datasets. With the increasing volume of health data emerging from EHRs, wearable devices, and genomic information, machine learning algorithms can be retrained and developed further to discern new layers of complex patterns in cardiovascular risk. Finally, other methodologies such as deep learning and reinforcement learning have recently been adapted to enhance the predictive performance of these models.

Another highly attractive avenue is to use transfer learning, fine-tuning models trained on one dataset (say the FHS dataset) for other populations with different cardiovascular risk profiles. This approach could increase generalizability across models and minimize the requirement for significant retraining in different healthcare settings. Finally, future research should investigate methods for better interpreting complex models. While SHAP values have great promise, the synthesizing of machine learning tools to ultimately simplify and intuitively distill complex models for clinicians will be critical for furthering the clinical translation of this work.

In conclusion, this research emphasizes the latent capabilities of machine learning approaches in the revolution of cardiovascular risk prediction, outperforming conventional models such as Logistic Regression and show enhanced predictive performance with models including Random Forest and Neural Networks. The integration of SHAP values to explain model predictions is one of the main obstacles for its adoption in a clinical setting, as it represents a vital step towards transferring sophisticated algorithms into valuable clinical interpretable insights. Although data quality, generalizability and model interpretability issues continue to plague the field of machine learning, these results serve

as powerful confirmation that machine learning can be an invaluable weapon in the ongoing fight against cardiovascular disease. With the healthcare industry focusing more on data-guided individualized care, machine learning is becoming a key component in predictive, preventive management of cardiovascular health.

---

## Compliance with ethical standards

### *Disclosure of conflict of interest*

No conflict of interest to be disclosed.

---

## References

- [1] Boulesteix, A. L., et al. (2018). Support Vector Machines for Classifying Patients at High Risk of Cardiovascular Disease. *BMC Medical Informatics and Decision Making*, 18(1), 5.
- [2] Choudhury, M., et al. (2019). Machine Learning Approaches for Cardiovascular Disease Risk Prediction: A Review. *Journal of Healthcare Engineering*, 2019.
- [3] Dey, D., et al. (2020). Random Forest Classifier for Cardiovascular Risk Assessment: A Study on Framingham Heart Study Data. *International Journal of Medical Informatics*, 134, 104053.
- [4] D'Agostino, R. B., et al. (2008). General Cardiovascular Risk Profile for Use in Primary Care: The Framingham Heart Study. *Circulation*, 117(6), 743-753.
- [5] Lee, J., et al. (2021). Big Data in Cardiovascular Research: Insights and Opportunities. *Journal of the American College of Cardiology*, 78(8), 789-804.
- [6] Breiman, L. (2001). Random Forests. *Machine Learning*, 45(1), 5-32.
- [7] Cortes, C., & Vapnik, V. (1995). Support-Vector Networks. *Machine Learning*, 20(3), 273-297.
- [8] Goodfellow, I., Bengio, Y., & Courville, A. (2016). *Deep Learning*. MIT Press.
- [9] Hanley, J. A., & McNeil, B. J. (1982). The Meaning and Use of the Area Under a Receiver Operating Characteristic (ROC) Curve. *Radiology*, 143(1), 29-36.
- [10] Hosmer, D. W., & Lemeshow, S. (2000). *Applied Logistic Regression*. Wiley.
- [11] Iglewicz, B., & Hoaglin, D. C. (1993). How to Detect and Handle Outliers. *SAGE Publications*.
- [12] Jain, A., & Singh, R. (2020). Data Preprocessing Techniques in Data Mining. *International Journal of Advanced Research in Computer Science*, 11(1), 1-5.
- [13] Lundberg, S. M., & Lee, S. I. (2017). A Unified Approach to Interpreting Model Predictions. In Proceedings of the 31st International Conference on Neural Information Processing Systems (Vol. 30, pp. 4765-4774).
- [14] Sokolova, M., & Lapalme, G. (2009). A Systematic Analysis of Performance Measures for Classification Tasks. *Information Processing and Management*, 45(4), 427-437.
- [15] World Health Organization (WHO). (2021). Cardiovascular Diseases (CVDs). Retrieved from [[https://www.who.int/news-room/fact-sheets/detail/cardiovascular-diseases-\(cvds\)](https://www.who.int/news-room/fact-sheets/detail/cardiovascular-diseases-(cvds))].