

Heart Disease Prediction: using Machine Learning

Manas Tripathi, Kshitiz Jain, Jeraldin PJ, Ayush Gupta, Mohd. Armaan

Raj Kumar Goel Institute of Technology, Ghaziabad, UP, India

Abstract: Cardiovascular disease remains one of the foremost global health challenges, contributing substantially to morbidity and mortality rates worldwide. Early detection and accurate prediction are critical in mitigating the impact of heart-related conditions. This study explores the integration of Machine Learning (ML) algorithms with Explainable Artificial Intelligence (XAI) methodologies to enhance the prediction and interpretability of heart disease risk. Utilizing a publicly available dataset from Kaggle, various classification models—such as Decision Trees, Random Forest, Logistic Regression, Support Vector Machines, K-Nearest Neighbours, and Naive Bayes—were trained and evaluated. Among these, the Random Forest model demonstrated the highest predictive performance, closely followed by Logistic Regression. To provide interpretability, SHAP (Shapley Additive explanations) was applied to the Logistic Regression model, offering intuitive visual insights into the influence of individual features on model predictions. The use of SHAP plots further enhances transparency by highlighting the contribution of specific variables to the prediction outcomes. This approach not only facilitates accurate diagnosis but also aids clinicians and stakeholders in understanding the model's decision-making process.

Keywords: Cardiovascular Disease, Machine Learning, Explainable AI, SHAP, Predictive Modelling

I. INTRODUCTION

Cardiovascular disease (CVD) is a broad classification that includes various disorders of the heart and circulatory system, such as coronary artery disease, arrhythmias, heart failure, and congenital heart conditions. These ailments collectively pose a significant challenge to public health, not only in developed nations but increasingly in developing countries like India. Rapid urbanisation, sedentary lifestyles, and dietary transitions have led to a surge in non-communicable diseases, with heart disease at the forefront. Despite significant strides in cardiology and clinical care, CVD remains the leading cause of death globally. According to the World Health Organization (WHO), cardiovascular conditions are responsible for over 17 million deaths annually—a figure that continues to rise.

A multitude of factors contribute to the prevalence of heart disease. These include behavioural risks such as tobacco use, physical inactivity, unhealthy diets, and the harmful use of alcohol. Furthermore, medical conditions like hypertension, diabetes, obesity, and dyslipidaemia significantly increase the risk. Genetic predisposition, age, and gender also influence the likelihood of developing heart-related issues. Early detection and timely intervention are essential in managing the disease burden. Traditionally, diagnosis relies on clinical tests, biomarkers, imaging, and physician expertise. However, these methods can be time-consuming, costly, and sometimes limited in detecting latent or emerging risks across large, diverse populations.

To bridge this gap, Artificial Intelligence (AI), particularly Machine Learning (ML), is emerging as a transformative force in healthcare. ML enables computers to learn patterns from historical data and make predictions or decisions without being explicitly programmed. When applied to medical datasets, ML algorithms can identify complex, non-linear relationships among variables, offering powerful tools for disease risk prediction and diagnosis. However, a significant challenge persists: many ML models operate as "black boxes," producing predictions without offering clear explanations, which hinders clinical trust and widespread adoption.

To address this, the field of Explainable Artificial Intelligence (XAI) has gained traction. XAI methods aim to increase the transparency of ML models by revealing how predictions are made. One prominent method is SHAP (Shapley Additive explanations), which assigns importance scores to input features, thus helping clinicians understand how



specific variables—such as cholesterol levels, age, or blood pressure—affect the model's output. This interpretability is crucial in healthcare, where decisions must be evidence-based and understandable by human experts.

In this context, the present study explores the effectiveness of various supervised ML algorithms in predicting heart disease using a well-established public dataset. By integrating SHAP with one of the top-performing models—Logistic Regression—this work not only seeks to enhance predictive accuracy but also strives to make the model's reasoning more interpretable and trustworthy. The overarching goals are twofold: (1) to develop reliable ML models capable of identifying individuals at risk of heart disease, and (2) to employ XAI tools that elucidate which clinical features most significantly contribute to each prediction. This dual approach promotes greater confidence in AI-assisted healthcare and supports proactive, data-driven medical decision-making.

II. MOTIVATION OF STUDY

Numerous diseases impact populations globally, with heart disease (HD) standing out as a critical health concern due to its substantial contribution to mortality rates among both men and women. According to the World Health Organization (WHO), heart disease is responsible for approximately 17.9 million deaths annually, accounting for nearly 31% of all global fatalities. Despite the availability of machine learning techniques and tools, there remains a lack of efficient and accurate predictive models specifically tailored for the early detection and prognosis of heart disease. At present, no dependable automated system exists that can substantially enhance diagnostic capabilities or mitigate the adverse outcomes associated with the condition.

This research is therefore motivated by the potential of machine learning algorithms to significantly reduce the burden of heart disease. By developing a robust predictive model, the study aims to contribute towards improving the quality of life for individuals at risk and possibly delaying the progression of the disease. The core objectives of this work are twofold: firstly, to design a model capable of predicting the presence of heart disease with high accuracy; and secondly, to identify the most effective classification algorithm among several options.

To achieve these objectives, a comparative analysis will be conducted involving multiple machine-learning classifiers, including Logistic Regression, K-Nearest Neighbours (KNN), Support Vector Machine (SVM), Naïve Bayes, Decision Tree, and Random Forest. The algorithm demonstrating the highest predictive accuracy will be considered the most suitable for heart disease detection in this context.

The structure of the paper is organised as follows: Section 1 provides the introduction. Section 2 presents a review of related works and existing methodologies. Section 3 outlines the flowchart of the proposed framework. Section 4 details the dataset used, and the methodological approach adopted. Section 5 discusses the experimental results and analysis. Finally, Section 6 concludes the study and suggests potential directions for future work.

III. RELATED WORK

Heart disease (HD) is a prevalent condition that commonly affects individuals in middle and later stages of life. A wide range of issues associated with heart disease have been addressed through machine learning (ML) methodologies. Marimuthu et al. reviewed several data-driven analytical techniques for heart disease prediction, highlighting the use of machine learning algorithms such as Decision Trees (DT), Naïve Bayes (NB), K-Nearest Neighbours (KNN), and Support Vector Machines (SVM) [1]. Battula et al. conducted an extensive survey that tabulated and compared various machine learning techniques employed for heart disease prediction since 2012 [2].

Numerous studies have carried out comparative analyses of machine learning algorithms for cardiac disorder classification. These evaluations have consistently demonstrated the predictive capabilities of ML techniques across various heart disease datasets [3]. For instance, one study proposed a decision support system utilising a logistic regression classifier, which achieved a classification accuracy of 77%, indicating its practical potential for clinical deployment.

Machine learning has emerged as a powerful tool for addressing complex healthcare problems by modelling relationships between dependent and independent variables. The healthcare industry, rich in data that is often unmanageable through manual processes, stands to benefit significantly from such automated approaches. Even in technologically advanced nations, heart disease continues to be a leading cause of death, often due to delayed detection



or misidentified risk factors. Machine learning can play a pivotal role in facilitating early risk identification and intervention.

Commonly employed classifiers for heart disease prediction include SVMs, Decision Trees, Logistic Regression, and Naïve Bayes. One study found that SVM achieved the highest predictive accuracy (92.1%), followed by neural networks (91%) and decision trees (89.6%) [4]. Gender and smoking habits were also recognised as key risk factors influencing heart disease occurrence [5].

Other research has further validated the efficacy of algorithms such as DT, NB, and associative classification methods in diagnosing cardiac conditions. Associative classification has demonstrated superior performance when handling unstructured data, offering greater flexibility and accuracy compared to traditional classifiers. Decision Trees were observed to be both user-friendly and precise, while Naïve Bayes emerged as one of the most effective models, closely followed by neural networks and decision trees [6].

Artificial Neural Networks (ANNs) have also been applied to disease prediction tasks. Supervised learning, particularly using the backpropagation algorithm, has shown satisfactory results in diagnostic accuracy. One such study introduced the Intelligent Heart Disease Prediction System (IHGPS), incorporating algorithms such as DT, NB, and Neural Networks (NN) [7]. The Naïve Bayes model demonstrated the highest accuracy at 86.1%, followed by NN with 86.12%, and DT with 80.4%.

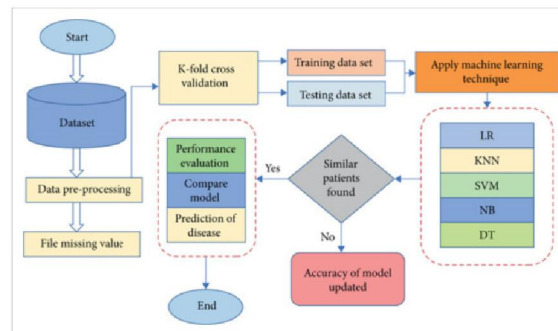
Recent high-accuracy studies frequently employ hybrid approaches combining multiple classification algorithms. The present study aims to enhance the effectiveness of classification algorithms using machine learning techniques. A comparative analysis of Logistic Regression (LR), KNN, SVM, NB, DT, and Random Forest (RF) is conducted, with results evaluated to determine the most accurate model. Furthermore, feature selection is applied to optimise performance, enhancing the applicability of these classifiers within the healthcare domain.

IV. FLOW CHART OF THE PROPOSED FRAMEWORK

The flow chart representing the entire experimental process, from data collection to the development of results, is illustrated in

Figure 1. Initially, data are gathered from the relevant sources and subsequently subjected to a preprocessing phase.

Figure 1: Proposed Framework for Heart Disease Prediction (image to be inserted here)



Data preprocessing aims to eliminate bias, reduce noise, and minimise inaccuracies within the dataset. Upon completion of this stage, the dataset is divided into training and testing subsets.

Various machine learning techniques are then employed to train and evaluate the model. The final step involves generating accurate predictions, which are subsequently compared across the applied algorithms to assess performance and effectiveness.

research aims to show how AI can make a positive difference in mental health and help to create new digital solutions in this area.



V. DATA COLLECTION AND METHODOLOGY

This section outlines the methodology adopted in this research, which includes the dataset selection, preprocessing procedures, and application of machine learning models.

A. DATASET

The study utilises the **Cleveland Heart Disease dataset** from the **UCI Machine Learning Repository**, comprising 520 records and 12 attributes. A detailed description of the dataset is provided in *Table 1*. This dataset has been adopted to develop a machine learning-based diagnostic framework for heart disease prediction.

Attribute	Description	Range
Age	Person's age in years	28–77
Sex	Person's gender (1=Male, 0=Female)	0,1
Chest Pain Type	Type of chest pain experienced	1,2,3,4
Resting BP	Resting blood pressure	0–200
Cholesterol	Serum cholesterol in mg/dl	0–603
Fasting BS	Fasting blood sugar (>120mg/dl = 1)	0,1
Resting ECG	Resting electrocardiographic results	0,1,2
MaxHR	Maximum heart rate achieved	60–202
Exercise Angina	Exercise-induced angina (1=Yes, 0=No)	0,1
Oldpeak	ST depression induced by exercise	-2–6
ST Slope	Slope of peak exercise ST segment	1,2,3
Heart Disease	Diagnosis (1=Disease, 0=Normal)	0,1

The target variable, labelled as *Outcome*, is binary:

"False" represents the absence of heart disease

"True" indicates the presence of heart disease

As with the overall framework, data preprocessing is conducted to reduce inconsistencies and improve model performance. The refined data are then split into training and testing sets for algorithmic processing. A variety of machine learning methods are applied and evaluated, with the aim of generating accurate predictions and identifying the most effective classification approach.

B. DATA PREPROCESSING

When applying machine learning algorithms, data cleaning is essential to optimise accuracy and efficiency. Proper data preparation ensures accurate data representation, enabling machine learning classifiers to be effectively trained and tested. To enable meaningful learning and validation, the data must undergo preprocessing.

StandardScaler is used to normalise features such that each has a mean of 0 and a standard deviation of 1. This ensures all features contribute equally during model training. Alternatively, the MinMaxScaler transforms the data such that all features fall within the range [0, 1]. Additionally, rows containing missing values are removed from the dataset. This study implements all these preprocessing techniques to improve data quality and model performance.

C. DATA CLEANING

The dataset initially contains raw, unprocessed information. Therefore, a series of data cleaning procedures are applied, including the removal of duplicate entries and irrelevant attributes, to enhance the dataset's quality and consistency.

D. FEATURE SELECTION

Feature selection is a crucial process of dimensionality reduction that involves selecting the most relevant variables for classification and prediction tasks. It is an essential step in many well-established classification systems [8]. To improve model accuracy, it is necessary to retain only the features that contribute meaningfully to prediction while eliminating redundant or irrelevant data [9].



By isolating the most informative features, this process facilitates improved classification outcomes. Feature selection is widely used across application domains as it effectively reduces data redundancy without compromising valuable information. It is applied in this research for the following reasons:

- (i) To reduce the complexity and duration of the training process
- (ii) To improve data interpretability by the algorithm
- (iii) To eliminate superfluous features from high-dimensional datasets
- (iv) To enhance prediction accuracy by focusing on essential variables

E. Correlation Matrix

Understanding the relationship between variables is often a useful step in constructing an effective dataset analysis. Correlation is a statistical measure that describes the degree to which two variables move in relation to one another. A positive correlation indicates that both variables increase together, whereas a negative correlation indicates an inverse relationship.

A correlation heatmap based on the dataset is presented in *Figure 2*, which visualises the inter-relationships among features. It reveals that features such as *age*, *gender*, and *Thalch* (maximum heart rate achieved) show a stronger association with the target variable (*outcome*). For instance, the correlation between *age* and *outcome* is 0.11, which is higher than for most other features.



F. K-FOLD CROSS VALIDATION AND DATA SPLITTING

K-fold cross-validation is a widely used method among researchers and practitioners for building reliable models while mitigating information bias. In this study, a 10-fold cross-validation approach has been adopted, where the entire dataset is randomly divided into ten equal-sized partitions. In each iteration, one partition is used as the validation (test) set, while the remaining nine serve as the training set. This process is repeated ten times, with each partition serving as the test set exactly once.

The results from each fold are aggregated using an accumulation function to assess overall model performance. This technique reduces the risk of both overfitting and underfitting by ensuring that each data point is used for both training and validation. Consequently, it enhances the robustness of the machine learning (ML) model by eliminating data bias. For additional analysis, a conventional data split of 70% training and 30% testing is also employed to compare results.

G. APPLICATION OF MACHINE LEARNING TECHNIQUES

Machine learning classification algorithms are employed to differentiate between individuals with and without heart disease. The entire experiment is conducted using Anaconda (2020), a widely adopted open-source Python distribution tailored for data science and machine learning applications. Anaconda facilitates efficient handling of large datasets, predictive analytics, and package management. Spyder (version 3.7.6), an integrated development environment (IDE), is used for code execution and computational tasks.



In machine learning, a model learns from training data and predicts outcomes for new, unseen data. Once the model is trained, it is validated using a separate testing dataset. The final product of this experiment includes the creation of a software tool capable of diagnosing heart disease based on the model's predictions.

Machine learning is broadly divided into two categories: supervised learning and unsupervised learning. In supervised learning, the model is trained using labelled data (i.e., with known outcomes), whereas in unsupervised learning, the model identifies patterns in unlabelled data.

Examples of supervised learning algorithms include:

- (i) Identifying spam emails using pre-labelled datasets
- (ii) Predicting cancer diagnoses from historical patient data
- (iii) Estimating property values based on location and size

Examples of unsupervised learning algorithms include:

- (i) Discovering hidden patterns in scientific datasets
- (ii) Reducing background noise from audio inputs
- (iii) Isolating background music from a song's chorus

In summary, **supervised learning** requires **labelled data**, while **unsupervised learning** operates on **unlabelled data**.

The machine learning algorithms applied in this study include:

Logistic Regression

K-Nearest Neighbour (KNN)

Support Vector Machine (SVM)

Naïve Bayes

Decision Tree (DT)

Random Forest (RF)

H. LOGISTIC REGRESSION

Logistic Regression is a supervised learning technique suitable for both classification and regression tasks. It outputs probabilities between 0 and 1 using the sigmoid function to map predicted values. This function estimates the probability of an instance belonging to a particular class. The model is based on maximum likelihood estimation.

Logistic regression confusion matrix

0	37	3
1	4	60
	0	1

I. K-NEAREST NEIGHBOR (KNN)

KNN classifies new instances by evaluating the distance—often Euclidean—to its 'k' closest data points in the training set. The majority class among the nearest neighbors determines the classification.

K nearest neighbors confusion matrix

0	35	5
1	1	63
	0	1



J. SUPPORT VECTOR MACHINE (SVM)

Support Vector Machines operate by finding the optimal boundary (hyperplane) that separates different class labels in a high-dimensional space. It is highly effective for classification tasks in complex, high-dimensional datasets, such as those found in sentiment analysis or medical diagnostics.

Support vector machine confusion matrix

0	39	1
1	1	63
	0	1

K. NAÏVE BAYES

Naïve Bayes uses Bayes' Theorem to calculate the likelihood of outcomes, assuming all features contribute independently. Classification is based on the posterior probability derived from the training dataset.

Naive bayes confusion matrix

0	35	5
1	4	60
	0	1

L. DECISION TREE (DT)

A Decision Tree is a tree-based structure where internal nodes reflect feature-based decisions, each branch denotes an outcome, and each leaf represents a class label. The root node is the first split based on the most informative feature. DTs are suitable for both categorical and continuous data and offer interpretability. However, they can be prone to overfitting.

Decision tree classifier confusion matrix

0	38	2
1	2	62
	0	1

M. RANDOM FOREST (RF)

Random Forest is an ensemble method that builds several decision trees and aggregates their outcomes to improve accuracy and reduce overfitting. Each tree is trained on a random subset of the dataset, and the final prediction is made by aggregating the outputs of individual trees.

Random forest confusion matrix

0	40	0
1	1	63
	0	1



N. PERFORMANCE EVALUATION

The performance of the before mentioned machine learning models is evaluated using the Cleveland Heart Disease dataset. The assessment criteria include standard classification metrics: **Accuracy**, **Precision**, **Recall**, **F1-Score**, and **Matthews Correlation Coefficient (MCC)**. These metrics are calculated using the following formulas:

$$\text{Accuracy} = (TP + TN) / (TP + TN + FP + FN)$$

$$\text{Precision} = TP / (TP + FP)$$

$$\text{Recall} = TP / (TP + FN)$$

$$F1 = 2 * (\text{Precision} * \text{Recall}) / (\text{Precision} + \text{Recall})$$

$$MCC = [(TP * TN) - (FP * FN)] / \sqrt{[(TP + FP)(TP + FN)(TN + FP)(TN + FN)]}$$

Where:

TP = True Positives

TN = True Negatives

FP = False Positives

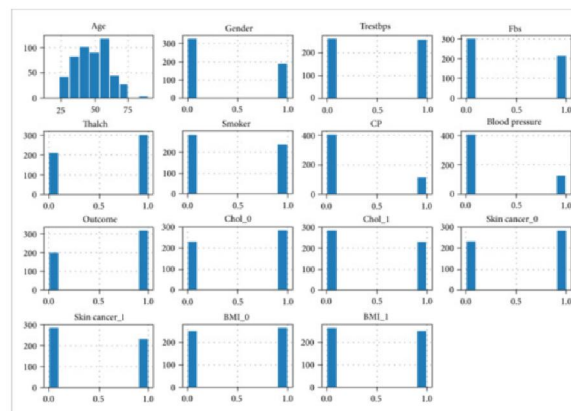
FN = False Negatives

These evaluation metrics allow for a comprehensive comparison of classification performance across different machine learning techniques.

VI. RESULT AND ANALYSIS

Numerous classification models and their statistical analyses are provided in this section of the research. On the Cleveland heart disease data, we assess the effectiveness of LR, KNN, SVM, NB, RF, and DT in the first stage. In this research, we investigated different machine learning algorithms for the prediction of cardiac disease using an experimental and analytical techniques.

Figure 3. displays the histogram that was created in addition to the plots that depict the distribution of each dataset attribute.



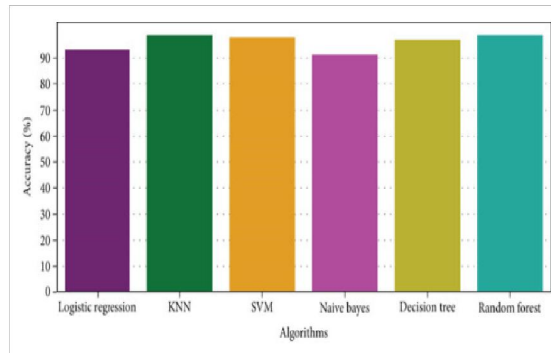
A. MODEL ACCURACY

The prediction models have been developed using twelve selected features from the dataset, and their performance has been assessed based on classification accuracy. To facilitate a comparative analysis of the employed machine learning algorithms, their respective accuracies are illustrated in Figure 4.

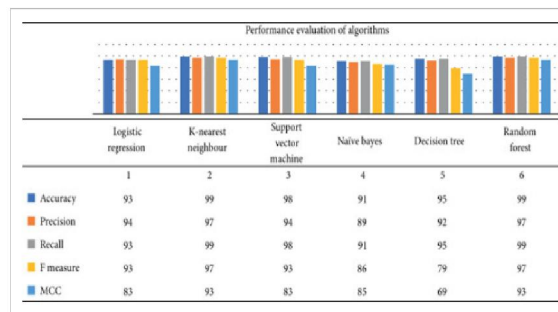
This comparison allows for a clearer understanding of the variations in predictive performance across different modelling techniques. It is evident from the results that the Random Forest (RF) and K-Nearest Neighbour (KNN) algorithms outperform the others in terms of accuracy, demonstrating superior consistency and reliability in classifying the dataset.

The accompanying bar graph (Figure 4) visually represents the accuracy of each algorithm, thereby highlighting the relative effectiveness of the models under study.





Six machine learning algorithms were employed in this study for the prediction of heart disease. The interrelationships among the features utilized in the dataset are illustrated in the scatterplot shown in Figure 5. Each point on the plot represents a data instance, with its position along the X and Y axes corresponding to the values of specific features. This visual representation facilitates an understanding of the distribution and correlation patterns within the dataset. Additionally, as illustrated in Figure 12, several other statistical metrics have been computed to evaluate the performance of the machine learning classifiers. These evaluation parameters include accuracy, precision, recall, F1-score, and Matthews Correlation Coefficient (MCC). Each metric provides a distinct perspective on the effectiveness of the models in predicting heart disease, thereby enabling a comprehensive performance comparison across the various algorithms employed



B. COMPARATIVE ANALYSIS

Table 2 presents a comparative evaluation of the proposed framework against several pertinent studies from the existing literature, with respect to methodologies applied, datasets utilised, and analytical approaches adopted. A consistent set of cardiac predictors was observed across all studies, enabling a meaningful comparison with the current investigation. The results indicate that the proposed framework yielded superior performance across multiple evaluation metrics, particularly in terms of prediction accuracy. This enhanced performance can be attributed to the integration of several advanced data processing techniques. Specifically, the use of data imputation addressed missing values effectively, while the implementation of scatterplot-based outlier detection and replacement improved data quality. Moreover, data transformation methods—namely standardisation and normalisation—ensured feature scaling was handled robustly. The application of the K-fold cross-validation technique during model development further contributed to the reliability and generalisability of the results, distinguishing the proposed approach from other comparable studies. These methodological improvements collectively underpin the observed performance gains in heart disease prediction.



Table 2. Comparison with the existing systems.

Authors	Techniques used	Dataset	Accuracy
Chauhan et al. [24]	RF, LR, SVM	PIDD data set	78%
Pouriyeh et al. [25]	KNN, SVM, RF, NB	Cleveland data	89%
Kedia et al. [26]	LR, KNN, SVM, DT	UCI data	90%
Atallah et al. [27]	SVM, DT, KNN, RF	Cleveland data	87%
Proposed method	LR, KNN, SVM, NB, DT, RF	Cleveland Heart from UCI's machine learning	99.04%

VII. CONCLUSION

Comparing various ML for the early detection of heart disease is the main contribution, preprocessing techniques were used to enhance the dataset's quality. With the primary objectives being the handling of corrupt and missing values as well as the removal of outliers to predict illness. Additionally, we used a variety of machine learning techniques, and the outcomes were compared using various statistical metrics. The experimental finding indicates a 70:30 ratios between testing and training the data. In this study, we perform 10-fold cross-validation to several machine learning methods, and we find that random forest and k-nearest neighbor are 99.04% accurate compared to other algorithms. Future work can be carried out using various combinations of machine learning methodologies to enhance prediction techniques. To better comprehension of the critical features and increase the precision of heart disease prediction, new feature selection approaches can also be developed.

VIII. FUTURE WORK

The future scope of this study spans multiple dimensions, including technological advancement, enhancement of patient outcomes, and overall improvement in healthcare delivery systems. As machine learning and artificial intelligence (AI) continue to evolve, there are several promising directions for future research and development.

One key area involves the integration of AI-powered algorithms to deliver personalised recommendations and treatment plans. By leveraging large volumes of patient data, future systems could be designed to predict individual health risks more accurately and propose tailored interventions. Such innovations could contribute to a reduction in hospital readmissions, better medication adherence, and more effective management of chronic conditions such as diabetes and cardiovascular diseases.

Additionally, future studies could explore the patient experience and satisfaction when using digital health technologies. This includes evaluating the usability and effectiveness of platforms that offer access to dietary guidance, direct communication with healthcare professionals, and health monitoring tools. Research methods such as patient surveys, interviews, or focus groups would provide valuable insight into user preferences, concerns, and acceptance of such technologies, especially within diverse socio-economic backgrounds in India.

Another significant area of interest lies in examining the readiness and adaptability of healthcare professionals towards adopting these technologies in clinical settings. Investigating the barriers—such as limited technical training, infrastructural constraints, or scepticism regarding AI decisions—as well as the enablers for successful integration, would help in designing better support systems. Equipping healthcare providers with intuitive tools and evidence-based practices will be crucial for widespread adoption.

In conclusion, future work should aim at building a more inclusive, intelligent, and patient-centric healthcare system. Combining explainable AI with ethical considerations, user education, and seamless healthcare infrastructure integration will be essential for making these advancements both effective and trustworthy.

REFERENCES

- [1]. Animesh Hazra, Arkomita Mukherjee, Amit Gupta, Asmita Mukherjee, "Heart Disease Diagnosis and Prediction Using Machine Learning and Data Mining Techniques: A Review", Research Gate Publications, July 2017, pp.2137-2159.



- [2]. V. Krishnaiah, G. Narsimha, N. Subhash Chandra, "Heart Disease Prediction System using Data Mining Techniques and Intelligent Fuzzy Approach: A Review", International Journal of Computer Applications, February 2016.
- [3]. Guizhou Hu, Martin M. Root, "Building Prediction Models for Coronary Heart Disease by Synthesizing Multiple Longitudinal Research Findings", European Science of Cardiology, 10 May 2005.
- [4]. TMythili, Dev Mukherji, Nikita Padaila and Abhiram Naidu, "A Heart Disease Prediction Model using SVM-Decision Trees- Logistic Regression (SDL)", International Journal of Computer Applications, vol. 68, 16 April 2013.
- [5]. <https://www.medicalnewstoday.com/articles/257484.php>.
- [6]. Nimai Chand Das Adhikari, Arpana Alka, and Rajat Garg, "HPPS: Heart Problem Prediction System using Machine Learning".
- [7]. Heart disease CSV file from <https://github.com/MainakRepositor/Cardiac-Arrest-Predictor>
- [8]. C. Feng, S. Wu, and N. Liu, "A user-centric machine learning framework for cyber security operations center," in 2017 IEEE International Conference on Intelligence and Security Informatics (ISI), Beijing, China, Jul. 2017, pp. 173–175, doi: 10.1109/ISI.2017.8004902.
- [9]. S. B. Kotsiantis, I. Zaharakis, and P. Pintelas, "Supervised machine learning: A review of classification techniques," Emerging artificial intelligence applications in computer engineering, vol. 160, no. 1, pp. 3–24, 2007.
- [10]. S. B. Kotsiantis, I. D. Zaharakis, and P. E. Pintelas, "Machine learning: a review of classification and combining techniques," ArtifIntell Rev, vol. 26, no. 3, pp. 159–190, Nov. 2006, doi: 10.1007/s10462-007-9052-3.
- [11]. C. Surv, M. N. Murty, P. J. Flynn, A. K. Jain, and P. J. Flynn, And. 1999.
- [12]. D. Sharma and N. Kumar, "A Review on Machine Learning Algorithms, Tasks and Applications," vol. 6, pp. 2278–1323, Oct. 2017.
- [13]. K. Pahwa and N. Agarwal, "Stock Market Analysis using Supervised Machine Learning," in 2019 International Conference on Machine Learning, Big Data, Cloud and Parallel Computing (COMIT Con), Faridabad, India, Feb. 2019, pp. 197–200, doi: 10.1109/COMITCon.2019.8862225.
- [14]. M. Pérez-Ortiz, S. Jiménez-Fernández, P. A. Gutiérrez, E. Alexandre, C. Hervás-Martínez, and S. Salcedo-Sanz, "A Review of Classification Problems and Algorithms in Renewable Energy Applications," Energies, vol. 9, no. 8, Art. no. 8, Aug. 2016, doi: 10.3390/en9080607.
- [15]. Anuradha and G. Gupta, "A selfexplanatory review of decision tree classifiers," in International Conference on Recent Advances and Innovations in Engineering (ICRAIE-2014), Jaipur, India, May 2014, pp. 1–7, doi: 10.1109/ICRAIE.2014.6909245.
- [16]. S. Patil and U. Kulkarni, "Accuracy Prediction for Distributed Decision Tree using Machine Learning approach," in 2019 3rd International Conference on Trends in Electronics and Informatics (ICOEI), Apr. 2019, pp. 1365–1371, doi: 10.1109/ICOEI.2019.8862580.
- [17]. N. S. Ahmed and M. H. Sadiq, "Clarify of the random forest algorithm in an educational field," in 2018 International Conference on Advanced Science and Engineering (ICOASE), 2018, pp. 179–184.
- [18]. D. Zeebaree, Gene Selection and Classification of Microarray Data Using Convolutional Neural Network. 2018.
- [19]. O. M. Salih Hassan, A. Mohsin Abdulazeez, and V. M. Tiryaki, "Gait-Based Human Gender Classification Using Lifting 5/3 Wavelet and Principal Component Analysis," in 2018 International Conference on Advanced Science and Engineering (ICOASE), Duhok, Oct. 2018, pp. 173–178, doi: 10.1109/ICOASE.2018.8548909.
- [20]. R. Zebari, A. Abdulazeez, D. Zeebaree, D. Zebari, and J. Saeed, "A Comprehensive Review of Dimensionality Reduction Techniques for Feature Selection and Feature Extraction," Journal of Applied Science and Technology Trends, vol. 1, no. 2, pp. 56–70, 2020.

