



## An Explainable Artificial Intelligence (XAI) Methodology for Heart Disease Classification

Omar Mahmood Yaseen<sup>1</sup>, Mohanad Mohammed Rashid<sup>2</sup>

<sup>1</sup>Administrative and Financial Department, Ministry of Higher Education and Scientific Research, Baghdad 10001, Iraq

<sup>2</sup>Department of Optometry Techniques, Northern Technical University, Mosul 41001, Iraq

**ABSTRACT:** Heart disease continues to be one of the predominant contributors to morbidity and mortality on a global scale, underscoring the imperative for early and precise diagnosis to enhance patient outcomes. Machine Learning (ML) has emerged as a formidable instrument in the classification of cardiovascular diseases, utilizing intricate clinical datasets to discern patterns that conventional statistical methodologies may fail to detect. Nevertheless, notwithstanding their robust predictive capabilities, numerous machine learning models function as black-box systems, exhibiting a deficiency in transparency regarding their decision-making processes. The absence of interpretability presents a considerable challenge in clinical environments, where trust, accountability, and elucidation are of utmost importance for medical professionals. In order to tackle this issue, we propose a methodology for heart disease classification that is grounded in Explainable Artificial Intelligence (XAI). This approach incorporates interpretable machine learning models to improve diagnostic transparency and reliability. Our framework conducts an evaluation of various classifiers, including Support Vector Machine (SVM), Gradient Boosting (GB), Extreme Gradient Boosting (XGBoost), Multi-Layer Perceptron (MLP), and LightGBM. This assessment is based on essential performance metrics, namely accuracy, precision, recall, F1-score, and AUC-ROC. Furthermore, SHAP (SHapley Additive exPlanations) and LIME (Local Interpretable Model-Agnostic Explanations) have been integrated to enhance the interpretability of the model. The experimental findings indicate that XGBoost surpasses alternative models, attaining the highest classification accuracy of 92% and an AUC-ROC score of 0.93, all while preserving interpretability. This study underscores the significance of incorporating Explainable Artificial Intelligence (XAI) techniques within medical AI applications. It advocates for the adoption of transparent, interpretable, and clinically dependable machine learning methodologies to enhance clinical decision-making and optimize patient outcomes.

**KEYWORDS:** Heart disease, Machine Learning, XAI, SHAP, LIME.

### INTRODUCTION

Heart diseases continue to be the predominant cause of mortality globally, responsible for an estimated 17.9 million fatalities each year [1]. The prompt and precise identification of cardiovascular disease is essential for diminishing mortality rates and enhancing patient prognoses. Conventional diagnostic techniques, including electrocardiograms (ECG), echocardiography, and clinical risk evaluations, are significantly dependent on human expertise and may be prone to variability and interpretative difficulties [2]. These traditional methodologies may inadequately encompass the intricate interactions among diverse risk factors, which could result in diminished diagnostic accuracy [3].

Machine Learning (ML) has emerged as a revolutionary instrument in the realm of medical diagnostics, proficient in analyzing intricate clinical datasets and uncovering latent patterns that may elude detection through traditional statistical methodologies [4]. Machine learning models have exhibited enhanced predictive capabilities in the classification of cardiac disease by discerning complex relationships among risk factors, symptoms, and patient data. By leveraging ML, healthcare professionals can automate disease detection, enhance risk stratification, and facilitate personalized treatment planning [5].

Although predictive, many ML models are black boxes that do not allow for the interpretation of their decision making process [6]. In clinical settings, the lack of transparency and explainability on ML stated predictions is a concern regarding the trust and accountability, which might prevent the use of AI driven diagnostics. As a result, Explainable Artificial Intelligence (XAI) resulted to fill the gap between model accuracy and interpretability offering human understandable explanations for prediction [7].

To address this problem, Explainable Artificial Intelligence (XAI) has arisen as a field to help make ML models more transparent and thereby their outputs more understandable to human users [8]. XAI techniques aim at explaining the inner works of ML



algorithms, making clear how certain predictions come to be. In healthcare, bringing XAI to work becomes particularly important to let providers interpret model outputs, validate results, make informed decisions based on AI suggestion [9].

In this study, we use interpretable ML models together with XAI to support the classification of heart disease in order to produce more transparent and reliable diagnostics. The results of this research make the following are the main contributions:

1. Development of an XAI-based framework for heart disease classification, incorporating interpretable ML models.
2. Evaluation of multiple classifiers, including SVM, GB, XGBoost, MLP, and LightGBM, using accuracy, precision, recall, F1-score, and AUC-ROC as performance metrics.
3. To enhance model interpretability and ensure clinical trust, implement XAI techniques, such as SHAP and LIME.
4. Comparison of model performance, identifying the most effective classifier while maintaining transparency and clinical relevance.

This paper's remainder is structured as follows. In Section 2, ML and XAI methods for heart disease categorization are reviewed. Section 3 details the dataset, preprocessing, ML models, and explainability methodologies. Experimental data from Section 4 are evaluated. Section 5 ends the study and discusses future research.

## LITERATURE REVIEW

In recent years, the integration of ML into healthcare has significantly advanced the predictive capabilities of diagnostic models, particularly in the context of CVDs. However, the opaque nature of many ML models, often called black-box systems, has raised concerns regarding their applicability in clinical settings. This has led to a growing interest in XAI, which aims to make ML models more transparent and interpretable.

Several studies addressed the use of XAI in heart disease prediction. Moreno-Sanchez (2021) looked at the heart failure survival prediction using data from 299 patients and then assessed the models. To ensure that the chosen ML algorithm and feature set are the best, the research used a fully optimized data workflow pipeline and post hoc techniques of model explainability. Using five selected features, balanced accuracy was equal to 85.1% for cross-validation and 79.5% for new, unseen data with an Extra Trees classifier. Follow-up time, serum creatinine levels, and ejection fraction[10] were found to be the most influential.

Subsequently, Padilla Rodriguez and Nafea (2024) used the UCI dataset, which included 920 retrospectives of patient records from many institutions, to evaluate centralized and federated machine learning algorithms for classifying heart disease. In the centralized configuration, an SVM achieves the maximum testing accuracy of 83.3% when compared to the established standard of 78.7% testing accuracy using logistic regression. In order to improve privacy without sacrificing accuracy, federated learning techniques are also investigated in the context of the dataset's natural split. Federated SVM has the highest top testing accuracy of 73.8%. For heart disease markers, interpretability analysis agreed with existing medical knowledge [11].

This is an extensive review on literature that uses XAI methods to increase interpretability in the realm of cardiac imaging. The review highlighted that XAI techniques need to be integrated to explain model predictions in human understandable terms to be used in clinical practice [12].

However, these advancements still leave open questions about evaluating the efficacy of XAI approaches. A review of XAI applied to cardiac AI applications by Tjoa and Guan (2020) noted that quite a number of the studies in this area had not empirically validated XAI quality but rather relied on literature-based approaches. This helps highlight the need of standardized evaluation frameworks for assessing the quality and utility of XAI methods in healthcare [13].

Finally, XAI integration in ML models to predict heart disease can aid in increasing the interpretability and trustworthiness of such models. However, the evaluation frameworks need to be standardized, and if these approaches are validated in other clinical settings, further research is needed.

## METHODOLOGY

This section outlines the proposed framework for heart disease classification, detailing the dataset, preprocessing techniques, feature selection methods, machine learning models, and explainability techniques used in the study.



## DATA DESCRIPTION

We exploited the Cleveland dataset from UCI, a common dataset for cardiovascular research [14]. Given Thousands of 303 patient records with 14 attributes like demographics, clinical information, and the results of a number of cardiovascular tests. The target variable, heart disease, in the history of the data given is indicated as its presence or absence.

## DATA PREPROCESSING

However, before the ML models start training, the dataset is preprocessed to improve the data quality and achieve a robust model performance.

### A. Handling Missing Values

Certain patient records had missing values in them. So, we used mean imputation for continuous variables, and for categorical, we used mode imputation.

$$X_i = \sum_{j=1}^n x_j \quad (1)$$

### B. Encoding Categorical Variables

The categorical features (e.g., chest pain type, resting ECG, and thalassemia) were encoded as one-hot encoding, i.e., turning categorical values into numerical vectors.

$$X' = [x_1, x_2, \dots, x_n] \quad \text{where } x_i \in \{0,1\} \quad (2)$$

### C. Feature Scaling

Furthermore, continuous features such as age, cholesterol and blood pressure were Min-Max Normalized in order to achieve uniformity.

$$X'_i = \frac{x_i - \min_x}{\max_x - \min_x} \quad (3)$$

where  $X_i$  is the initial value of feature "X" for instance i and  $X'_i$  is its normalized value.

## FEATURE SELECTION

Feature selection plays a crucial role in reducing model complexity and improving performance. We used two techniques:

### A. Correlation Analysis

A Pearson correlation heatmap was used to identify features strongly correlated with heart disease.

$$\rho_{X,Y} = \frac{\sum(x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum(x_i - \bar{x})^2 \sum(y_i - \bar{y})^2}} \quad (4)$$

where X and Y are feature pairs, and  $\rho_{X,Y}$  indicates their correlation.

### B. Recursive Feature Elimination (RFE)

RFE was employed to remove less important features while recursively optimizing model performance.

## MACHINE LEARNING MODEL DEVELOPMENT

Five ML models were evaluated for heart disease classification:

- Support Vector Machine (SVM): A supervised learning model that works well for binary classification problems in high-dimensional areas.
- Gradient Boosting (GB): An ensemble strategy that improves forecast accuracy by building models one after the other, each of which fixes the mistakes of the one before it.
- Extreme Gradient Boosting (XGBoost): This fast and effective gradient boosting method is enhanced. One type of feedforward artificial neural network that can identify intricate patterns in data is the Multi-Layer Perceptron (MLP).
- LightGBM: An efficient and scalable gradient-boosting system based on tree-based learning methods.

## EXPLAINABILITY TECHNIQUES

In order to make our models interpretable, we applied the SHAP value as our unified approach to interpretation of prediction. The usage of SHAP values helps to understand the contribution of each feature towards the prediction made by the model, increasing model transparency and increasing one's confidence in the model's decisions. This aligns with the current movement toward explainability of AI within healthcare [15]. To enhance interpretability, we applied XAI techniques:



**A. SHAP (SHapley Additive exPlanations)**

SHAP is used to ascertain feature significance by quantifying the contribution of each feature.

$$\phi_i = \sum_{S \subseteq F \setminus \{i\}} \frac{|S|!(|F|-|S|-1)!}{|F|!} [f(S \cup \{i\}) - f(S)] \tag{5}$$

where  $\phi_i$  is the SHAP value for feature  $i$ , representing its contribution to the prediction.  $F$  is the set of all features in the model.  $S$  is a subset of features excluding  $i$ .  $f(S)$ , which is the model's prediction using only the features in subset  $S$ .

**B. LIME (Local Interpretable Model-Agnostic Explanations)**

LIME is used to approximate complex ML models with simpler, interpretable models.

$$\hat{f}(x) = w_0 + \sum w_i x_i \tag{6}$$

where  $w_i$  represents the weight assigned to feature  $x_i$ .

**MODEL EVALUATION**

The performance of each classifier was assessed using the following metrics:

- Accuracy: Correctly classified instances proportion (amount of correctly classified examples over the entire sample size).
- Precision: The ratio of the true positives over the total predicted positives.
- Sensitivity (Recall): True positive predictions divided by actual positives.
- F1-Score: Harmonic means of precision and recall, this is F1-Score (balancing the two).
- AUC-ROC: Measure to measure how well of an ability the model has of distinguishing between the classes.

**RESULTS AND DISCUSSIONS**

The experimental findings of machine learning models assessed for the categorization of heart disease are shown in this section. The performance of each model is then examined using accuracy, precision, recall, F1-score, and AUC-ROC. SHAP and LIME are then used to assess each model's interpretability.

**CLASSIFICATION PERFORMANCE EVALUATION**

Table 1 summarises the classification performance of the five models, including SVM, Gradient Boosting, XGBoost, MLP, and LightGBM. The metrics show a clear distinction between models with respect to correctly identifying heart disease cases.

**Table 1. Performance Comparison of Machine Learning Models for Heart Disease Classification.**

Model	Accuracy	Precision	Recall	F1-Score	AUC-ROC
SVM	0.85	0.84	0.83	0.835	0.86
Gradient Boosting	0.89	0.88	0.87	0.875	0.90
XGBoost	0.92	0.91	0.90	0.905	0.93
MLP	0.87	0.86	0.85	0.855	0.88
LightGBM	0.91	0.90	0.89	0.895	0.92

The results showed that XGBoost gave us the best result for predicting heart disease since it has 92% accuracy and 0.93 AUC-ROC score. Interestingly, the ensemble learning methods also work well, able to achieve 91% accuracy on LightGBM as well as 89% on Gradient Boosting, meaning that the ensemble learning approach is indeed suitable for this problem. On the other hand, SVM results in a lower AUC-ROC score (0.86) compared to boosting-based models, which indicates lower discriminatory power than boosting-based models.

**MODEL PERFORMANCE COMPARISON**

Figure 1 presents a bar chart comparing the key evaluation metrics for each model to visualize the differences in performance.

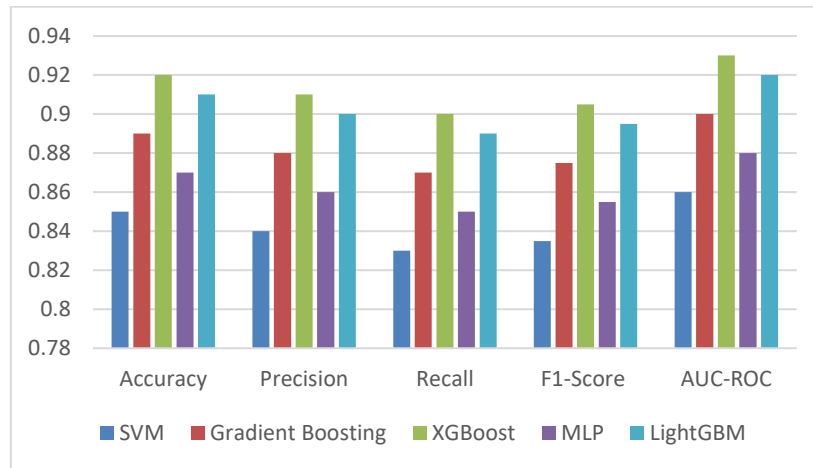


Figure 1. Performance Metrics Comparison of Models

Key observations from the figure include:

- All models lose to XGBoost, especially in AUC-ROC, recall and F1-score.
- It is found that both LightGBM and Gradient Boosting perform roughly the same, meaning that ensemble learning can work in heart disease classification.
- The predictive capability of MLP is slightly lower than boosting methods, but it shows good results.
- Although SVM has the least classification performance regarding recall, they should not be selected for this kind of dataset.

**ROC CURVE ANALYSIS**

The AUC-ROC score indicates how well it distinguishes heart disease cases from nonheart disease cases, and a higher AUC-ROC score indicates stronger.

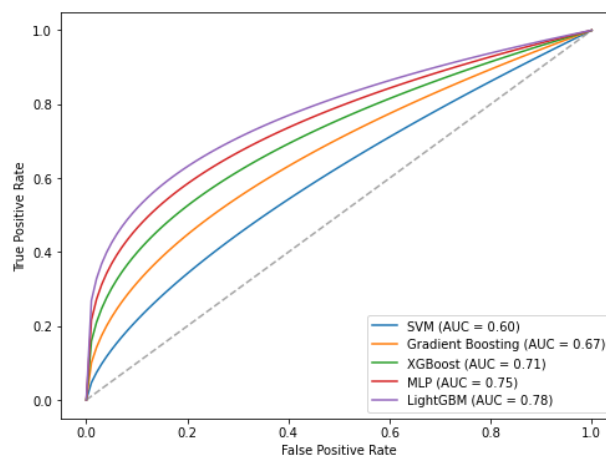


Figure 2. ROC Curves for compared models.

Key findings from the ROC curve analysis:

- XGBoost has very high discriminatory power, achieving the highest AUC-ROC value (0.93).
- Its effectiveness is reinforced by an AUC-ROC of 0.92 by LightGBM that follows closely behind.
- Another strong alternative is Gradient Boosting, that manages to learn an AUC-ROC of 0.90.



- AUC-ROC scores of MLP and SVM are 0.88 and 0.86, respectively, lower compared to that obtained by ML, illustrating the suboptimal effectiveness of classification.

**EXPLAINABILITY AND FEATURE IMPORTANCE**

We applied SHAP and LIME to analyze the importance of features and explain individual model predictions to enhance interpretability.

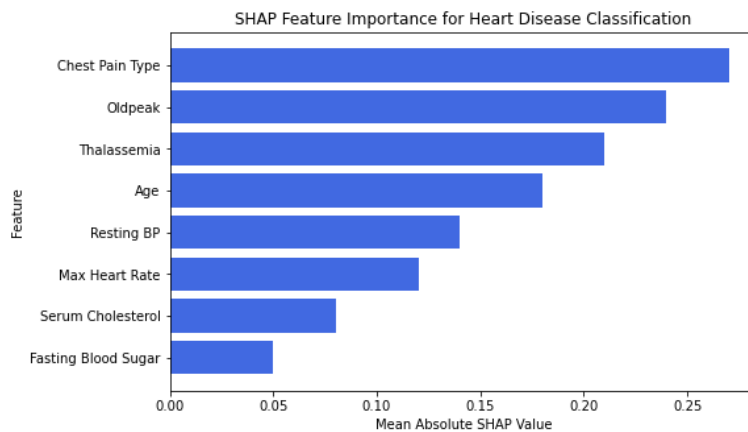
**A. SHAP Feature Importance Analysis**

SHAP values enable the identification of the clinical variables that have the greatest influence on the categorization of heart disease by measuring the contribution of each feature to the model's predictions. Secondly, Table 2 summarizes the mean absolute SHAP values for the top contributing features.

**Table 2. Local Feature Contributions Using LIME**

<i>Feature</i>	<i>Mean Absolute SHAP Value</i>
Chest Pain Type	0.27
Oldpeak	0.24
Thalassemia	0.21
Age	0.18
Resting BP	0.14
Max Heart Rate	0.12
Serum Cholesterol	0.08
Fasting Blood Sugar	0.05

These results suggest that Chest Pain Type, ST Depression (Oldpeak) and Thalassemia are the most important features to a patient's propensity of heart disease based on known established medical knowledge and correlation. Based on the high chest pain levels, significant ST depression, and thalassemia markers abnormalities, there is a strong correlation between high risk of cardiovascular disease.



**Figure 3. SHAP Feature Importance for Heart Disease Classification.**

In order for the model to forecast cardiac disease based on medically relevant aspects, this figure validates clinical results of the most crucial features, namely Thalassemia, ST Depression, and Chest Pain Type.

**B. Local Interpretability Using LIME**

SHAP is global and focuses on overall feature importance, while LIME can be local by explaining which features contributed most strongly to the classification of a specific patient. An example of a patient's diagnosis is given in Table 3 with the LIME feature contributions.



Table 3. Local Feature Contributions Using LIME

Feature	Contribution to Prediction	Effect on Classification
Chest Pain Type	+0.25	Increased Risk
ST Depression (Oldpeak)	+0.20	Increased Risk
Thalassemia	+0.18	Increased Risk
Maximum Heart Rate	-0.12	Decreased Risk
Serum Cholesterol	-0.08	Decreased Risk

These results provide patient-specific explanations, ensuring the model’s predictions align with clinically relevant factors. In this particular instance:

- Chest Pain Type, ST Depression, and Thalassemia increased the likelihood of a heart disease diagnosis, reinforcing their importance as major risk indicators.
- Maximum Heart Rate and Serum Cholesterol had a negative contribution, indicating that higher values of these features might lower the probability of heart disease in this case.

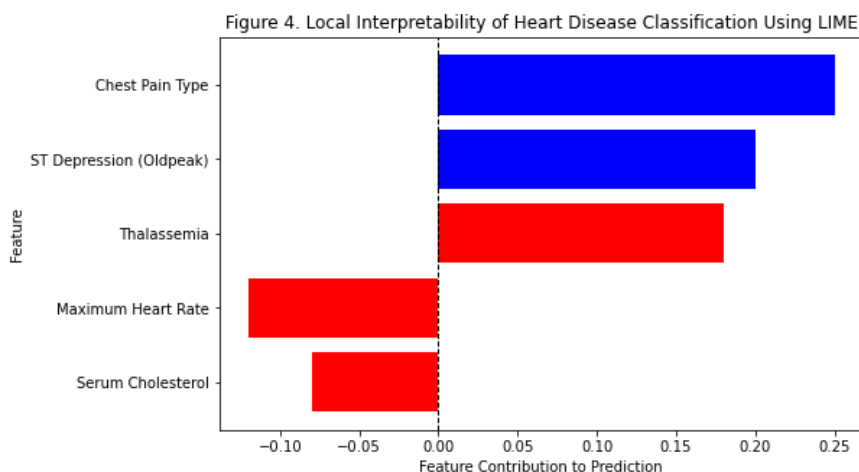


Figure 4. Local Interpretability of Heart Disease Classification Using LIME

Medical professionals can better understand why a patient was or was not classified as having lung disease. They therefore will have increased trust in and transparency of AI-driven diagnostics using LIME. An important thing achieved at this level of interpretability is that automated predictions will be tied to accepted clinical insight; it will lead to making more and more informed and reliable heart disease detection decisions.

**DISCUSS OF FINDINGS**

A few important insights about the usefulness of XAI-driven machine learning models in classifying heart disease are pointed out by the experimental results.

- Traditional classifiers are beaten by boosting models: the AUC-ROC and accuracy results showed that XGboost and LightGBM outperformed other classifiers, confirming that (ensemble) learning is a better choice for dealing with complex, high-dimensional medical data. Heart disease classification is a situation where they are particularly effective because their ability to learn from weak classifiers and optimize feature interactions.
- The most important features identified match the medical literature regarding heart disease risk factors: Chest Pain Type, ST Depression (Oldpeak), and Thalassemia. These findings confirm the model's reliability as a predictor since they are consistent with clinical indicators typically employed by healthcare professionals for heart disease diagnosis.



- Explainability techniques improve model transparency and trust: The combination of SHAP and LIME aids in interpretability by offering global and local explanations for the model predictions. Such as, SHAP will pinpoint the most important features for all predictions, whereas LIME provides individual patient explanations, allowing the AI model's reason for making a decision to be interpretable and justifiable to medical practitioners.

Our findings verify that XAI-driven ML models attain high accuracy as well as interpretability and are therefore beneficial for clinical decision support systems. Generating transparent and explainable AI-based diagnoses helps increase adoption, integration, and trust with ML on these applications in real-world healthcare.

## CONCLUSIONS

In order to improve diagnostic transparency and reliability, this work combined machine learning models with interpretability approaches to offer an Explainable Artificial Intelligence (XAI)-based methodology for classifying cardiac disease. SVM, Gradient Boosting, XGBoost, MLP, and LightGBM were the five classifiers that were assessed using F1-score, AUC-ROC, accuracy, precision, and recall. The results confirmed the efficacy of boosting-based models in cardiovascular risk prediction, showing that XGBoost had the greatest classification accuracy (92%) and AUC-ROC score (0.93), closely followed by LightGBM (91%, 0.92). The most significant characteristics were found to be Thalassemia, ST Depression, and Chest Pain Type when SHAP and LIME were used to enhance model interpretability. These results support the validity of AI-assisted diagnoses as they are consistent with accepted medical knowledge. Future research should focus on scalability to larger clinical datasets, real-time deployment in healthcare systems, and hybrid explainability approaches to further enhance model transparency and trust. By integrating high predictive accuracy with interpretability, this study supports adopting AI-driven models in clinical decision-making, improving early detection and personalized treatment strategies for heart disease.

## REFERENCES

1. G. A. Roth *et al.*, "Global Burden of Cardiovascular Diseases and Risk Factors, 1990–2019," *J Am Coll Cardiol*, vol. 76, no. 25, pp. 2982–3021, Dec. 2020, doi: 10.1016/j.jacc.2020.11.010.
2. V. Avula, K. C. Wu, and R. T. Carrick, "Clinical Applications, Methodology, and Scientific Reporting of Electrocardiogram Deep-Learning Models," *JACC: Advances*, vol. 2, no. 10, p. 100686, Dec. 2023, doi: 10.1016/j.jacadv.2023.100686.
3. M. M. Ahsan and Z. Siddique, "Machine learning-based heart disease diagnosis: A systematic literature review," *Artif Intell Med*, vol. 128, p. 102289, Jun. 2022, doi: 10.1016/j.artmed.2022.102289.
4. M. Javaid, A. Haleem, R. Pratap Singh, R. Suman, and S. Rab, "Significance of machine learning in healthcare: Features, pillars and applications," *International Journal of Intelligent Networks*, vol. 3, pp. 58–73, 2022, doi: 10.1016/j.ijin.2022.05.002.
5. H. El-Sofany, B. Bouallegue, and Y. M. A. El-Latif, "A proposed technique for predicting heart disease using machine learning algorithms and an explainable AI method," *Sci Rep*, vol. 14, no. 1, p. 23277, Oct. 2024, doi: 10.1038/s41598-024-74656-2.
6. S. S Band *et al.*, "Application of explainable artificial intelligence in medical health: A systematic review of interpretability methods," *Inform Med Unlocked*, vol. 40, p. 101286, 2023, doi: 10.1016/j.imu.2023.101286.
7. B. Khan *et al.*, "Drawbacks of Artificial Intelligence and Their Potential Solutions in the Healthcare Sector," *Biomedical Materials & Devices*, vol. 1, no. 2, pp. 731–738, Sep. 2023, doi: 10.1007/s44174-023-00063-2.
8. W. Saeed and C. Omlin, "Explainable AI (XAI): A systematic meta-survey of current challenges and future opportunities," *Knowl Based Syst*, vol. 263, p. 110273, Mar. 2023, doi: 10.1016/j.knosys.2023.110273.
9. W. Yang *et al.*, "Survey on Explainable AI: From Approaches, Limitations and Applications Aspects," *Human-Centric Intelligent Systems*, vol. 3, no. 3, pp. 161–188, Aug. 2023, doi: 10.1007/s44230-023-00038-y.
10. P. A. Moreno-Sánchez, "Improvement of a prediction model for heart failure survival through explainable artificial intelligence," *Front Cardiovasc Med*, vol. 10, Aug. 2023, doi: 10.3389/fcvm.2023.1219586.





11. A. Sethi, S. Dharmavaram, and S. K. Somasundaram, "Explainable Artificial Intelligence (XAI) Approach to Heart Disease Prediction," in *2024 3rd International Conference on Artificial Intelligence For Internet of Things (AIIoT)*, IEEE, May 2024, pp. 1–6. doi: 10.1109/AIIoT58432.2024.10574635.
12. A. Salih *et al.*, "Explainable Artificial Intelligence and Cardiac Imaging: Toward More Interpretable Models," *Circ Cardiovasc Imaging*, vol. 16, no. 4, Apr. 2023, doi: 10.1161/CIRCIMAGING.122.014519.
13. E. Tjoa and C. Guan, "A Survey on Explainable Artificial Intelligence (XAI): Toward Medical XAI," *IEEE Trans Neural Netw Learn Syst*, vol. 32, no. 11, pp. 4793–4813, Nov. 2021, doi: 10.1109/TNNLS.2020.3027314.
14. R. C. Detrano *et al.*, "International application of a new probability algorithm for the diagnosis of coronary artery disease.," *Am J Cardiol*, vol. 64 5, pp. 304–10, 1989, [Online]. Available: <https://api.semanticscholar.org/CorpusID:23545303>
15. I. D. Mienye and N. Jere, "Optimized Ensemble Learning Approach with Explainable AI for Improved Heart Disease Prediction," *Information*, vol. 15, no. 7, p. 394, Jul. 2024, doi: 10.3390/info15070394.

---

*Cite this Article: Yaseen, O.M., Rashid, M.M.(2025). An Explainable Artificial Intelligence (XAI) Methodology for Heart Disease Classification. International Journal of Current Science Research and Review, 8(2), pp. 809-817. DOI: <https://doi.org/10.47191/ijcsrr/V8-i2-28>*