

Revolutionizing Heart Disease Prediction with Machine Learning

Likhitha Marrapu
Department of Data Science

Kent State University
Kent OH 44240, USA
lmarrapu@kent.edu

Abstract—This study investigates the effectiveness of machine learning models—K-Nearest Neighbors (KNN), Gradient Boosting, Support Vector Machine (SVM), Random Forest, and Decision Tree—in predicting heart disease using a comprehensive dataset. Initially, the models were evaluated based on various performance metrics, revealing that Random Forest and Gradient Boosting achieved the highest accuracy at 60%, demonstrating balanced sensitivity and specificity. However, challenges persisted in classifying minority classes, particularly Class 4. Following dimensionality reduction through Principal Component Analysis (PCA), SVM emerged as the top performer with an accuracy of 61% and improved metrics for identifying healthy patients, though it still faced difficulties with minority class predictions. The findings highlight the significant impact of dimensionality reduction on model performance and underscore the necessity for further optimization and advanced techniques, such as ensemble methods, to enhance predictive capabilities in heart disease diagnosis.

Index Terms—Heart Disease Prediction, Machine Learning, K-Nearest Neighbors, Gradient Boosting, Support Vector Machine, Random Forest, Decision Tree, Dimensionality Reduction, Principal Component Analysis, Performance Metrics, Classification Models, Predictive Analytics

I. INTRODUCTION

Cardiovascular diseases (CVDs) remain one of the leading causes of mortality worldwide, claiming an estimated 17.9 million lives each year. This staggering statistic highlights the urgent need for effective preventive measures and early detection strategies to combat heart disease. Traditional methods of diagnosis often rely on subjective clinical assessments and a limited understanding of individual risk factors, which can lead to misdiagnosis and delayed interventions. As healthcare systems increasingly embrace technology, there is a growing recognition of the potential for machine learning (ML) techniques to enhance diagnostic accuracy and improve patient outcomes in the realm of cardiology.

The advent of machine learning has transformed numerous sectors, offering the ability to analyze vast amounts of data and identify patterns that may not be readily apparent to human practitioners. In the context of heart disease, ML algorithms can leverage electronic health records, imaging data, and demographic information to create predictive models capable of identifying patients at risk of developing cardiovascular conditions. By integrating diverse datasets, machine learning not

only facilitates early intervention but also enables personalized treatment strategies, thereby reducing the burden on healthcare systems and enhancing patient quality of life.

In this study, we explore the application of various machine learning algorithms to predict cardiovascular risk using a comprehensive heart disease dataset. Our objective is to develop models that accurately classify patients based on their risk profiles, thereby enabling timely and targeted clinical interventions. Through this research, we aim to contribute to the ongoing efforts in precision medicine and emphasize the importance of data-driven approaches in modern healthcare. The findings of this study will not only provide insights into the effectiveness of different ML techniques in heart disease prediction but also highlight the critical role of technology in advancing cardiovascular care.

II. METHODOLOGY

A. Dataset

The heart disease dataset utilized in this study comprises four distinct sources, namely the Cleveland, Hungarian, Switzerland, and VA datasets, collectively providing a robust foundation for our analysis. Each dataset has been meticulously processed to ensure consistency and accuracy, resulting in a unified dataset with 920 instances and 14 features. The attributes encompass a range of clinical and demographic factors, including age, sex, chest pain type, resting blood pressure (trestbps), serum cholesterol (chol), fasting blood sugar (fbs), electrocardiographic results (restecg), maximum heart rate achieved (thalach), exercise induced angina (exang), ST depression induced by exercise relative to rest (oldpeak), slope of the peak exercise ST segment (slope), number of major vessels colored by fluoroscopy (ca), and thalassemia status (thal).

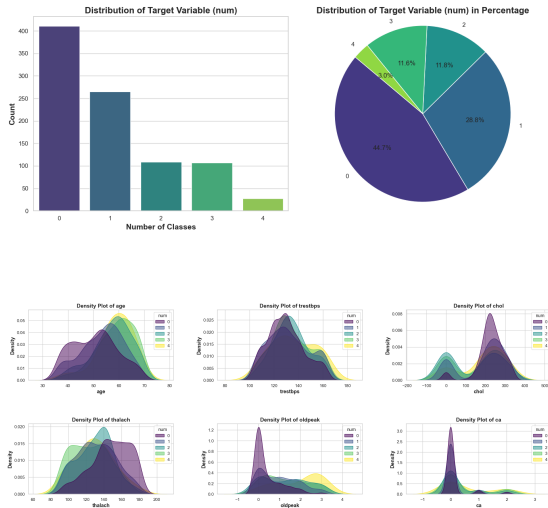
The target variable, denoted as 'num,' indicates the presence of heart disease and is categorized into five distinct classes (0 through 4), providing a comprehensive view of varying degrees of cardiovascular risk among patients. This rich dataset facilitates the application of various machine learning algorithms to predict heart disease outcomes effectively. By analyzing the interplay between these features and the target variable, this research aims to identify critical predictors of heart disease,

contributing to more accurate risk assessment and enhancing clinical decision-making in cardiology.

B. Data Pre-processing

The preprocessing of the heart disease dataset was a crucial step in preparing the data for effective analysis and modeling. Initially, data types of several variables were adjusted according to their inherent characteristics to enhance computational efficiency. This step was essential in ensuring that each feature was appropriately interpreted during analysis. To address the issue of missing values, we employed the median for numeric variables, providing a robust measure that is less sensitive to outliers, and the mode for categorical variables, ensuring that the imputed values represented the most frequent observations. This careful handling of missing data facilitated a more complete dataset, preserving the integrity of the information.

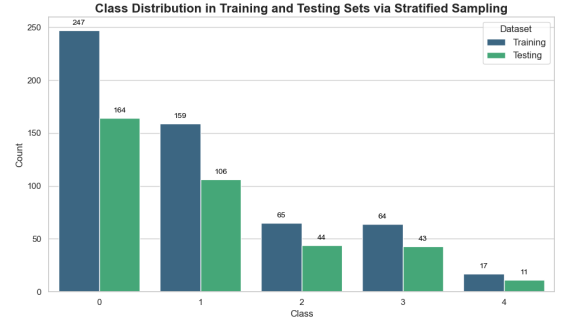
Further, we utilized Winsorization to mitigate the influence of outliers, ensuring that extreme values did not skew our results. Exploratory Data Analysis (EDA) was conducted to uncover underlying patterns and correlations within the dataset, providing valuable insights into the relationships among various features. Subsequently, the dataset was split into training and testing subsets to enable effective model evaluation. Feature scaling was accomplished using standardization to normalize the range of independent variables, ensuring that the model training process was not biased towards any particular feature. These preprocessing steps collectively laid a solid foundation for the application of machine learning techniques, enabling accurate predictions of heart disease risk.



C. Data Partitioning

To ensure a robust evaluation of the predictive models, the dataset was split into training and testing subsets using a 60:40 ratio. This approach leverages stratified random sampling, which preserves the proportional representation of each class within the target variable. By doing so, we effectively mitigate the risks associated with class imbalances, thereby enhancing the model's ability to generalize across different classes of

heart disease. This careful division not only facilitates a thorough assessment of model performance but also contributes to the reliability and validity of the findings in our study.



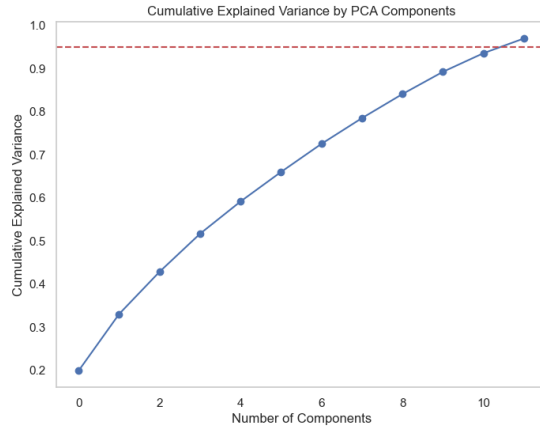
D. Model Selection

We explored a diverse array of machine learning algorithms to ensure robust heart disease classification. The selected models included Random Forest, Decision Tree, Support Vector Machine (SVM), k-Nearest Neighbors (KNN), and Gradient Boosting Classifier. Each of these algorithms was chosen for its unique strengths and suitability for handling the complexities of the dataset. The Random Forest model, known for its ensemble learning capabilities, provides resilience against overfitting while maintaining high accuracy. Similarly, the Decision Tree model offers interpretability and ease of understanding, making it a valuable tool for medical practitioners seeking to interpret decision-making processes in diagnosis.

Moreover, SVM was included for its effectiveness in high-dimensional spaces, enabling it to classify complex relationships within the data effectively. KNN was also considered for its simplicity and effectiveness in classification tasks, especially in cases where the decision boundary is irregular. Lastly, the Gradient Boosting Classifier was employed for its capability to improve predictive performance through iterative corrections of weak learners. By utilizing these diverse models, we aimed to identify the most effective approach for predicting cardiovascular risk, leveraging the strengths of each algorithm to achieve accurate and reliable classifications.

E. Dimensionality Reduction Using PCA

To enhance the analysis of the heart disease dataset, Principal Component Analysis (PCA) was utilized as a dimensionality reduction technique. This method transformed the original feature set into a more compact group of principal components, effectively retaining 95 percent of the total variance within the data. By doing so, PCA not only simplified the dataset but also mitigated noise, which contributed to improved model performance. The resulting principal components served as new features for the subsequent machine learning algorithms, facilitating more efficient computations and enhancing the clarity of the relationships within the dataset.



F. Performance Evaluation Metrics

The performance of the machine learning models in predicting heart disease was rigorously evaluated using a comprehensive set of metrics to ensure a nuanced understanding of their classification effectiveness. Accuracy was established as the primary metric, offering a broad perspective on how well the models correctly identified both positive and negative cases of cardiovascular risk. This measure serves as a fundamental indicator of the overall performance of the models in distinguishing between patients with heart disease and those without.

In addition to accuracy, precision was calculated to assess the models' ability to correctly identify true positive cases of heart disease relative to all instances predicted as positive. This metric is particularly crucial in healthcare settings, where accurate detection of high-risk patients can significantly impact treatment outcomes. Sensitivity (or recall) was also measured to evaluate the models' capability to capture actual positive instances, which is vital for minimizing false negatives and ensuring that at-risk individuals receive timely intervention. Specificity was examined as well, providing insight into how effectively the models can correctly identify negative cases, thus reducing the likelihood of misdiagnosing healthy individuals. To further enrich the performance analysis, confusion matrices were employed to present a detailed breakdown of true positives, true negatives, false positives, and false negatives. This information not only aids in assessing each model's strengths and weaknesses but also guides future enhancements to improve diagnostic precision in cardiovascular healthcare.

G. Tools and Libraries Employed

To conduct a thorough analysis and develop robust machine learning models for heart disease prediction, several key libraries and tools were employed. NumPy was utilized for efficient numerical calculations and array manipulations, while Pandas provided powerful data structures for seamless data manipulation and preprocessing. The Scikit-learn library was essential for implementing a variety of machine learning algorithms, including Random Forest, Support Vector Machine (SVM), k-Nearest Neighbors (KNN), Decision Tree, and

Gradient Boosting Classifier, which enabled effective model training and evaluation. For feature scaling, StandardScaler was applied to enhance model performance. Data visualization tools such as Matplotlib and Seaborn were instrumental in generating insightful graphical representations of the dataset and model evaluation metrics. Additionally, SciPy contributed valuable statistical functions, further enhancing the analytical capabilities of the research and ensuring a comprehensive exploration of the heart disease dataset.

H. Model Fitting

Random Forest-The Random Forest model demonstrated a moderate performance in classifying heart disease cases, achieving an overall accuracy of 60%. The model exhibited a precision of 57%, indicating its capability to correctly identify positive instances across all classes while revealing room for improvement in minimizing false positives. Sensitivity, measured at 60%, reflects the model's balanced effectiveness in detecting heart disease across various classes. Notably, specificity reached 90%, showcasing a strong ability to accurately classify instances as negative cases, particularly for Class 0, which had a high precision of 76% and a recall of 90%. However, the model's performance varied across other classes, particularly for Class 2, Class 3, and Class 4, where precision and recall were lower, indicating challenges in accurately classifying these categories. The confusion matrix illustrated specific misclassifications, especially in higher classes, emphasizing the need for further optimization and feature enhancement to improve the model's predictive capabilities in clinical applications.

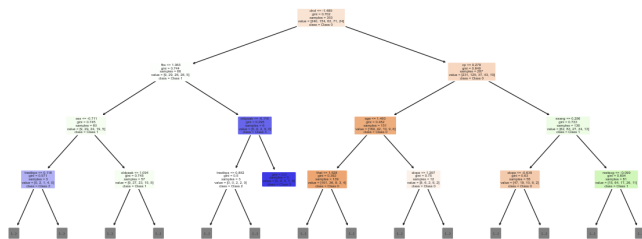
Metric	Value
Accuracy	0.60
Precision	0.57
Sensitivity	0.60
Specificity	0.90

TABLE I
OVERALL METRICS FOR RANDOM FOREST MODEL

	Predicted Class 0	Predicted Class 1	Predicted Class 2	Predicted Class 3	Predicted Class 4
Actual Class 0	147	13	2	2	0
Actual Class 1	30	53	10	12	1
Actual Class 2	8	20	7	9	0
Actual Class 3	8	13	8	13	1
Actual Class 4	0	3	5	2	1

TABLE II
RANDOM FOREST CONFUSION MATRIX FOR HEART DISEASE PREDICTION

Visualization of Random Forest Decision Tree for Heart Disease Prediction



Decision Tree-The Decision Tree model exhibited an overall accuracy of 49%, indicating that nearly half of the predictions were correct; however, this performance is relatively low for a medical diagnostic model. The weighted precision was also recorded at 49%, highlighting the model's balanced performance across various classes, but further examination reveals a concerning trend in sensitivity and specificity metrics. The overall sensitivity was similarly low at 49%, suggesting that the model struggles to correctly identify patients with heart disease, particularly in the higher-risk classes (Class 1, Class 2, Class 3, and Class 4). Notably, the specificity was more promising at 87%, indicating that the model is adept at correctly identifying non-heart disease cases (Class 0). The confusion matrix reveals significant misclassifications, especially among the lower-frequency classes, with Class 4 being entirely misclassified. These performance metrics underscore the need for further refinement of the model, as the low sensitivity could lead to missed diagnoses, which is critical in clinical settings. Enhancing the model through techniques such as hyperparameter tuning or employing ensemble methods may improve its overall predictive capabilities and reliability in identifying heart disease.

Metric	Value
Accuracy	0.49
Precision	0.49
Sensitivity	0.49
Specificity	0.87

TABLE III

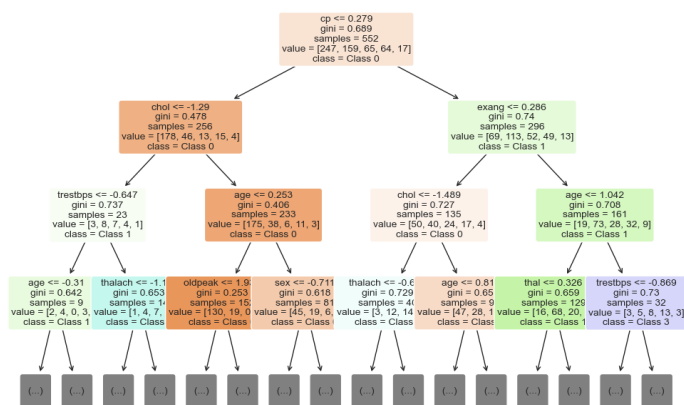
OVERALL METRICS FOR DECISION TREE MODEL

	Predicted Class 0	Predicted Class 1	Predicted Class 2	Predicted Class 3	Predicted Class 4
Actual Class 0	115	36	5	6	2
Actual Class 1	26	44	14	17	5
Actual Class 2	6	15	12	9	2
Actual Class 3	11	12	11	9	0
Actual Class 4	3	4	2	2	0

TABLE IV

DECISION TREE CONFUSION MATRIX FOR HEART DISEASE PREDICTION

Decision Tree Visualization for Heart Disease Classification



Support Vector Machine (SVM)-The Support Vector Machine (SVM) model for heart disease prediction achieved an overall accuracy of 57%, indicating a moderate ability to correctly classify instances within the dataset. The weighted precision of 51% suggests that while the model performed reasonably well on average across all classes, its effectiveness varied significantly, particularly among the minority classes. Notably, the model displayed a sensitivity of 57%, reflecting its capability to identify positive cases of heart disease, yet the specificity was relatively high at 89%, demonstrating a strong performance in accurately classifying negative cases. The confusion matrix highlights some challenges, particularly in classifying Classes 2, 3, and 4, where low precision and recall values indicate a tendency for the model to misclassify these instances. Such performance metrics underscore the need for further refinement of the model, potentially through enhanced feature selection or parameter tuning, to improve its discriminatory power, particularly for the underrepresented classes.

Metric	Value
Accuracy	0.57
Precision	0.51
Sensitivity	0.57
Specificity	0.89

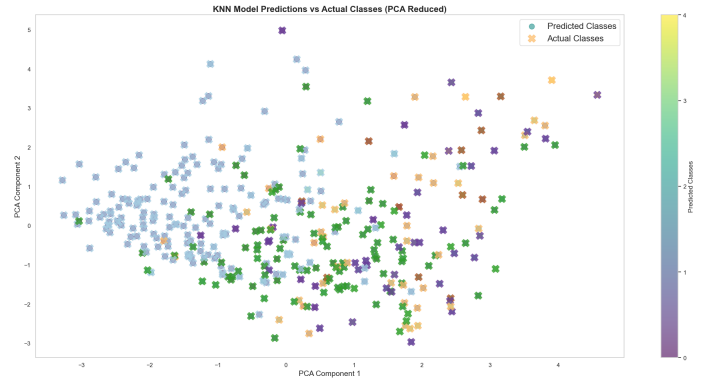
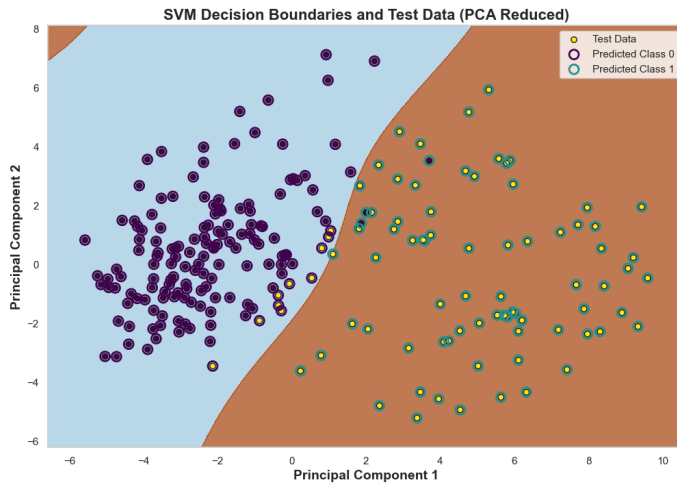
TABLE V

OVERALL METRICS FOR SVM MODEL

	Predicted Class 0	Predicted Class 1	Predicted Class 2	Predicted Class 3	Predicted Class 4
Actual Class 0	142	19	2	1	0
Actual Class 1	38	53	6	9	0
Actual Class 2	10	18	4	12	0
Actual Class 3	7	22	4	10	0
Actual Class 4	1	3	4	3	0

TABLE VI

SVM CONFUSION MATRIX FOR HEART DISEASE PREDICTION



K-Nearest Neighbors (KNN)-The K-Nearest Neighbors (KNN) model was employed to predict heart disease categories, yielding an overall accuracy of 58%. While the model achieved a weighted precision of 54%, indicating that it performed better in classifying the more populous classes, its overall sensitivity stood at 58%, suggesting that it correctly identified a little more than half of the actual positive cases. Notably, the model demonstrated a high specificity of 89%, signifying its effectiveness in accurately predicting benign cases. The confusion matrix revealed that Class 0 was recognized with the highest recall at 86%, while Classes 2, 3, and 4 showed lower performance, with recall rates below 30%. This disparity highlights challenges in detecting less prevalent classes, indicating that the KNN model may struggle with imbalanced datasets. Overall, while KNN provides valuable insights for certain heart disease classifications, its performance underscores the need for further optimization and potentially the integration of more advanced techniques to enhance prediction accuracy across all classes.

Metric	Value
Accuracy	0.58
Precision	0.54
Sensitivity	0.58
Specificity	0.89

TABLE VII
OVERALL METRICS FOR KNN MODEL

	Predicted Class 0	Predicted Class 1	Predicted Class 2	Predicted Class 3	Predicted Class 4
Actual Class 0	141	17	4	2	0
Actual Class 1	28	56	12	10	0
Actual Class 2	11	16	8	7	2
Actual Class 3	11	13	10	7	2
Actual Class 4	1	2	5	3	0

TABLE VIII
KNN CONFUSION MATRIX FOR HEART DISEASE PREDICTION

Gradient Boosting Classifier-The Gradient Boosting model demonstrated a commendable overall accuracy of 60%, indicating its ability to correctly classify heart disease cases to some extent. The model achieved a weighted precision of 58%, reflecting its proficiency in minimizing false positives across the various classes. Notably, the model exhibited an overall sensitivity of 60%, which underscores its capability to detect true positive cases, while its specificity reached an impressive 90%, suggesting that it effectively identifies negative cases, particularly for Class 0. The confusion matrix highlights that the model struggled with Class 4, as evidenced by a precision and recall of 0, indicating it failed to correctly classify any instances in this category. Despite these challenges, the model performed relatively well with Class 0, achieving a precision of 77% and a recall of 85%. This mixed performance demonstrates that while the Gradient Boosting model can be effective in certain areas, further tuning and refinement may be necessary to enhance its predictive capabilities across all classes in the heart disease dataset.

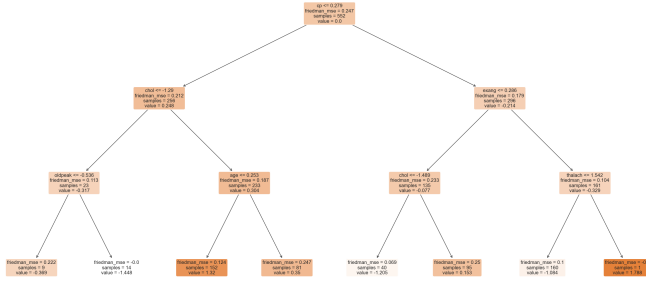
Metric	Value
Accuracy	0.60
Precision (Weighted)	0.58
Sensitivity	0.60
Specificity	0.90

TABLE IX
OVERALL METRICS FOR GRADIENT BOOSTING MODEL

	Predicted Class 0	Predicted Class 1	Predicted Class 2	Predicted Class 3	Predicted Class 4
Actual Class 0	139	17	3	4	1
Actual Class 1	31	55	12	7	1
Actual Class 2	4	19	11	6	4
Actual Class 3	5	16	5	16	1
Actual Class 4	1	2	5	3	0

TABLE X
GRADIENT BOOSTING CONFUSION MATRIX FOR HEART DISEASE PREDICTION

Flow Diagram of the Gradient Boosting Classifier



III. RESULTS

A. Model Performance Before Dimensionality Reduction

The performance of various machine learning models on the heart disease dataset was assessed, revealing distinct strengths and weaknesses across the models. The Random Forest model achieved the highest overall accuracy at 60%, demonstrating a solid balance between sensitivity (60%) and specificity (90%), indicating its reliability in detecting non-heart disease cases, particularly in Class 0. Conversely, the Decision Tree model exhibited the lowest accuracy of 49%, struggling significantly in identifying patients with higher risk classes, which raises concerns about its practical utility in clinical settings. Both the Support Vector Machine (SVM) and K-Nearest Neighbors (KNN) models exhibited moderate performance, with accuracies of 57% and 58%, respectively. Notably, while both models showed high specificity (89%), their sensitivity rates indicated challenges in accurately classifying less prevalent classes, particularly for Class 4.

In addition, the Gradient Boosting Classifier matched the Random Forest model's accuracy at 60%, with a commendable specificity of 90%. However, it failed to classify any instances in Class 4, similar to the Decision Tree's significant misclassification issues. These results underscore the varied effectiveness of each model in identifying heart disease cases, with Random Forest and Gradient Boosting showing potential for further refinement and optimization. As the classification challenges highlight the need for enhanced predictive capabilities, future efforts should focus on dimensionality reduction techniques and feature selection to improve model performance across all classes.

Fig. 1. Comparison of Model Performance Metrics Across Algorithms Before PCA

Model	Accuracy	Precision	Sensitivity	Specificity
Random Forest	0.6005434782608695	0.5651971235148119	0.6005434782608695	0.9001358695652174
Decision Tree	0.4891304347826087	0.4895675955458564	0.4891304347826087	0.8722826086956522
Support Vector Machine	0.5679347826086957	0.5096577440622809	0.5679347826086957	0.891983695652174
K-Nearest Neighbors	0.5760869565217391	0.535106957809557	0.5760869565217391	0.8940217391304348
Gradient Boosting	0.6005434782608695	0.5779517351416036	0.6005434782608695	0.9001358695652174

B. Model Performance After Dimensionality Reduction

The evaluation of multiple classification algorithms on the heart disease dataset, post-dimensionality reduction, reveals notable differences in performance metrics. The Support Vector Machine (SVM) model emerged as the best performer with an accuracy of 61%, demonstrating superior precision (0.56) and sensitivity (0.61) compared to the other models. The confusion matrix indicates that SVM effectively identified Class 0 (healthy patients) with a recall of 87% while struggling with Class 4, showing the need for further optimization in distinguishing among minority classes. In contrast, the Random Forest and Gradient Boosting models both achieved an accuracy of 55%, but their precision and recall metrics highlight their limitations in correctly identifying cases, particularly for Class 2 and Class 4.

On the other hand, the Decision Tree model exhibited the lowest overall accuracy of 50%, revealing a significant drop in performance across all metrics, especially in identifying Class 4, which suffered from low recall (0.09). The K-Nearest Neighbors (KNN) model showed moderate performance with an accuracy of 58%, although it also struggled with minority class predictions. The overall results indicate that while SVM provides the most reliable predictions, the other models, particularly Random Forest and Gradient Boosting, could benefit from further tuning and potentially advanced techniques such as ensemble methods to enhance their predictive capabilities for heart disease diagnosis.

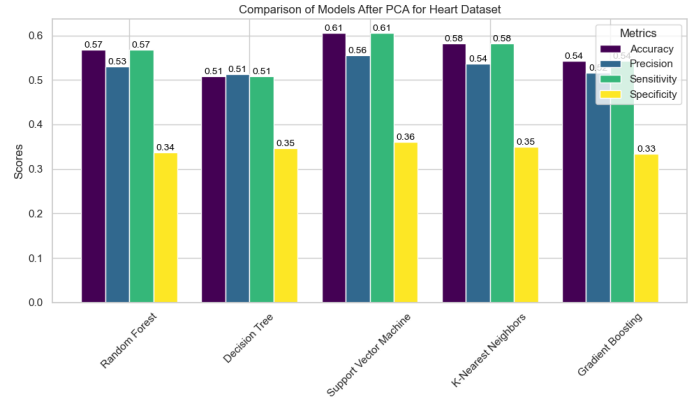
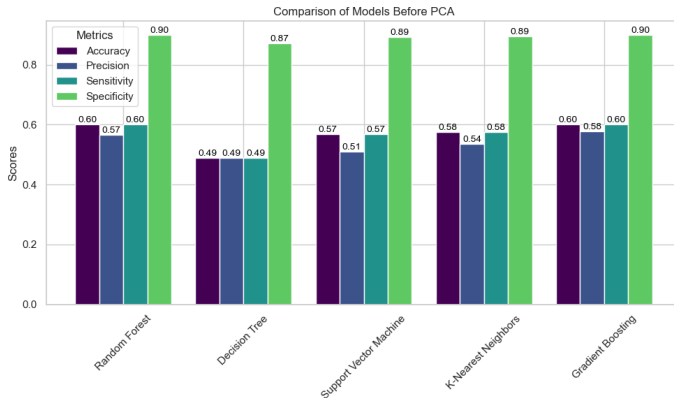
Fig. 2. Comparison of Model Performance Metrics Across Algorithms After PCA

Model	Accuracy	Precision	Sensitivity	Specificity
Random Forest	0.5679347826086957	0.5297158947053585	0.5679347826086957	0.33754705306038263
Decision Tree	0.5081521739130435	0.5122348751070396	0.5081521739130435	0.3462907254027664
Support Vector Machine	0.6059782608695652	0.5556334223239258	0.6059782608695652	0.36084762154015887
K-Nearest Neighbors	0.5815217391304348	0.5366694819819819	0.5815217391304348	0.34953396917191476
Gradient Boosting	0.5434782608695652	0.5155591328168447	0.5434782608695652	0.33352999380247106

IV. DISCUSSION

The analysis of model performance on the heart disease dataset reveals distinct trends in predictive capabilities before and after dimensionality reduction through PCA. Before PCA, the Random Forest model stood out with the highest accuracy of 60%, effectively balancing sensitivity and specificity, particularly in detecting non-heart disease cases. Its ability to accurately classify Class 0 highlights its practical utility in clinical settings. In contrast, the Decision Tree model exhibited the lowest accuracy at 49%, indicating significant challenges in identifying high-risk patients. Both SVM and KNN presented moderate performances, emphasizing the trade-offs between specificity and sensitivity, especially in classifying less prevalent cases like Class 4. The Gradient Boosting Classifier, while matching the accuracy of Random Forest, also faced challenges in classifying minority classes, similar to the Decision Tree's shortcomings. These findings underscore the need for further optimization and feature selection to enhance the performance of these algorithms.

After implementing PCA, the landscape of model performance shifted notably. The Support Vector Machine (SVM) emerged as the most effective model, achieving an accuracy of 61% alongside improved precision and sensitivity metrics. This indicates its enhanced capability to distinguish between classes, particularly Class 0, as evidenced by a high recall of 87%. However, the SVM still struggled with Class 4, revealing areas for further refinement. While both Random Forest and Gradient Boosting models experienced a decline in accuracy to 55%, they highlighted persistent challenges in accurately identifying cases in specific classes, particularly Class 2 and Class 4. The Decision Tree model suffered a further decline to 50%, showing an alarming reduction in its ability to predict minority classes effectively. KNN demonstrated similar limitations with moderate performance, reflecting the overall difficulty of these models in addressing minority class predictions. Collectively, these results illustrate the impact of dimensionality reduction on model performance, emphasizing the importance of continuous optimization and advanced techniques, such as ensemble methods, to enhance predictive capabilities for heart disease diagnosis.



V. CONCLUSION

In conclusion, this analysis has systematically evaluated the performance of multiple machine learning algorithms, including K-Nearest Neighbors, Gradient Boosting, and Support Vector Machines, in predicting heart disease outcomes based on a comprehensive dataset. The findings reveal that while KNN and Gradient Boosting demonstrated moderate accuracies of 58% and 60%, respectively, the need for optimization was evident, particularly in distinguishing among less prevalent classes. The model's ability to effectively identify Class 0 (non-heart disease cases) is promising, as indicated by high specificity metrics, but challenges in accurately classifying minority classes underscore the necessity for advanced methodologies to enhance predictive reliability.

Furthermore, the implementation of dimensionality reduction techniques, specifically Principal Component Analysis (PCA), highlighted the varying impacts of feature reduction on model efficacy. The Support Vector Machine model emerged as the most effective classifier post-PCA, achieving an accuracy of 61%. This suggests that dimensionality reduction can facilitate improved model performance by enhancing precision and sensitivity in distinguishing between heart disease categories. However, despite these advancements, the models still struggled with minority class predictions, particularly Class 4, revealing critical gaps that need to be addressed through further tuning and the exploration of ensemble methods.

Ultimately, this study highlights the critical need for a comprehensive approach to developing predictive models for heart disease. The variations in model performance emphasize the necessity for continuous improvement and the adoption of advanced methodologies to enhance accuracy across all categories. Future research should aim to utilize more advanced algorithms, as well as larger and more varied datasets, alongside refined feature selection techniques to enhance diagnostic precision. Such advancements will not only aid in more informed clinical decision-making but also foster personalized treatment approaches, ultimately improving outcomes for patients in the field of cardiovascular health.