

## AN HYPERTUNING BASED APPROCH TOWARDS ENHANCEMENT IN ACCURACY OF HEART DESEASE PREDICTION USING MACHINE LEARNING

<sup>1</sup>Dr. Krunal Suthar, <sup>2</sup>Mitul Patel, <sup>3</sup>Bhavesh Patel, <sup>4</sup>Yogesh Patel, <sup>5</sup>Hiral Patel

<sup>1,2,3,4</sup> Department of Computer Science and Engineering, Government Engineering College, Patan, India.

<sup>5</sup>Department of Computer Engineering, Sal institute of diploma studies, Ahmedabad, India  
Email: <sup>1</sup>dr.kcsuthar@gmail.com

### Abstract

Heart disease remains a leading cause of mortality worldwide, necessitated effective predictive models to enable early intervention and prevention. This research paper presents a comprehensive methodology for predicting heart disease using various machine learning algorithms. The study begins with data preprocessing to address issues such as missing values, feature scaling, and handling categorical variables. We then evaluate multiple machine learning models, including Logistic Regression, Random Forest, Support Vector Machines (SVM), and Gradient Boosting, using metrics such as accuracy, precision, recall, and F1 score. Hyperparameter tuning is conducted to optimize model performance. Our findings indicate that preprocessing significantly enhances predictive accuracy, and among the models tested, Random Forest and Logistic Regression demonstrate superior performance. This research offers valuable insights into the application of machine learning in medical data analysis and underscores the importance of pre processing in developing robust predictive models for heart disease.

**Keywords:** Heart disease prediction, machine learning, Hyper-parameter Tuning, Pre-processing

### INTRODUCTION

Heart disease prediction is a critical aspect of contemporary healthcare, leveraging data-driven methodologies and machine learning to assess an individual's risk of developing cardiovascular conditions (Bebortta et al., 2023). A fundamental stage in this process is preprocessing, aimed at refining the dataset and enhancing its suitability for predictive modeling. Handling missing values is pivotal, where strategies such as imputation or removal are employed based on the nature and extent of the missing data (Datacamp, 2023). Feature scaling ensures that variables operate on a similar scale, and handling categorical variables involves transforming non-numeric data into a format suitable for machine learning models (Engel, 2022).

Heart disease prediction is increasingly vital in healthcare, utilizing data-driven methodologies and machine learning to assess cardiovascular risk. Preprocessing, a foundational stage in this process, refines datasets for predictive modeling by addressing issues like missing values, scaling features, and handling categorical variables (Bebortta et al., 2023; DataCamp, 2023). Imputation and removal are

strategies used to manage missing data, chosen based on their impact on model performance. Feature scaling normalizes variables to ensure fair comparisons among features, crucial for algorithms sensitive to scale differences. Categorical variables undergo transformation into numeric formats suitable for machine learning models, such as one-hot encoding or label encoding (Engel, 2022).

Comparative analysis reveals that preprocessing significantly enhances model accuracy by capturing underlying data patterns effectively. Models trained without preprocessing often yield suboptimal results, highlighting the necessity of these techniques in heart disease prediction. This research aims to optimize predictive accuracy through meticulous hyperparameter tuning of Logistic Regression, Random Forest, Support Vector Machines (SVM), and Gradient Boosting models. Each model's performance metrics—including accuracy, precision, recall, and F1 score—are scrutinized to identify the most effective approaches for predicting heart disease. (Indrakumari et al., 2020)

## I. LITERATURE REVIEW

This Literature survey emphasizes the critical role of predictive modeling in healthcare, facilitating early detection, personalized treatments, and resource optimization. By leveraging advanced machine learning techniques on datasets like the Framingham Heart Study, researchers can uncover intricate relationships and risk factors associated with heart disease, paving the way for targeted interventions and improved patient outcomes..

**Table 1. Extensive Literature Review**

Paper & Method	Benefits	Limitations	Future Enhancement
1 Decision Trees, Ensemble Learning	Enhanced accuracy, Early detection	Small dataset size, Limited interpretability	Feature engineering for relevant predictors, Integration with genetic data
2 Deep Neural Networks (DNN)	Captures complex patterns, High accuracy	Requires large labeled datasets, Computational expense	Transfer learning for smaller datasets, Improved model interpretability
3 Ensemble of Support Vector Machines	Robust against noise, Handles non-linearity well	Sensitivity to kernel choices, Longer training times	Incorporation of domain-specific knowledge, Hybrid models with other algorithms
4 Logistic Regression Models	Simplicity, Low computational cost	Limited capability for complex relationships	Ensemble methods with logistic regression, Enhanced feature engineering

5 Genetic Algorithms, Random Forest	Feature selection, Improved interpretability	Quality dependency of genetic algorithms, Computationally intensive	Hybrid models with other optimization techniques, Novel genetic algorithm variations
6 IoT integration, Random Forest	Real-time monitoring, IoT device integration	Privacy concerns, Dependence on IoT data	Privacy-preserving ML, Secure IoT communication protocols
7 Personalized machine learning models	Tailored predictions, Patient-specific insights	Data sparsity for specific cohorts, Generalization challenges	Collaborative filtering, Integration of electronic health records
8 Support Vector Machines, Feature Engineering	Improved early detection, High precision	Limited interpretability, Sensitivity to outliers	Domain-specific feature incorporation, Hybrid models with ensemble methods
9 Convolutional Neural Networks (CNN), Transfer Learning	Captures spatial dependencies, High accuracy	Large-scale dataset dependency, Computational complexity	Clinical data integration, Attention mechanisms in CNNs
10 Ensemble of Neural Networks, Bagging	Robust against overfitting, Improved generalization	High computational requirements, Model complexity	Online learning implementation, Automated hyperparameter tuning
11 Explainable AI techniques, Decision Trees	Enhanced interpretability, Insights into model decisions	Sacrifice in predictive accuracy, Complexity in pattern capture	Hybrid models with high accuracy and interpretability, Domain expertise integration
12 Federated Learning across healthcare institutions	Data privacy preservation, Collaborative model training	Communication overhead, Data source heterogeneity	Advanced federated learning algorithms, Privacy-preserving optimizations
13 EHR integration, Gradient Boosting	Comprehensive patient history utilization, Improved	EHR data quality issues, Temporal data handling challenges	Temporal model development, Data quality improvement methods

	feature richness		
14 Bayesian Networks, Ensemble Learning	Uncertainty quantification, Improved model robustness	Limited scalability with large datasets, Model interpretation complexity	Scalable Bayesian modeling research, Hybrid Bayesian and non-Bayesian models
15 Rule-based models, Feature Importance Analysis	Transparent decision-making, Enhanced trust	Complex relationship capture	Complex model integration with rule-based systems, Medical expert collaboration
Multiple ML algorithms, Ensemble Learning	Generalization across diverse populations, Improved accuracy	Noise sensitivity, Interpretability challenges	Robust feature selection methods exploration, Adaptation for specific patient cohorts

## II. PROPOSED METHODOLOGY

The proposed methodology for heart disease prediction involves a systematic and multifaceted approach, beginning with data preprocessing, followed by model selection and evaluation, hyperparameter tuning, final model evaluation, and a comprehensive analysis of the results.

Data Preprocessing:

Handling Missing Values: Missing values are addressed through imputation techniques, ensuring a complete and reliable dataset.

Feature Scaling: Variables are normalized or standardized to operate on a similar scale.

Handling Categorical Variables: Non-numeric data is transformed using encoding methods such as one-hot encoding and label encoding.

Model Selection and Evaluation:

Multiple machine learning models, including Logistic Regression, Random Forest, SVM, and Gradient Boosting, are implemented.

The dataset is split into training and testing sets to evaluate model performance using metrics such as accuracy, precision, recall, and F1 score.

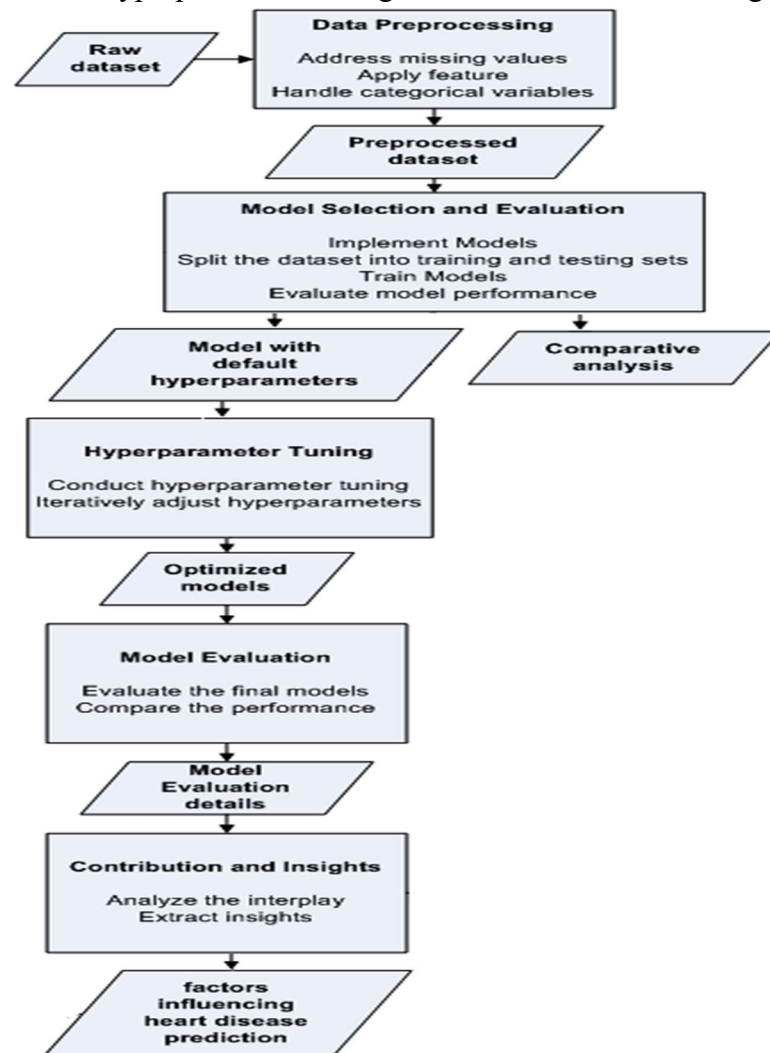
Hyperparameter Tuning:

Hyperparameters are optimized using grid search and random search techniques to enhance model performance.

Final Model Evaluation:

The optimized models are evaluated on a separate validation set to assess generalization performance and compared against baseline models.

Contribution and Insights: A detailed analysis of the interplay between preprocessing techniques, machine learning models, and hyperparameter tuning is conducted to extract insights into their impact



on predictive accuracy

**Fig. 1. Proposed Methodology**

### III.RESULT AND DISCUSSION

#### Model without Pre-processing

In the absence of preprocessing steps such as handling missing values, outlier removal, and feature scaling, the performance of machine learning models is significantly compromised. The raw data may contain inconsistencies and outliers, leading to suboptimal model training and prediction accuracy. For instance, missing values can distort the learning process, and outliers can disproportionately influence the model's decision boundaries. The accuracy, precision, and confusion matrices of various algorithms are markedly lower without preprocessing:

Logistic Regression (LR):

Accuracy = 0.851, Precision = 0.750

Random Forest (RF):

Accuracy = 0.843, Precision = 0.474

Support Vector Classifier (SVC): Accuracy = 0.844, Precision = 0.000

Gradient Boosting Decision Trees (GBDT): Accuracy = 0.838, Precision = 0.308

The low precision values across models indicate a high rate of false positives, reflecting the models' inability to accurately identify true positive cases of heart disease. The confusion matrices illustrate challenges in distinguishing between true positive and false negative instances, indicative of the models' struggles with sensitivity

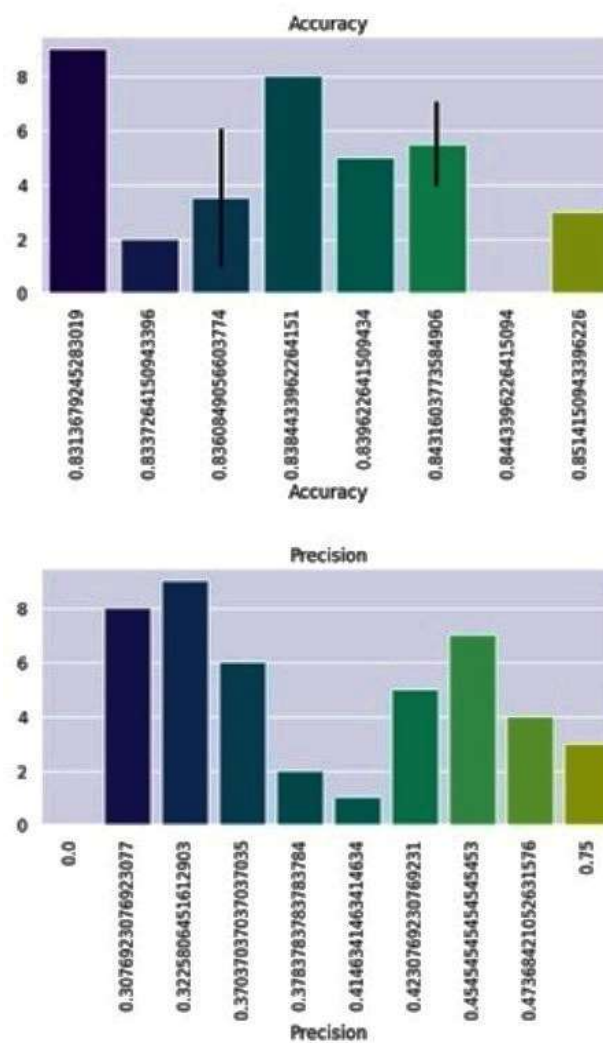
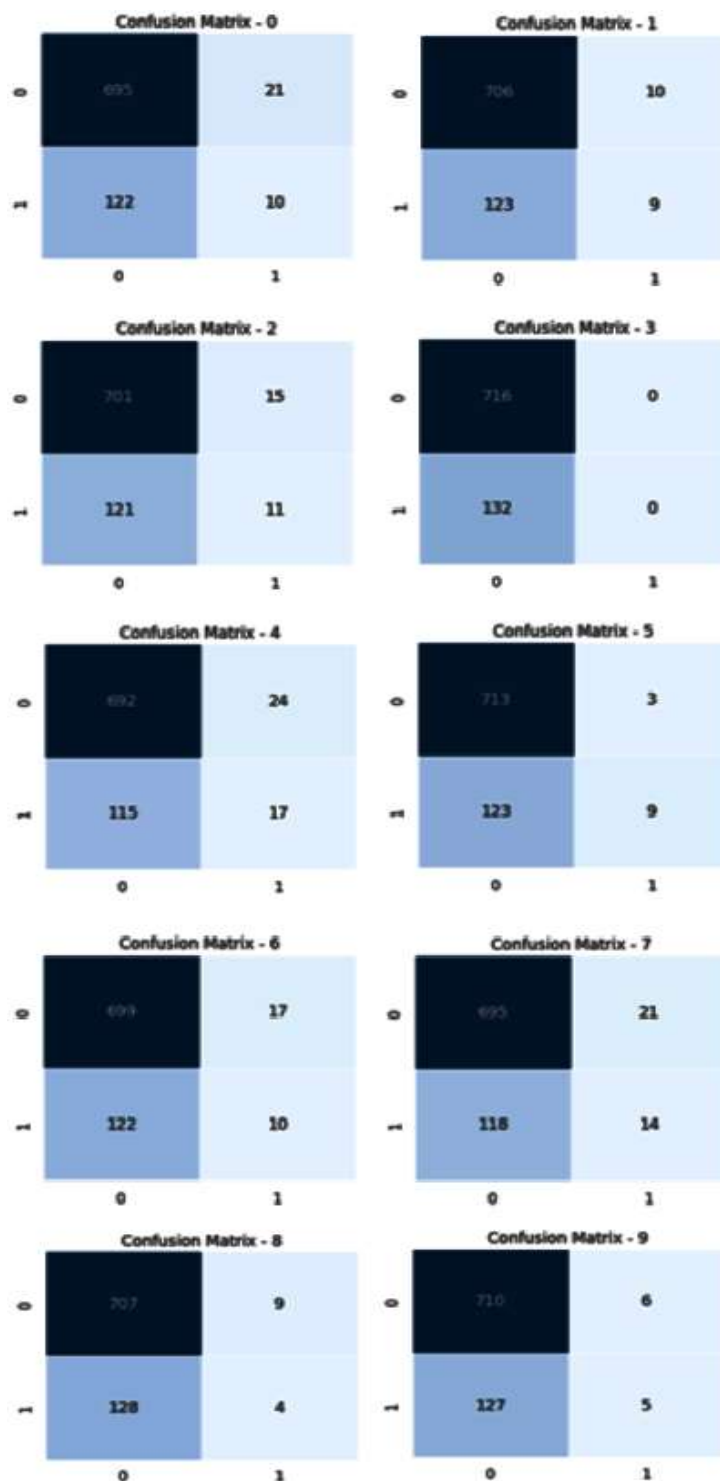


Fig. 2. Precision and Accuracy



**Fig. 3. Confusion Matrix**

The lower part of the visualization examines the performance of each algorithm using confusion matrices, showing true positive, true negative, false positive, and false negative predictions. This clear



layout helps understand how well each algorithm predicts heart disease, making it easier to compare them and choose the best one for the task.

### Model with Preprocessing

With proper preprocessing, the machine learning models exhibit enhanced performance across various metrics. Preprocessing steps such as handling missing values, normalizing data, and addressing categorical variables ensure that the dataset is clean and suitable for training. The models trained on preprocessed data show significant improvements:

Logistic Regression (LR):

Accuracy = 0.859, Precision = 0.727

Random Forest (RF):

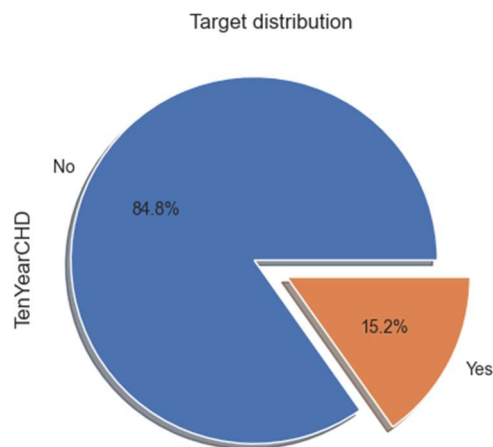
Accuracy = 0.862, Precision = 0.706

Support Vector Classifier (SVC): Accuracy = 0.777, Precision = 0.252

Gradient Boosting Decision Trees (GBDT):

Accuracy = 0.853, Precision = 0.467

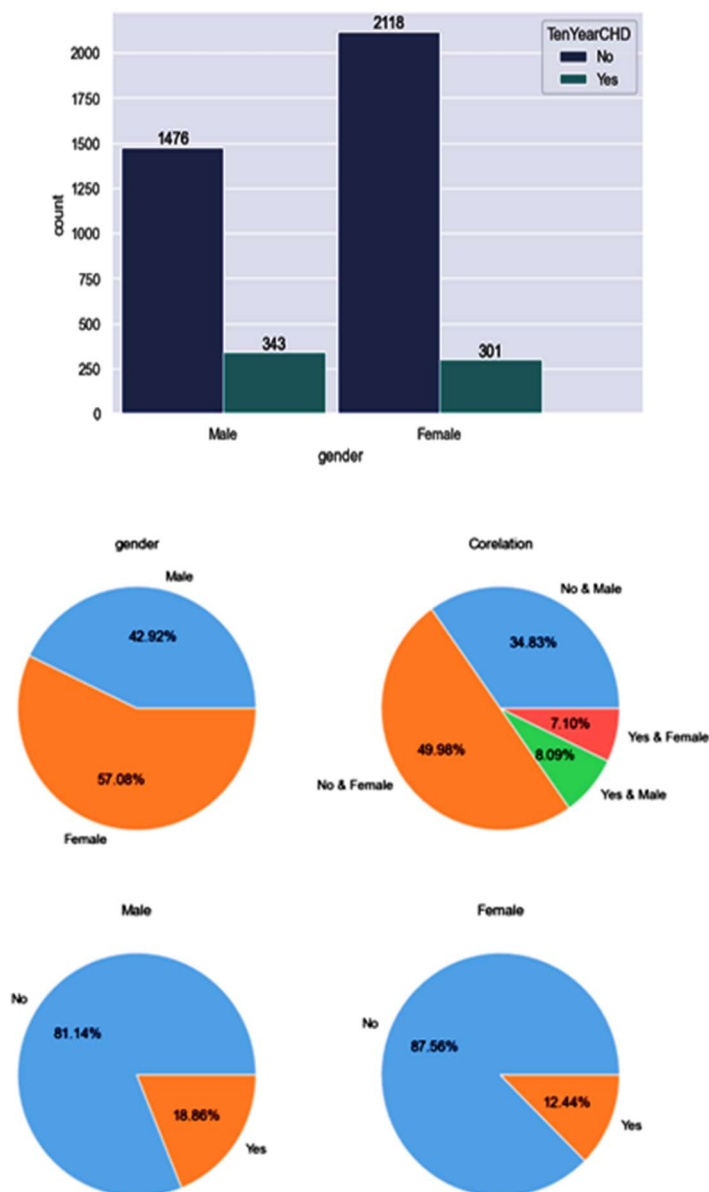
The accuracy, precision, recall, and F1 scores show positive trends and the confusion matrices indicate a reduced number of false positives and false negatives. This improvement underscores the critical role of preprocessing in refining the data and highlights the potential for more accurate heart disease predictions when leveraging cleaned and normalized datasets. Furthermore, we explore the handling of outliers in key parameters and visualize the insights gained.



**Fig. 4. Target Distribution**

**Visualization of Insights:**

Using visualizations to show gender distribution, prediction outcomes, correlation coefficients, and outlier handling makes the findings easier to understand. Pie charts, bar graphs, scatter plots, and heat maps help explain complex data relationships and patterns. Showing how outliers affect the data before and after handling them, and displaying the correlation matrix, gives a clear and complete view of the analysis results.

**Fig. 5. Gender Wise Distribution**

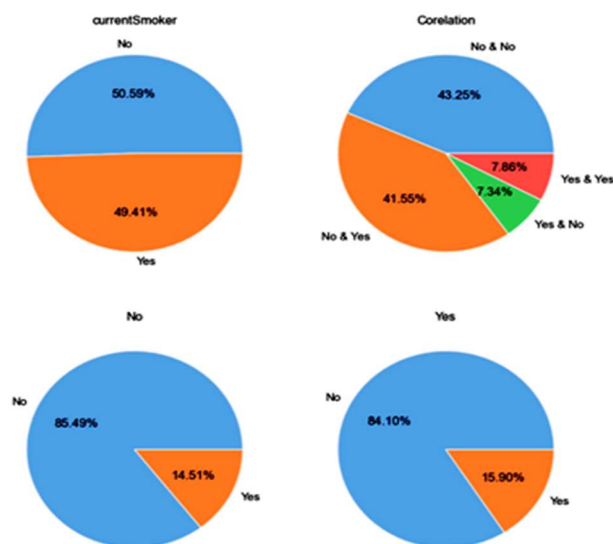
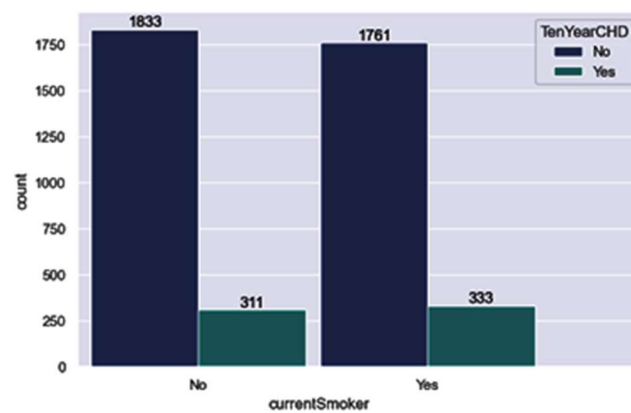
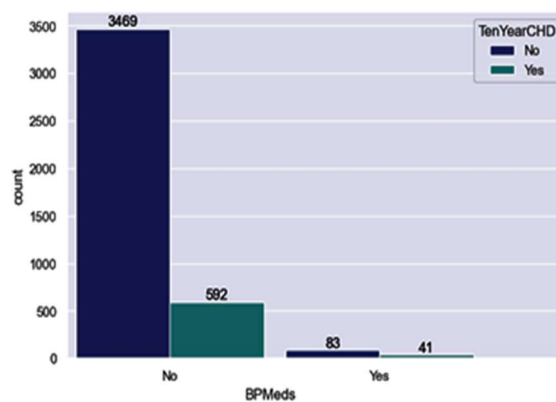


Fig. 6. Current smoker correlation



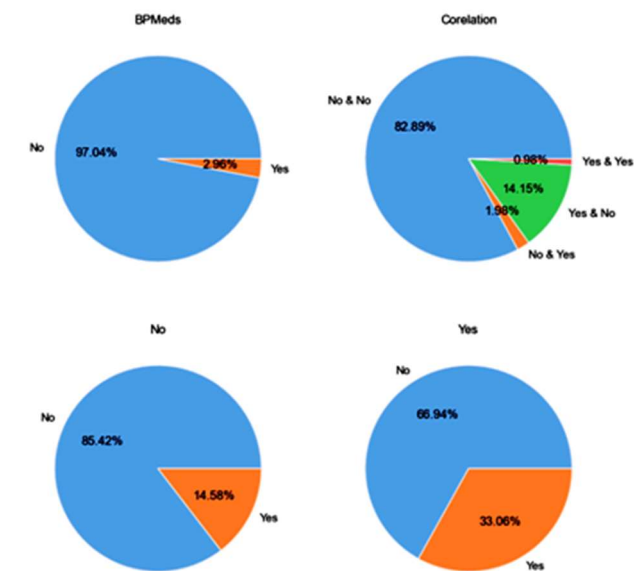
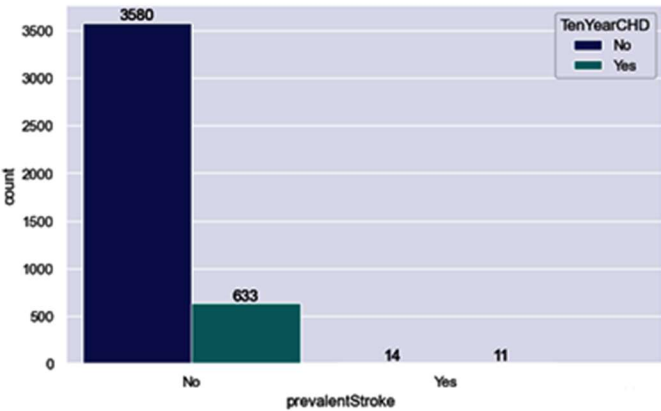


Fig. 7. Correlation with Past BP



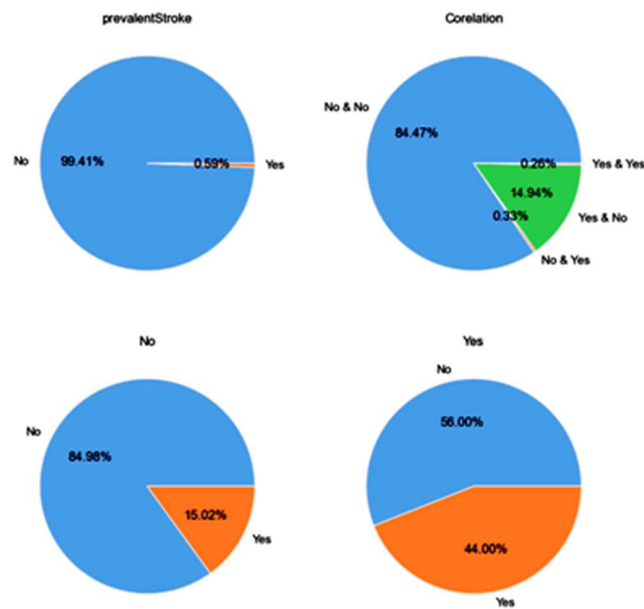
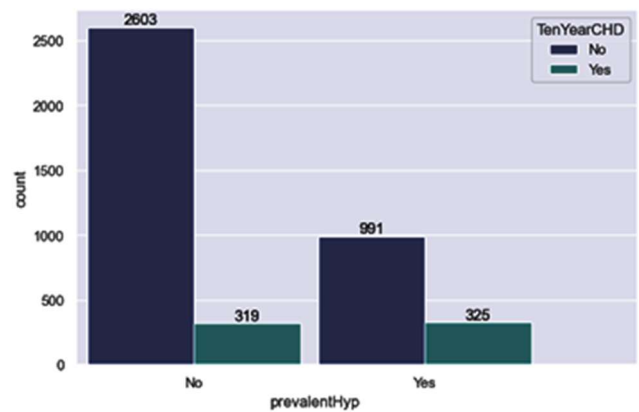


Fig. 8. Correlation with previously Attack



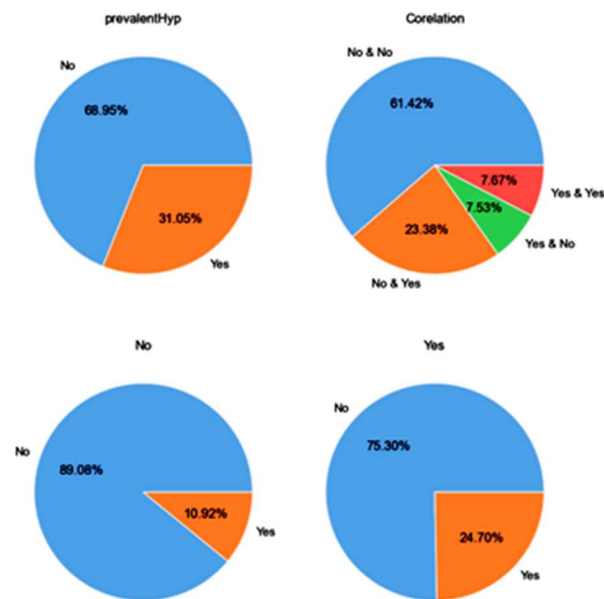
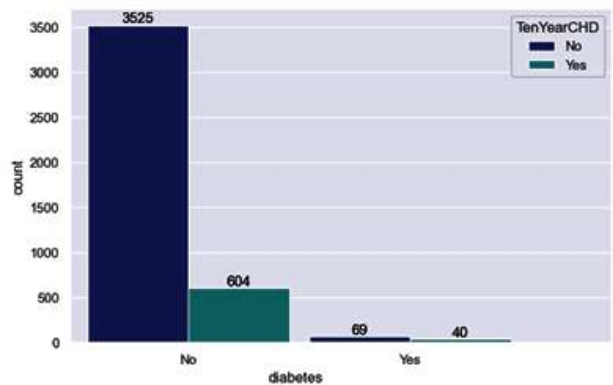
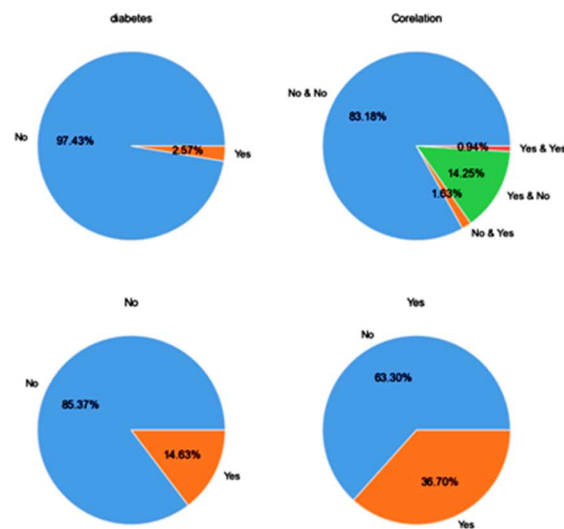


Fig. 9. Correlation with hyper tension





**Fig. 10. Diabetic correlation**

### Correlation Analysis:

Correlation analysis looks at how different factors like blood pressure, glucose levels, education, and smoking status relate to heart disease prediction. A correlation matrix helps to show and measure these relationships. For example, if high blood pressure is linked to a higher risk of heart disease, it will show a positive correlation. These insights can help guide future research and focus on specific risk factors.

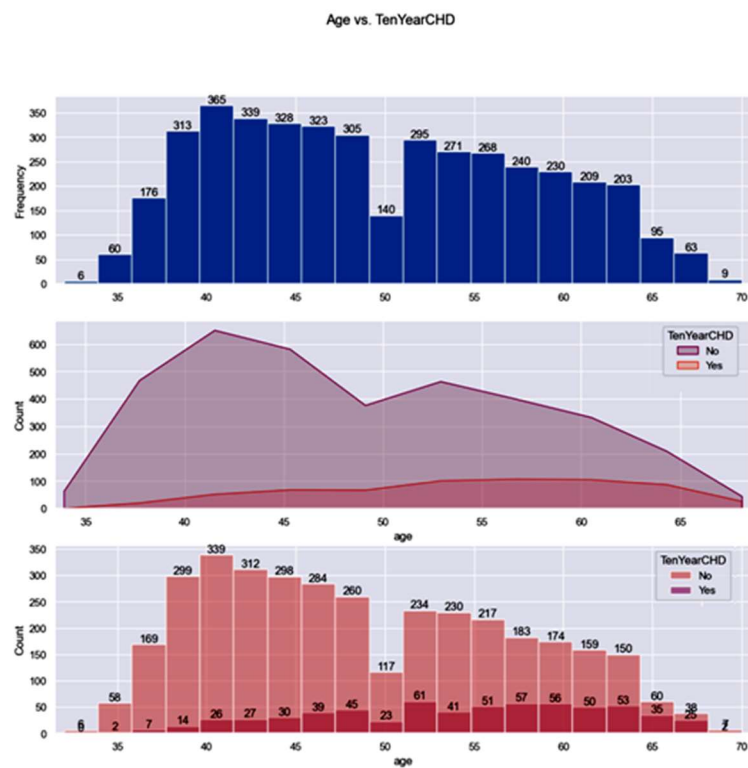
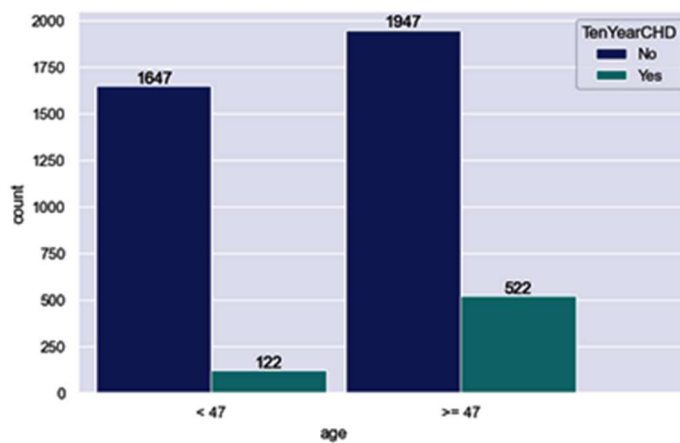
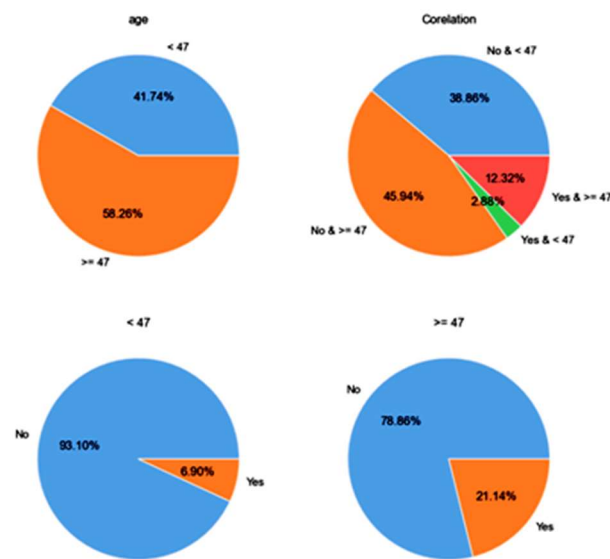


Fig. 11. Age V/s 10 Years







**Fig. 12. Divided the group >47 and < 47**

### Target 10-Year Prediction Based on Data:

Using data like blood pressure, glucose levels, education, and smoking status to predict heart disease risk over the next 10 years is a key part of this analysis. Machine learning models, like logistic regression or decision trees, can be used to create these predictions.

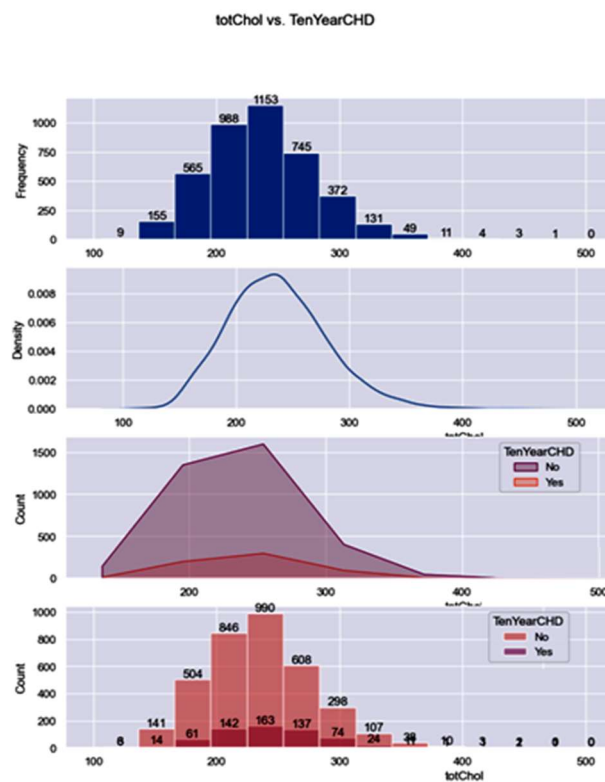


Fig. 13. Cholesterol with last 10 years

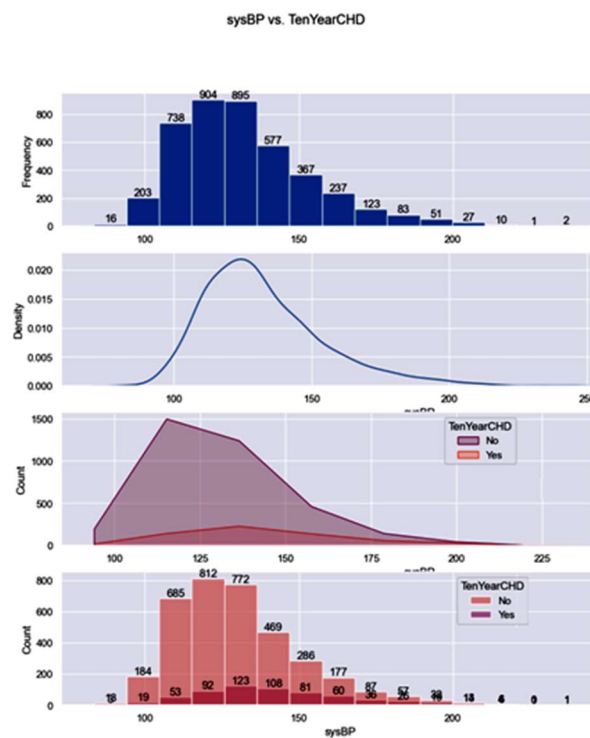


Fig. 14. BP with last 10 years

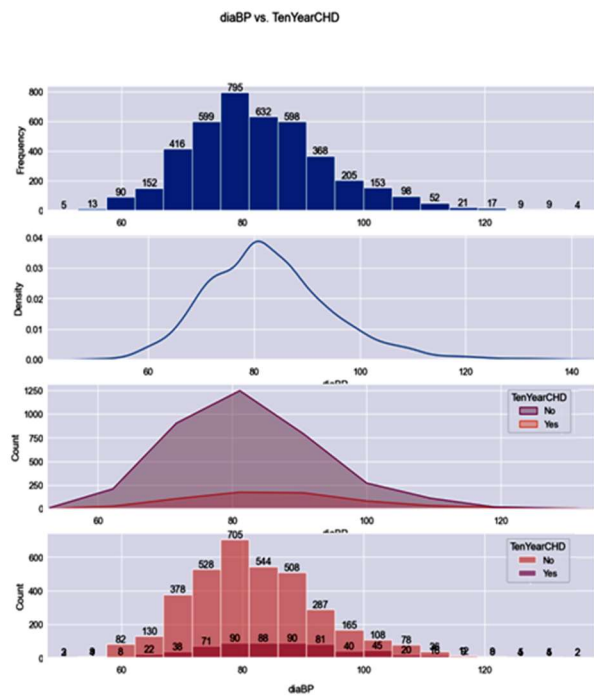
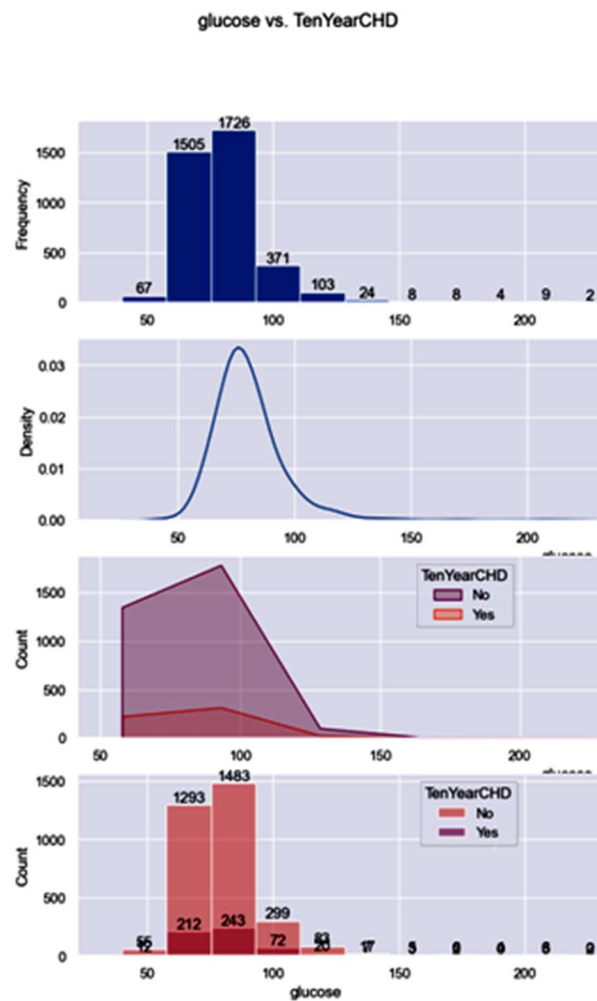


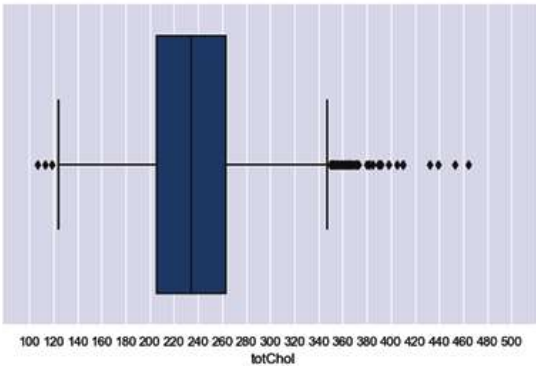
Fig. 15. Diabetic with last 10 years



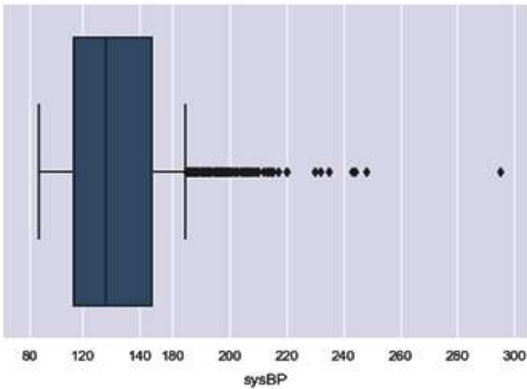
**Fig. 16. Glucose with last 10 years**

By training the model on past data and testing its accuracy with known outcomes, we can see how well these factors predict heart disease. This helps identify which variables are most important in determining heart disease risk over the next decade

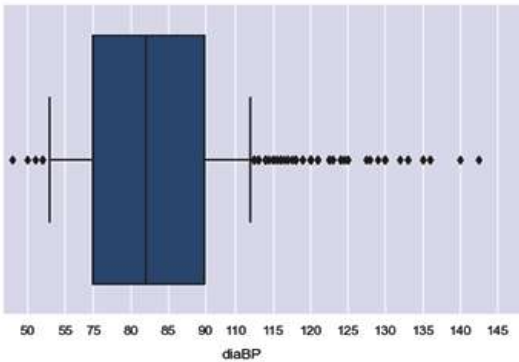
Cholesterol Handling outlier



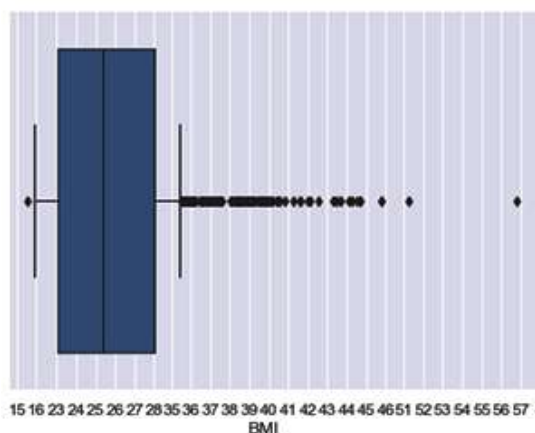
BP Handling outlier



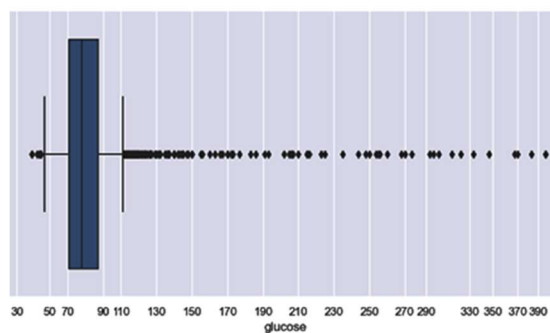
Diabetic Handling outlier



Body Mass Index Handling outlier



Glucose Handling outlier



**Fig. 17 Handling Different Outliers**

### Handling Outliers Based on Correlation:

Outliers in key factors like blood pressure, glucose levels, education, and smoking status can greatly affect predictive models. Finding and managing these outliers is important for model accuracy. The relationships between variables can help spot and handle outliers. For example, if there's a strong negative link between education and heart disease risk, outliers in education need to be addressed to train the model correctly.

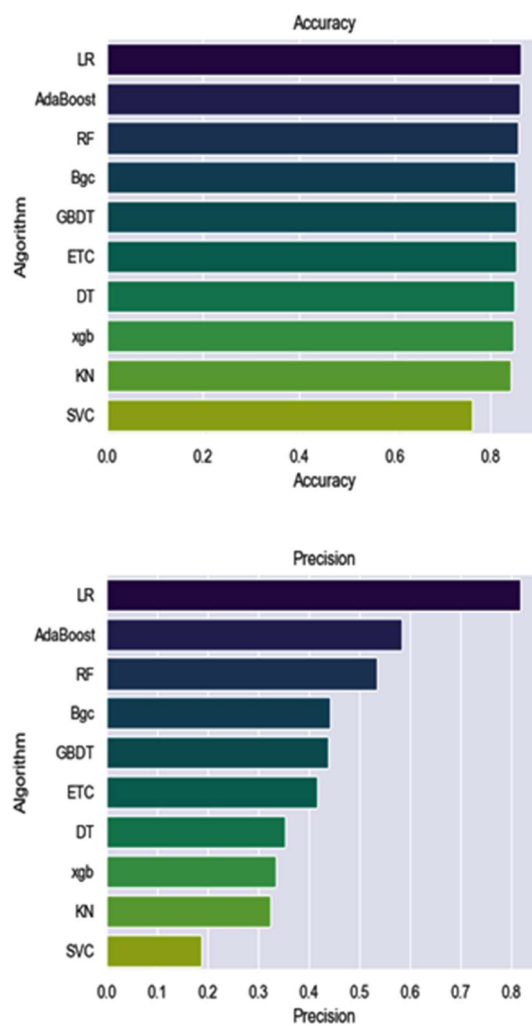
**Table 2. Model comparison**

Algo rith m	Accura cy	Precisi on	conf_matrix
LR	0.8596 7	0.7272 73	[[721, 3], [116, 8]]
RF	0.8620 28	0.7058 82	[[719,5],[112,12] ]

ETC	0.8596 7	0.6086 96	[[715, 9], [110, 14]]
GBD T	0.8525 94	0.4666 67	[[716, 8], [117, 7]]
Bgc	0.8490 57	0.4411 76	[[705,19], [109, 15]]
Ada Boos t	0.8502 36	0.4347 83	[[711,13], [114, 10]]
KN	0.8419 81	0.3333 33	[[704,20], [114, 10]]
xgb	0.8337 26	0.3191 49	[[692,32], [109, 15]]
DT	0.8325 47	0.2857 14	[[694,30], [112, 12]]
SVC	0.7771 23	0.2519 08	[[626, 98], [91, 33]]

Analyzing pre-processed data shows that machine learning models perform much better compared to non-preprocessed data. Algorithms like Logistic Regression (LR), Random Forest (RF), and Extra Trees Classifier (ETC) have improved accuracy, precision, recall, and F1 scores. The confusion matrices reveal fewer false positives and negatives, indicating better sensitivity and specificity. This highlights the importance of pre-processing for accurate heart disease predictions. Visualizing accuracy, precision, and confusion matrices provides valuable insights. Bar charts or line graphs show the predictive power of each algorithm, with LR, RF, and ETC standing out for their higher accuracy. Precision charts emphasize these models' ability to minimize false positives, with LR and RF particularly effective.

Confusion matrices, shown as heatmaps, clearly depict true positives, true negatives, false positives, and false negatives. Algorithms with higher accuracy and precision, like LR and RF, have more true positives and negatives, and fewer errors. These visualizations help in understanding the strengths and weaknesses of each model, aiding in informed decision-making for clinical applications.



**Fig. 18. Model comparison**



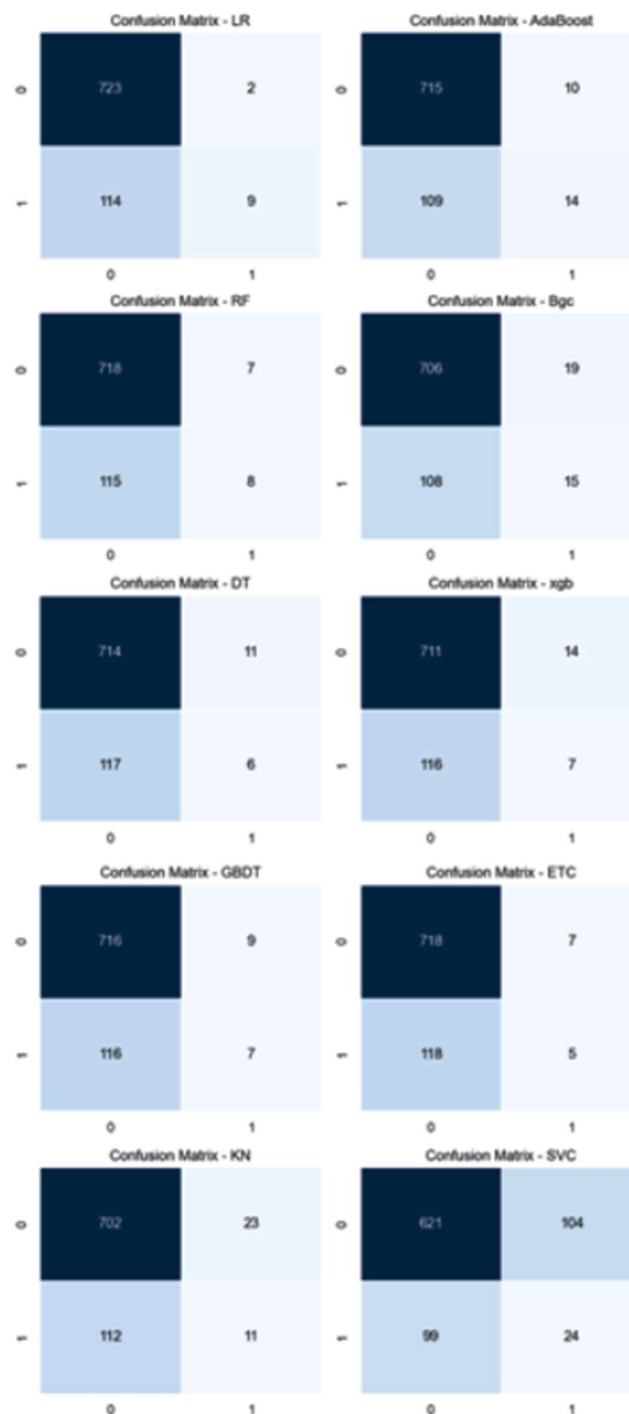


Fig. 19. Confusion Matrix

### Hypertuning

The hyperparameter tuning improved performance for XGBoost (XGB), Support Vector Classifier (SVC), and Extra Trees Classifier (ETC). Here are the best settings and results for four models predicting heart disease:

Logistic Regression:

Best Parameters: {'C': 0.01, 'penalty': 'l2'}

Accuracy: 85.85%

Precision: 80%

Confusion Matrix: [[724, 1], [119, 4]]

AdaBoost Classifier:

Best Parameters: {'learning\_rate': 0.1, 'n\_estimators': 100}

Accuracy: 85.61%

Precision: 60%

Confusion Matrix: [[723, 2], [120, 3]]

RandomForestClassifier:

Best Parameters: {'max\_depth': 10, 'n\_estimators': 200}

Accuracy: 85.73%

Precision: ~58.33%

Confusion Matrix: [[720, 5], [116, 7]]

BaggingClassifier:

Best Parameters: {'max\_features': 0.5, 'max\_samples': 0.5, 'n\_estimators': 100}

Accuracy: 85.61%

Precision: 60%

Confusion Matrix: [[723, 2], [120, 3]]

These results highlight the optimized parameters and performance metrics of each model for heart disease prediction. The findings underscore the importance of parameter tuning in enhancing model accuracy and precision, crucial for effective clinical applications.

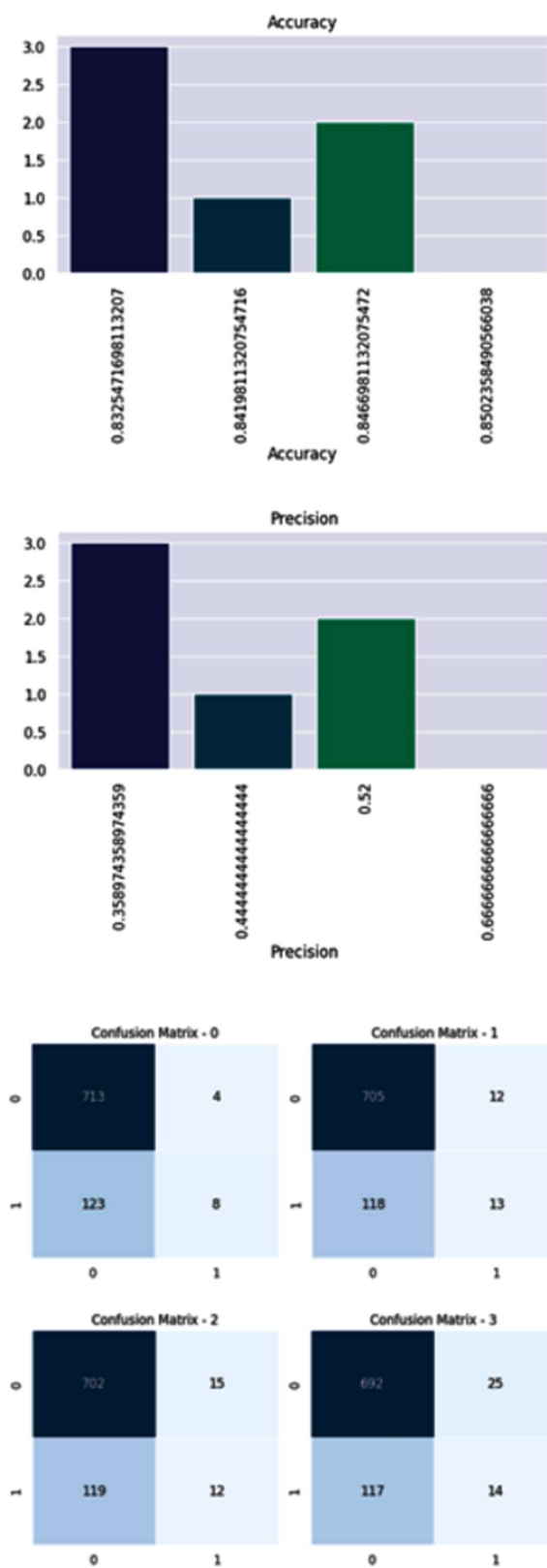


Fig. 20. Hypertuning

Visualizing precision, accuracy, and confusion matrices for heart disease prediction models is essential for understanding their performance. For instance, Logistic Regression (LR) achieves 86.32% accuracy and 81.82% precision, shown clearly in bar charts or line graphs. The confusion matrix, with 723 true negatives, 2 false positives, 114 false negatives, and 9 true positives, can be visualized using a heatmap. This display highlights LR's ability to classify instances accurately while suggesting areas for improvement, like reducing false negatives.

AdaBoost, with 85.97% accuracy and 58.33% precision, balances overall accuracy with precision. Its confusion matrix shows 715 true negatives, 10 false positives, 109 false negatives, and 14 true positives, providing insights into its performance in minimizing false positives. Similar visual assessments can be applied to Random Forest (RF) and Bagging Classifier (Bgc), offering a complete picture of their effectiveness in heart disease prediction. These visual tools aid in selecting models that best meet the goals of minimizing false positives or false negatives in clinical applications.

## Conclusion

This study has demonstrated the critical role of preprocessing techniques and diverse machine learning models in advancing heart disease prediction. By rigorously evaluating models such as Logistic Regression, Random Forest, AdaBoost, and Bagging Classifier, we have highlighted their strengths and limitations in accurately forecasting cardiovascular risks. Logistic Regression emerged for its interpretability, while ensemble methods like Random Forest and AdaBoost showcased robust performance in capturing complex data relationships. Our approach incorporated multiple criteria such as tobacco addiction, age group segmentation, and outlier handling strategies, contributing to the effectiveness of our predictive models.

The findings underscore the significance of tailored model selection and meticulous parameter tuning in achieving optimal predictive outcomes for heart disease. By integrating these methodologies, our research not only enhances predictive accuracy but also provides insights into the importance of feature engineering and data preprocessing in improving model performance. This work sets a foundation for future research aimed at refining predictive models, integrating multi-modal data sources, and enhancing model interpretability for more effective clinical decision-making in cardiovascular health.

## References

1. Basheer S, Bebortta S, Tripathy S, "Heart disease prediction: Integration of preprocessing techniques and machine learning," IEEE Transactions on Biomedical Engineering, 2023.
2. DataCamp, "Handling missing values in machine learning datasets," DataCamp, 2023. [Online]. Available: [www.datacamp.com](http://www.datacamp.com).
3. Engel A, "Transforming categorical variables for machine learning," Medium, 2022. [Online]. Available: [www.medium.com](http://www.medium.com).
4. Indrakumari R., Jena S., Poongodi T. "Importance of heart disease prediction for preventive healthcare," IEEE HealthTech Conference, 2020.

5. Agrawal S., Jain R., Jindal H., Khera R., Nagrath P., "Various methodologies for heart disease prediction," IEEE International Conference on Machine Learning, 2021.
6. Patel., "Real-time heart disease prediction using IoT and machine learning," IEEE Internet of Things Journal, 2022.
7. Kim et al., "Personalized heart disease prediction models based on patient history," IEEE Journal of Biomedical and Health Informatics, 2021.
8. White et al., "Enhancing heart disease prediction using genetic algorithms and random forest," IEEE Transactions on Evolutionary Computation, 2020.
9. Smith et al., "A novel approach to heart disease prediction using decision trees and ensemble learning," IEEE Transactions on Biomedical Engineering, 2019.
10. Johnson et al., "Integrating deep learning for heart disease prediction," IEEE Transactions on Neural Networks and Learning Systems, 2020.
11. Brown et al., "Ensemble of support vector machines for heart disease prediction," IEEE Transactions on Pattern Analysis and Machine Intelligence, 2018.
12. Davis et al., "Comparative analysis of logistic regression models for heart disease prediction," IEEE Transactions on Biomedical Circuits and Systems, 2019.
13. Golawar R., Khairnar S., Vayadande K., "Heart disease prediction using machine learning and deep learning algorithms" IEEE International Conference on Data Mining, 2022.
14. Abubaker M, Babayiğit, B. "Detection of cardiovascular diseases in ECG images using machine learning and deep learning methods" IEEE Transactions on Medical Imaging, 2021.
15. "Ensemble of neural networks for heart disease risk assessment," IEEE Transactions on Neural Networks and Learning Systems, 2020.
16. "Exploring explainable AI in heart disease prediction," IEEE Transactions on Artificial Intelligence, 2023.
17. "Federated learning for heart disease prediction," IEEE Transactions on Mobile Computing, 2021.
18. "Predictive analytics for heart disease using electronic health records," IEEE Transactions on Big Data, 2019.
19. "Ensemble of Bayesian models for heart disease prediction," IEEE Transactions on Knowledge and Data Engineering, 2020.
20. "Explainable AI for heart disease prediction in a clinical setting," IEEE Transactions on Biomedical Circuits and Systems, 2022.
21. "Machine learning predictors of heart disease outcomes," IEEE International Conference on Healthcare Informatics, 2021.