



**Applying Machine Learning Techniques and Optimization to Diagnose Ischemic Heart Disease (IHD) in Diabetes Patients in Bangladesh: Investigations from a Cross-Sectional Study in 2024**

**Submitted By Mir khalid Hassan**

**ID: 20231167**

**Batch: 11th**

**Section: B**

# Index

Topics	Page
<b>1. Introduction</b>	<b>1-2</b>
<b>2. Literature Review</b>	<b>3-4</b>
<b>3. Methods</b>	<b>5-7</b>
<b>3.1 Study Design and Data Collection</b>	<b>5</b>
3.1.1 Variables and Data Description	5
3.1.2 Target Variable Definition	
<b>3.2 Data Preprocessing</b>	<b>6-7</b>
3.2.1 Handling Missing Values	6
3.2.2 Removing Redundant Features and Duplicates	6
3.2.4 Feature Scaling	6
3.2.5 SMOTE (Synthetic Minority Over-sampling Technique)	6
3.2.3 Outlier Detection and Treatment	6
<b>3.3 Feature Selection</b>	<b>8-12</b>
3.3.1 Mutual Information (MI) Score Analysis	8
3.3.2 Random Forest Feature Importance	9
<b>3.4 Machine Learning Models</b>	<b>12-13</b>
3.4.1 Baseline Models	12
3.4.2 Advanced Models	13
<b>3.5 Hyperparameter Tuning and Ensemble Learning</b>	<b>13</b>
3.5.1 Hyperparameter Optimization	13
3.5.2 Ensemble Learning	13
<b>3.6 Performance Metrics</b>	<b>14</b>
3.6.1 Classification Performance	14
3.6.2 Probabilistic and Ranking Metrics	14
<b>4. Results</b>	
<b>4.1 Descriptive Statistics</b>	<b>15</b>
4.1.1 Summary Statistics of Numerical Features	15
4.1.2 Summary of Categorical Features	16
4.1.3 Outlier Analysis	16
<b>4.2 Inferential Statistics</b>	<b>17</b>
4.2.1 Relationship Between Diabetes Mellitus (DM) and ischemic heart disease (IHD)	17
<b>4.3 Model Performance</b>	<b>18-27</b>
4.3.1 Model Comparison	18
4.3.2 Initialize the model with the best parameters	21
4.3.3 Confusion Matrix	21
4.3.4 ROC Curve and Precision-Recall Analysis	22
4.3.5 Cross-Validation Performance	24
4.3.6 Best Model Recommendation	25
4.3.7 Soft voting ensemble method	26
4.3.8 Model Performance Comparison	27
<b>5. Discussion</b>	
<b>5.1 Interpretation of Results</b>	<b>29</b>
<b>5.2 Comparison with Previous Studies</b>	<b>31</b>
<b>5.3 Implications for Public Health</b>	<b>33</b>
<b>5.4 Limitations of the Study</b>	<b>33</b>
<b>5.5 Future Work</b>	<b>33</b>
<b>6. Conclusion</b>	<b>35</b>
<b>7. References</b>	<b>36</b>
<b>8. Appendices</b>	<b>39</b>
<b>Code Snippets</b>	<b>42</b>

## Letter of Endorsement

### To Whom It May Concern:

The research project called "**Applying Machine Learning Techniques and Optimization to Diagnose Ischemic Heart Disease (IHD) in Diabetes Patients in Bangladesh: Investigations from a Cross-Sectional Study in 2024**" was turned in by **Mir Khalid Hassan** (ID: 20231167, Batch: 11th, Section: B) and has been approved because it is academically sound. I oversaw the research at **Jahangirnagar University's Department of Statistics and Data Science**.

I have meticulously observed the advancement of this project and have actively contributed guidance at each phase of its evolution. The student has demonstrated exceptional commitment and diligence in conducting this study. This letter constitutes formal acceptance of the research, and I encourage its submission as a requisite for the course.

Please don't hesitate to contact me if you need more explanation or information.

**Respectfully,**

**Md. Habibur Rahman,**

Associate Professor, Department of Statistics and Data Science,

Jahangirnagar University, Savar, Dhaka-1342, Bangladesh.

Cell Phone: +8801942246023.

Email: [habib.drj@juniv.edu](mailto:habib.drj@juniv.edu), [habib.drj@gmail.com](mailto:habib.drj@gmail.com).

## Recognition

I wish to convey my deep appreciation to my respected supervisor, **Md. Habibur Rahman, Associate Professor, Department of Statistics and Data Science, Jahangirnagar University**, for his outstanding supervision, steadfast support, and essential insights during this project. His proficiency and thorough methodology have been important in developing this work and achieving its completion.

I am profoundly grateful to the faculty and staff of the **Department of Statistics and Data Science** at **Jahangirnagar University** for supplying essential resources and for their unwavering support throughout my research.

I would like to express my profound gratitude to all participants in this study, whose collaboration and engagement have been essential to the successful completion of this project.

Finally, I express my gratitude to my family and friends for their unwavering encouragement, patience, and understanding, which have served as a continual source of motivation during my academic pursuits.

### **Submitted By:**

**Mir Khalid Hassan**  
**ID: 20231167**  
**Batch: 11th**  
**Section: B**

## Abstract

**Background:** Ischemic heart disease (IHD) is a predominant cause of morbidity and mortality, especially among diabetic individuals in Bangladesh. Early detection of ischemic heart disease is essential for appropriate intervention and improved patient outcomes. **Machine learning (ML)** methodologies have demonstrated potential in improving diagnostic precision for ischemic heart disease (IHD). This project is to utilize machine learning methodologies and optimization techniques to enhance the identification of ischemic heart disease in diabetic patients in Bangladesh. The main goal of this research is to examine the efficacy of machine learning techniques, encompassing algorithm optimization and ensemble learning, for identifying ischemic heart disease in diabetic patients via a cross-sectional study done in 2024.

**Methods:** A dataset comprising clinical, demographic, and laboratory data from diabetic patients was analyzed utilizing fourteen distinct machine learning algorithms: **Logistic Regression (LR)**, **k-Nearest Neighbors (kNN)**, **Naive Bayes (NB)**, **Decision Tree (DT)**, **Support Vector Machine (SVM)**, **Ridge Classifier (RC)**, **Random Forest (RF)**, **Quadratic Discriminant Analysis (QDA)**, **AdaBoost**, **Gradient Boosting (GB)**, **Linear Discriminant Analysis (LDA)**, **Extra Trees Classifier (ETC)**, **Classifier Chain (CC)**, and **Decision Forest (DF)**. **Five-fold cross-validation** was utilized for **hyperparameter adjustment** to enhance the model's predictive accuracy. The effectiveness of various models was evaluated based on accuracy and precision, with a specific focus on assessing their generalization abilities using accuracy analysis for each fold.

**Results:** The evaluation metrics indicate that Gradient Boosting is the most effective model, with an accuracy of 0.910 and a ROC AUC of 0.9694. It consistently outperforms other models in these areas, indicating that it is the most reliable model for classification tasks where predicted accuracy and class distinction are critical. However, in terms of log loss, Random Forest demonstrates higher performance with a lower value of 0.2493, signifying its enhanced reliability in probabilistic predictions. Nevertheless, Gradient Boosting's exceptional performance in accuracy and ROC AUC makes it the preferable choice for most applications.

**Conclusions:** The study highlights the potential of **Machine Learning Techniques** and Optimization, with Gradient Boosting being the most effective model for overall classification performance. Random Forest is a viable alternative for reducing log loss. Soft Voting Ensemble, while successful, falls short of the top models. The study also suggests that integrating GridSearchCV, five-fold cross-validation, and soft voting ensemble classifiers can help in early diagnosis and treatment planning for diabetes patients with ischemic heart disease.