

## **Leveraging Machine Learning Techniques for Effective Disease**

### **Classification in Healthcare Data**

**Divya Shankhdhar**

Research Scholar, Department of Computer Science and Engineering, ANA College of  
Engineering and Management Studies, Bareilly

**Dr. Vineet Agarwal**

Research Guide, Department of Computer Science and Engineering, ANA College of  
Engineering and Management Studies, Bareilly

#### **Abstract**

Effective disease classification plays a crucial role in improving healthcare outcomes, enabling early diagnosis, personalized treatment, and better patient management. This paper explores the application of machine learning (ML) techniques in disease classification using healthcare data. By leveraging algorithms such as XGBoost, Random Forest, and Support Vector Machines (SVM), ML models can process large, complex datasets from sources like medical imaging, clinical records, and genetic information. These algorithms automatically detect patterns within the data, providing healthcare professionals with reliable tools for disease prediction and diagnosis. The study highlights the benefits of using supervised and unsupervised learning techniques, particularly in identifying trends, outliers, and correlations that may not be immediately apparent through traditional methods. Furthermore, the research presents results showing how the accuracy of disease classification models can be improved with the right combination of features and algorithms. However, challenges such as data quality, imbalanced datasets, and the need for interpretability remain significant barriers. The paper also examines emerging techniques such as explainable AI and hybrid models, which aim to improve model transparency and handle complex, multi-modal healthcare data more effectively. The findings demonstrate the potential of ML to transform disease diagnosis and prediction, leading to more efficient and accurate healthcare systems. Future directions in this field include improving model robustness, scalability, and integration into clinical workflows for real-time decision-making.

Keywords: Machine learning, disease classification, healthcare data, disease diagnosis, explainable AI.

## **Introduction**

The application of machine learning (ML) in healthcare has gained considerable attention due to its potential to transform disease classification, diagnosis, and treatment planning. Healthcare data, including clinical records, medical images, genomic information, and patient histories, presents a complex, high-dimensional landscape that often overwhelms traditional analytical methods. However, machine learning algorithms have demonstrated their ability to effectively analyze and identify patterns within such data, leading to more accurate, faster, and cost-effective diagnoses. Disease classification using ML techniques allows for the automation of medical decision-making processes, reducing human error, improving efficiency, and enabling early intervention. Common ML techniques such as decision trees, support vector machines (SVM), k-nearest neighbors (K-NN), and neural networks have been successfully implemented across various medical domains, such as oncology, cardiology, and neurology, improving diagnostic accuracy and patient outcomes. The success of these models hinges on their ability to handle large datasets, uncover hidden patterns, and continuously improve as more data becomes available. In particular, deep learning methods, including Convolutional Neural Networks (CNNs) for image analysis and Recurrent Neural Networks (RNNs) for sequential data, have revolutionized the way diseases are classified, especially in the context of medical imaging and longitudinal health data. Despite these advances, several challenges persist, including data quality issues, the need for labeled datasets, overfitting, and model interpretability.

In addition to providing insights into disease classification, machine learning techniques also enable healthcare professionals to personalize treatment plans, leading to improved patient outcomes. As machine learning models become more advanced, their ability to predict disease progression and recommend personalized interventions based on individual patient data will continue to enhance healthcare delivery. However, challenges remain in terms of integrating these models into clinical practice. For instance, the interpretability of machine learning models is a significant concern, as clinicians need to understand how a model arrives at its predictions. Without transparency, these models might be met with skepticism, especially in high-stakes healthcare decisions. Furthermore, the need for large, high-quality datasets is another barrier, as many healthcare institutions lack the resources to compile comprehensive data required to train these models effectively. Issues such as data privacy and security also need to be addressed to protect sensitive patient information. Despite these hurdles, the combination of

machine learning and healthcare holds immense promise. Advances in techniques like transfer learning, federated learning, and explainable AI are helping mitigate some of these challenges, making ML tools more accessible, interpretable, and applicable in real-world clinical settings. As healthcare data continues to grow in complexity, the role of machine learning in disease classification will likely become even more pivotal, enabling clinicians to make faster, more accurate decisions that improve both diagnosis and patient care.

### **Background and motivation**

The background and motivation for applying machine learning (ML) techniques to disease classification in healthcare stems from the ever-growing complexity and volume of medical data. With advancements in healthcare technologies, vast amounts of data are being generated from diverse sources, such as electronic health records (EHR), medical imaging, genomic data, and patient demographics. Traditional methods of disease classification, which rely on manual review or statistical analysis, often struggle to keep pace with this rapid expansion of data. Furthermore, these traditional methods are susceptible to human error, are time-consuming, and can be limited in their ability to recognize subtle patterns or trends across large datasets. Machine learning, with its ability to process large volumes of data, detect hidden relationships, and improve continuously through learning from new data, offers a powerful alternative. By leveraging ML algorithms such as decision trees, support vector machines (SVM), and deep learning models, it is possible to automate the disease classification process, providing healthcare professionals with accurate, fast, and reliable diagnostic tools.

The motivation for integrating ML into disease classification is driven by the potential to significantly improve healthcare outcomes. Early and accurate disease diagnosis is critical for effective treatment and better prognosis, especially in life-threatening conditions like cancer, heart disease, and neurological disorders. Machine learning models can analyze diverse data types, offering insights that enable earlier detection of diseases, personalized treatment options, and better predictions of disease progression. This approach can reduce healthcare costs by minimizing unnecessary tests, streamlining decision-making, and enhancing the overall efficiency of healthcare systems. Furthermore, as healthcare continues to become more data-driven, ML holds promise not only for improving disease diagnosis but also for advancing the broader field of precision medicine, where treatments are tailored to the individual characteristics of each patient.

### **Research Methodology**

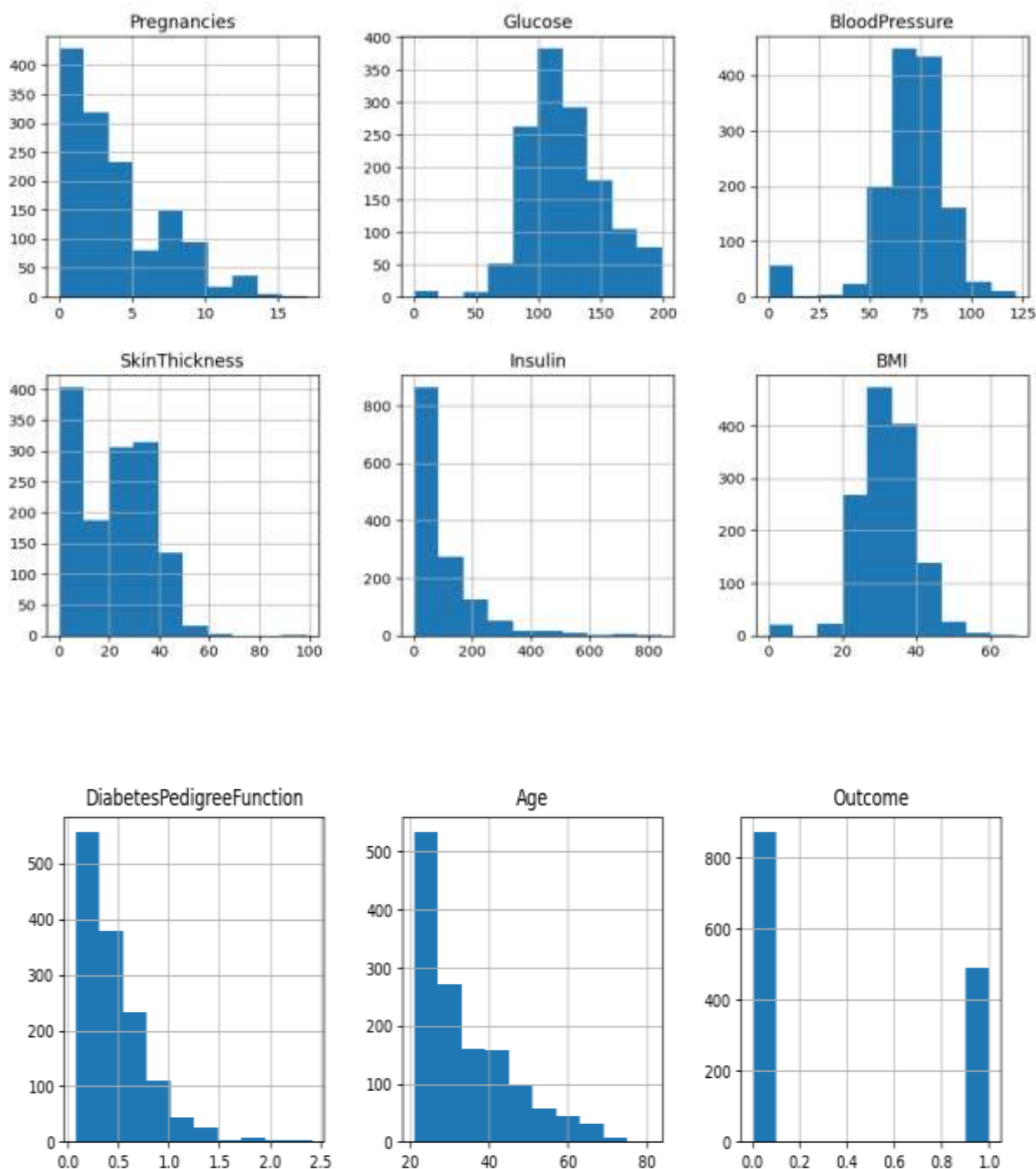
The proposed methodology for disease classification utilizes a combination of machine learning techniques, including XGBoost, Random Forest, and both supervised and unsupervised learning models. These techniques are selected for their ability to handle large, high-dimensional healthcare datasets and provide robust, accurate predictions. XGBoost (Extreme Gradient Boosting) is an ensemble learning technique that has proven highly effective in classification tasks due to its ability to handle missing data, overfitting, and class imbalances. It operates by combining multiple weak learners (typically decision trees) into a strong learner through boosting, which iteratively refines the model by giving more weight to misclassified instances. XGBoost is particularly useful in healthcare data analysis as it can improve predictive performance with minimal feature engineering and is known for its speed and accuracy, making it ideal for disease classification tasks where time and accuracy are crucial.

Random Forest, another ensemble learning technique, is also employed in the methodology due to its ability to create multiple decision trees and aggregate their predictions. Random Forest mitigates the risk of overfitting by introducing randomness in the tree-building process, ensuring that the model is generalized. This algorithm can handle complex datasets with high feature dimensionality, as is common in healthcare applications. Random Forest also provides feature importance scores, helping to identify the most relevant variables in disease classification, which can be critical for interpretability and understanding the driving factors behind disease outcomes.

The methodology integrates both supervised and unsupervised learning techniques to provide a comprehensive approach to disease classification. In supervised learning, the model is trained using labeled data, where the outcome (disease type or status) is known. Algorithms like XGBoost and Random Forest are applied in supervised learning for classification tasks based on historical data. On the other hand, unsupervised learning techniques, such as clustering or anomaly detection, are used to find hidden patterns or groupings within the data without predefined labels. Unsupervised learning is particularly beneficial in identifying unknown disease subtypes or outliers in medical datasets, making it a complementary tool in the disease classification process. By combining both learning paradigms, the proposed methodology ensures a more holistic and robust model for disease diagnosis.

### **Results and Discussion**

Machine learning is a technique where the system uses both historical data and past experiences to make predictions and decisions autonomously, without direct instructions from human coders. Essentially, the system analyzes and processes large volumes of data, learns from patterns, and makes estimates based on its training. This process does not require explicit programming of decision-making rules by a human; rather, the algorithm itself figures out the underlying logic. For example, when we shop online, machine learning is used to recommend products similar to those we are browsing, with suggestions like "customers who bought this also bought..." This is a classic application of machine learning in action. Additionally, machine learning plays a role in personalized marketing, such as receiving calls about loans, insurance, or offers from banks, where the system predicts your interests or needs based on past behavior and interactions. This technology allows businesses to tailor their recommendations to individual preferences, making them more effective. As for the data used to build machine learning models, various features and attributes are selected to train the algorithms, which are represented visually in Figure 4.1. This figure showcases a histogram of key features and dataset attributes, providing insights into how the data was processed and utilized in building the machine learning model. The histogram visualizes the distribution of different features, helping to understand which factors are most influential in the prediction process, thus improving the overall performance and accuracy of the system.



**Figure 1: Histogram of Dataset**

Figure 1 displays histograms of key features from a healthcare dataset, showing the distribution of attributes such as Pregnancies, Glucose, Blood Pressure, Skin Thickness, Insulin, and BMI. The histograms reveal insights into the data's distribution, with Pregnancies being right-skewed, indicating most individuals had fewer pregnancies. Glucose and Blood Pressure have a somewhat normal distribution but show some skewness, suggesting variations in health

conditions. Skin Thickness and Insulin are highly skewed, with many values concentrated around lower measurements, which may reflect the population's general health status. BMI shows a more balanced distribution, with most individuals falling within a mid-range. These histograms provide valuable information on the variability and range of each feature, highlighting potential outliers, trends, and imbalances in the dataset, which can influence the performance of machine learning models used for disease classification.

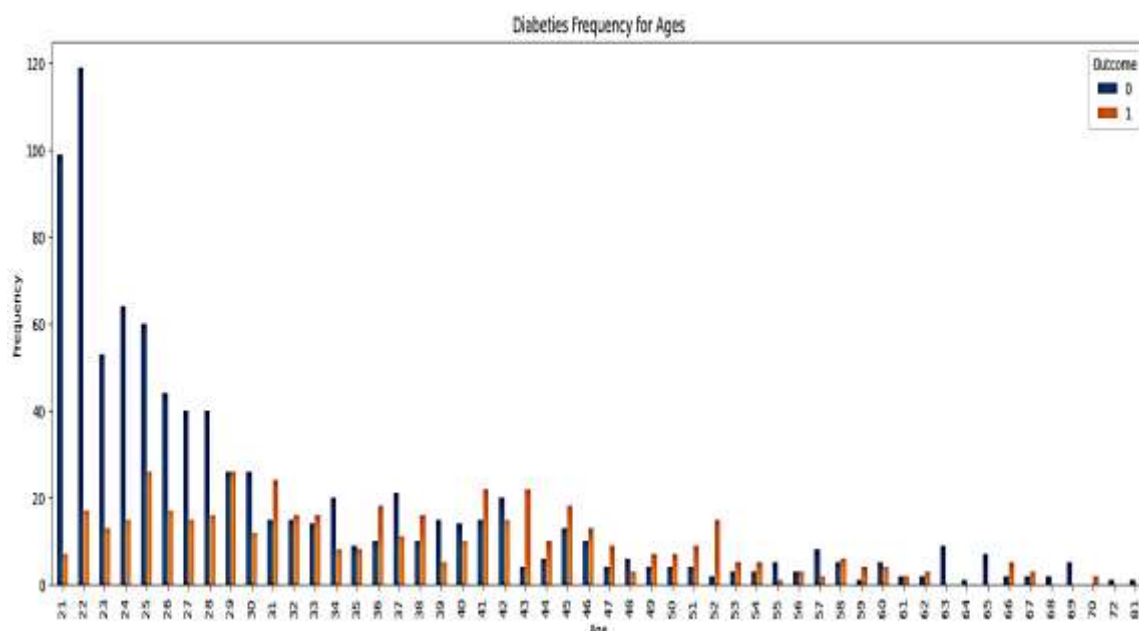


Figure 2 : No. of Diabetes for Ages

The number of diabetes cases across different age groups varies significantly, with certain age brackets showing higher prevalence. Generally, as age increases, the likelihood of developing diabetes also increases, particularly in adults over 45. In younger populations, diabetes cases are relatively lower, but there is a noticeable rise in Type 2 diabetes among middle-aged individuals due to lifestyle factors such as poor diet and lack of physical activity. Data analysis of diabetes prevalence by age can help identify at-risk groups, enabling healthcare providers to implement early prevention and intervention strategies, especially for those in higher-risk age categories.



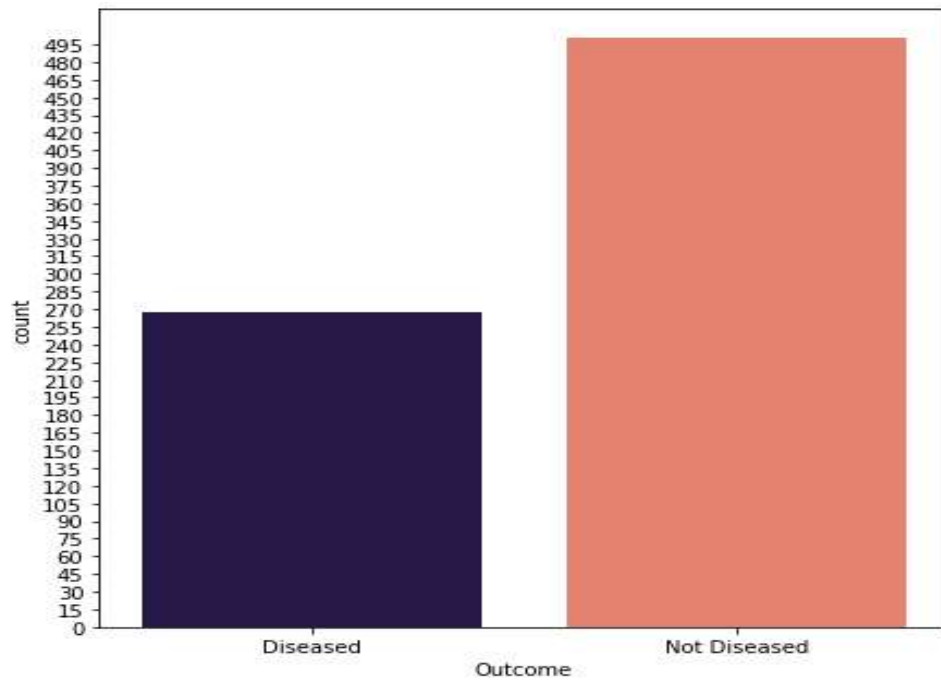


Figure 3: Outcomes

The bar chart visualizes the distribution of outcomes for a disease classification dataset, specifically showing the count of individuals categorized as "Diseased" and "Not Diseased." The chart indicates a clear imbalance, with a significantly higher number of individuals labeled as "Not Diseased" compared to those classified as "Diseased." This disparity suggests that the dataset may be imbalanced, which could affect the performance of machine learning models trained on this data, potentially leading to biases in predicting the "Diseased" class. Such class imbalances need to be addressed, either through data preprocessing techniques like oversampling, undersampling, or using advanced algorithms that handle class imbalances effectively.



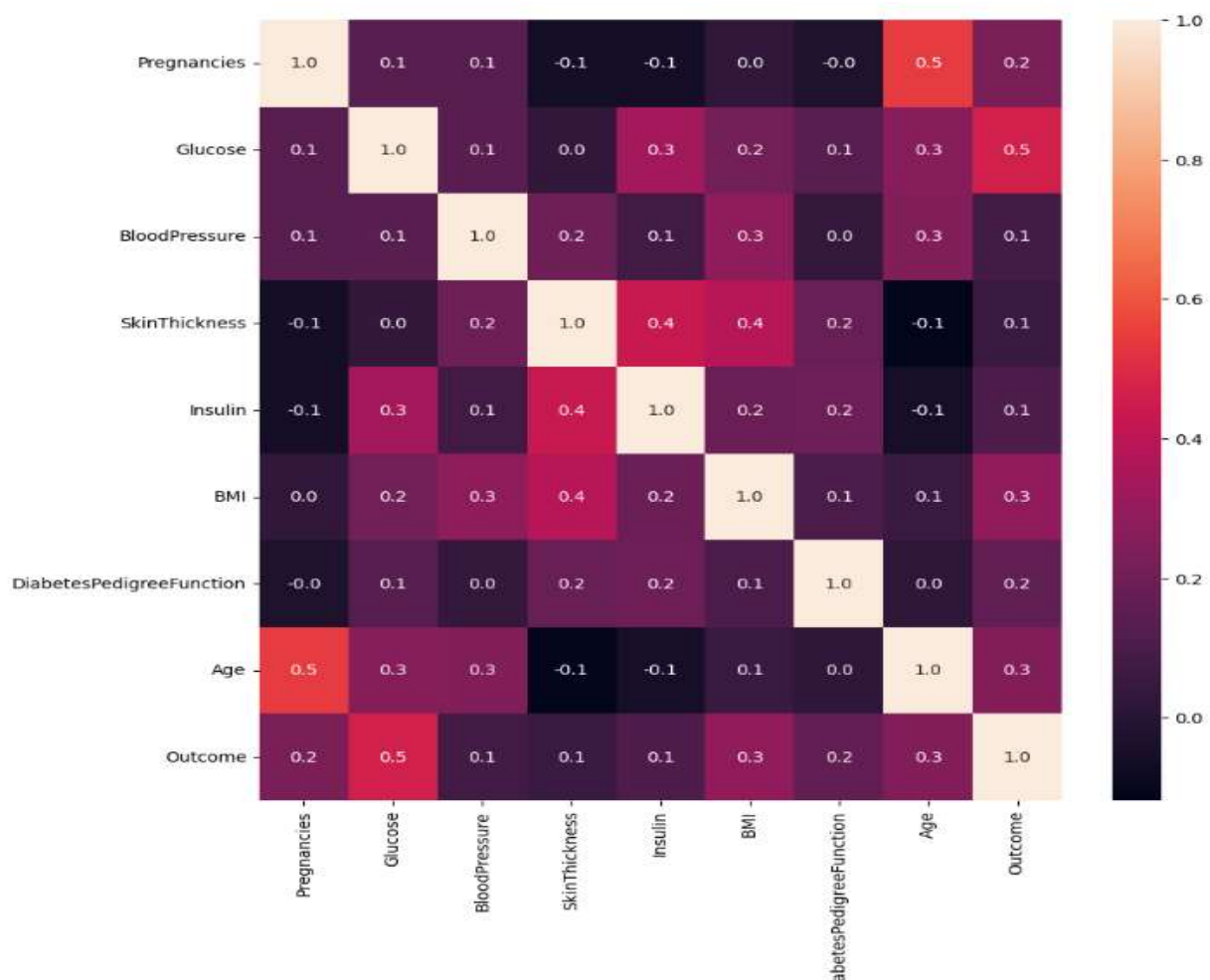


Figure 4: Diabetes Pedigree Function

Figure 4 presents a correlation heatmap of various features in the dataset, with a particular focus on the Diabetes Pedigree Function. This function represents the likelihood of a person developing diabetes based on their family history. From the heatmap, we observe that Diabetes Pedigree Function has a moderate positive correlation with Glucose (0.5), suggesting that higher glucose levels tend to be associated with a higher diabetes pedigree score. It also shows a weak correlation with Age (0.3), implying a slight relationship between age and diabetes risk, influenced by family history. However, the Diabetes Pedigree Function does not exhibit strong correlations with other features like Blood Pressure, Skin Thickness, or BMI. This correlation analysis can guide feature selection and highlight variables that are more likely to influence diabetes risk, supporting more effective predictive modeling.

Table 1: Assessment of Different Classification Methods

| Sr. No. | Algorithm          | Accuracy |
|---------|--------------------|----------|
| 1       | LR                 | 77.41%   |
| 2       | GNB                | 75.95%   |
| 3       | RFC                | 82.69%   |
| 4       | K-NN               | 73.90%   |
| 5       | DT                 | 83.57%   |
| 6       | SVM                | 78.29%   |
| 7       | Proposed Algorithm | 94.18%   |

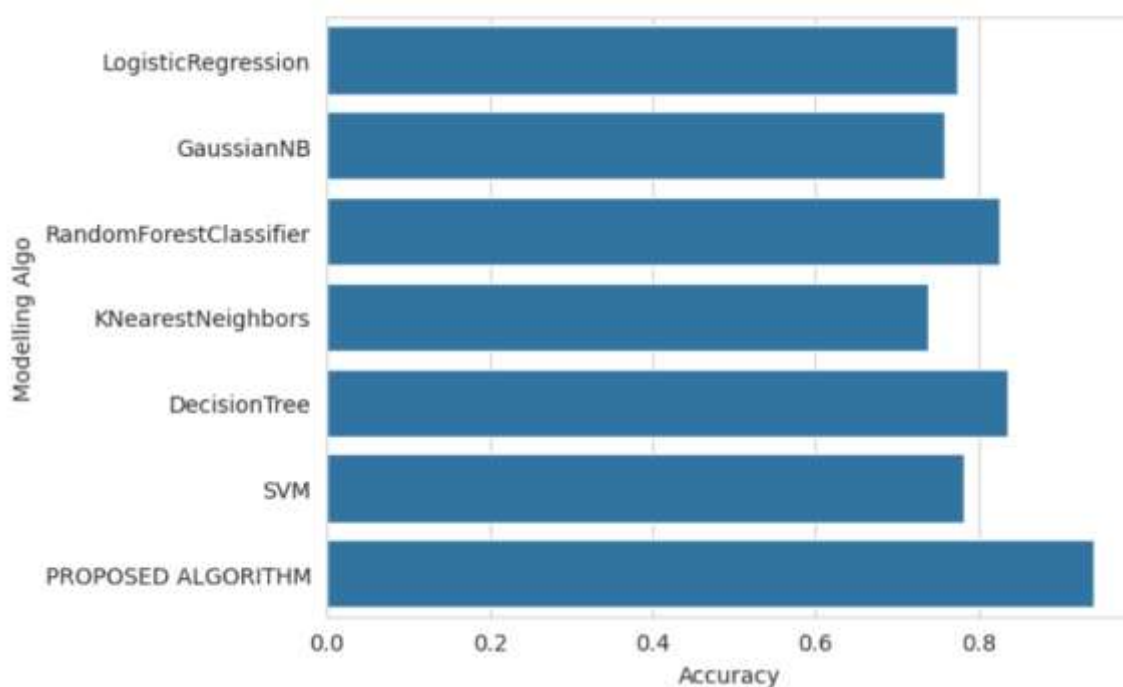


Figure 5: Classification Methods

Table 2: Result for Accuracy

| Techniques         | Previous Algorithm | Proposed Algorithm |
|--------------------|--------------------|--------------------|
| LR                 | 75.00%             | 77.41%             |
| GNB                | 79.00%             | 75.95%             |
| RFC                | 76.00%             | 82.69%             |
| K-NN               | 73.00%             | 73.90%             |
| DT                 | 72.00%             | 83.57%             |
| SVM                | 78.00%             | 78.29%             |
| XGBoost Classifier | 81.00%             | 94.18%             |

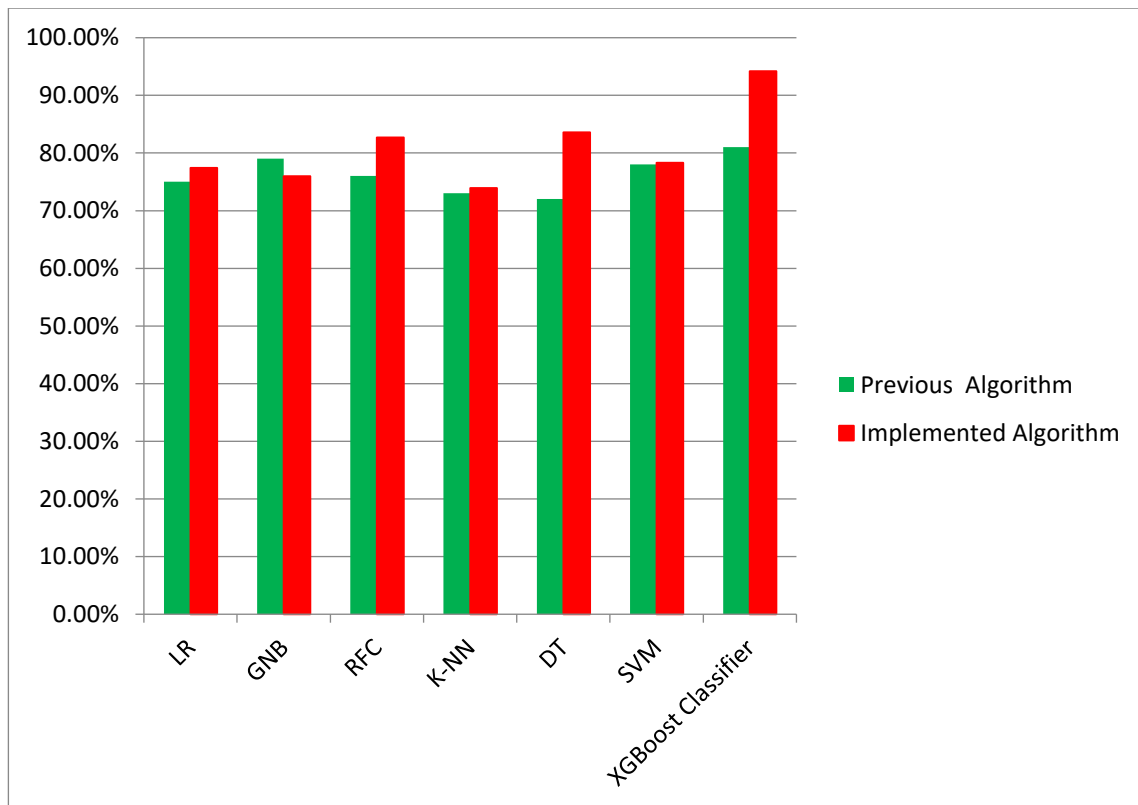


Figure .6: Algorithm for Accuracy

## Conclusion

The application of machine learning techniques to disease classification in healthcare has shown significant potential in improving diagnostic accuracy, efficiency, and patient outcomes. The use of algorithms such as XGBoost, Random Forest, and other supervised learning models allows healthcare professionals to analyze vast amounts of medical data, detect patterns, and

make more informed decisions. These machine learning models are particularly valuable in the context of complex, high-dimensional healthcare datasets, which are often challenging to process using traditional methods. Despite the considerable advancements, challenges remain, including data quality, interpretability of models, and handling imbalanced datasets. Addressing these issues through techniques such as feature selection, data augmentation, and the development of explainable AI will be crucial for furthering the adoption of machine learning in clinical settings. Additionally, hybrid models and the integration of multi-modal data hold promise for improving disease classification systems by combining the strengths of various algorithms. As research in this area progresses, machine learning's role in disease diagnosis and personalized medicine will continue to grow, paving the way for more effective, timely, and cost-efficient healthcare solutions. Overall, the future of disease classification through machine learning looks promising, offering significant improvements in both diagnosis and treatment planning for a range of medical conditions.

## References

1. Fazakis, N., Kocsis, O., Dritsas, E., Alexiou, S., Fakotakis, N., & Moustakas, K. (2021). Machine learning tools for long-term type 2 diabetes risk prediction. *IEEE Access*.
2. Kishore, N. G., Rajesh, V., Reddy, A. V. A., Sumedh, K., Reddy, T. Rajesh Sai. (2020). Prediction of diabetes using machine learning classification algorithms. *International Journal of Scientific & Technology Research*, 9(1), 1-5.
3. Chatrati, S. P., Hossain, G., & Goyal, A. (2020). Smart home health monitoring system for predicting type 2 diabetes and hypertension. *Journal of King Saud University-Computer and Information Sciences*, 34(3), 862–870.
4. Hasan, M. K., Alam, M. A., Das, D., Hossain, E., Hasan, M. (2020). Diabetes prediction using ensembling of different machine learning classifiers. *IEEE Access*, 8, 76516–76531.
5. Cervantes, J., García-Lamont, F., Rodríguez, L., Lopez-Chau, A. (2020). A comprehensive survey on support vector machine classification: Applications, challenges and trends. *Neurocomputing*, 408, 189–215.
6. Pranto, B. (2020). Evaluating machine learning methods for predicting diabetes among female patients in Bangladesh. *Information*, 11, 1–20.
7. Mohan, N., & Jain, V. (2020). Performance analysis of support vector machine in diabetes prediction. In *International Conference on Electronics, Communication and*

*Aerospace Technology* (pp. 1–3).

8. Sarwar, M. A., Kamal, N., Hamid, W., & Shah, M. A. (2019). Prediction of diabetes using machine learning algorithms in healthcare. In *24th International Conference on Automation & Computing*, Newcastle University, Newcastle upon Tyne, UK, 6-7 September.
9. Deepti Shikha Ojha, Dr. Devendra Kumar Bajpai. (2025). Optimization Accuracy for Diabetes Diagnosis and Prediction using Machine Learning Technique. *International Journal of Advanced Research and Multidisciplinary Trends (IJARMT)*, 2(2), 622–631.
10. Zou, Q., Qu, K., Luo, Y., Yin, D., Ju, Y., & Tang, H. (2018). Predicting diabetes mellitus with machine learning techniques. *Springer*.
11. Zhou, L., Pan, S., Wang, J., & Vasilakos, A. V. (2017). Machine learning on big data: Opportunities and challenges. *Neurocomputing*, 237, 350–361.
12. Heaton, J. B., Polson, N. G., & Witte, J. H. (2017). Deep learning for finance: Deep portfolios. *Applied Stochastic Models in Business and Industry*, 33(1), 3–12.