

Theoretical Perspectives on Machine Learning in the Diagnosis of Coronary Heart Disease: A Comparative Framework for Clinical Decision Support

Md Rahat Hossain¹, Azad Rahman^{2*}

¹Department of Mechanical Engineering, Yangzhou University, China

Email: mdrahathossain74@gmail.com

²Department of Electrical and Electronics Engineering, Daffodil International University, Dhaka, Bangladesh

Email: mrzad.eee@gmail.com

Abstract

Coronary Heart Disease (CHD) remains one of the most prevalent causes of global mortality, posing a continuous challenge to modern healthcare systems. Traditional diagnostic methods, while effective, are often limited by subjectivity, time constraints, and the increasing complexity of patient data. Machine learning (ML) offers a theoretical framework that can revolutionize CHD diagnosis by enabling automated, data-driven decision-making. This paper explores the theoretical basis of using ML algorithms—specifically Logistic Regression, Random Forest, and Support Vector Machine (SVM)—for CHD diagnosis. We present a conceptual comparison of these models in terms of learning mechanisms, data interpretability, clinical applicability, and model complexity. Rather than focusing on numerical results, this paper provides a high-level analysis of how these models theoretically contribute to improving diagnostic accuracy, patient stratification, and clinical workflow efficiency.

Keywords

Coronary Heart Disease, Machine Learning, Classification.

1. Introduction

The growing incidence of coronary heart disease (CHD) has highlighted the need for improved diagnostic tools capable of early identification and prevention. Traditional diagnostic techniques, including stress testing, angiography, and physician assessment, require manual interpretation and clinical judgment. While effective, these methods may not scale efficiently in the face of large patient volumes and increasing complexity in medical data [1].

Machine learning (ML), an evolving field within artificial intelligence, offers theoretical constructs for automated learning from data. ML models can identify complex patterns and

relationships in clinical variables that are not easily captured through conventional methods. In this theoretical analysis, we explore three prominent ML algorithms—Logistic Regression, Random Forest, and Support Vector Machine—and examine their conceptual roles in supporting the diagnosis of CHD. This work does not aim to present empirical performance results, but rather to understand the underlying theoretical advantages and limitations of each model in a healthcare context [16][17].

2. Theoretical Foundations of Machine Learning in Healthcare

The theoretical underpinnings of machine learning (ML) in the healthcare sector represent a significant paradigm shift from traditional statistical inference to intelligent, data-driven decision-making. While conventional statistical models often rely on assumptions of linearity, normality, and independence among variables, machine learning is designed to discover patterns from data without requiring such strict assumptions. This makes ML particularly suitable for clinical environments where data is often heterogeneous, incomplete, and nonlinear in nature [1]. The application of ML in healthcare is supported by a foundational understanding of how these models learn from data to make predictions, classify outcomes, and even recommend treatment pathways.

In the context of coronary heart disease (CHD), the predictive task is typically a binary classification problem—identifying whether a patient is likely to develop or currently has the disease based on a set of observed clinical features. These features may include age, sex, cholesterol levels, blood pressure, blood sugar, and ECG results, among others [2]. ML models use such input features to establish complex mappings between these attributes and the likelihood of disease presence. Unlike rule-based systems or manually crafted expert systems, ML models adapt and optimize their structure based on patterns in the data, improving over time with exposure to more examples [3].

Three foundational ML algorithms are commonly referenced in the theoretical exploration of CHD diagnosis: Logistic Regression (LR), Random Forest (RF), and Support Vector Machine (SVM). Each of these models represents a different school of thought in machine learning theory—ranging from classical statistical inference to modern ensemble learning and kernel-based geometric learning [4].

Logistic Regression, although considered a traditional statistical method, is deeply embedded in the theoretical landscape of machine learning as a fundamental classification technique. It uses a sigmoid function to map input features to a probability between 0 and 1, allowing clear interpretation of the impact of each variable [5]. Its coefficients represent the log-odds of the dependent variable (in this case, CHD presence), making it highly interpretable—an essential trait in clinical applications. However, logistic regression assumes a linear relationship between the independent variables and the log-odds of the outcome, which can

be a limiting factor in medical domains where feature interactions and nonlinear patterns are prevalent [6].

Random Forest introduces the concept of ensemble learning, where multiple weak learners (decision trees) are combined to form a more powerful and accurate model. The theory behind Random Forests lies in the principle of "bagging" (bootstrap aggregating), which reduces variance and avoids overfitting by training each tree on a random subset of the data [7]. In medical applications, where noise and missing values are common, Random Forests demonstrate strong resilience. Moreover, they are capable of modeling complex interactions between variables without requiring prior knowledge of the data structure [8]. Another theoretical advantage of RF is its ability to measure feature importance, giving insights into which clinical factors most influence predictions. This aligns well with the needs of healthcare professionals who seek to understand not only what the model predicts but also why [9].

Support Vector Machine is rooted in statistical learning theory and offers a fundamentally different approach to classification. SVMs aim to find the optimal separating hyperplane that maximizes the margin between data points of different classes [10]. For linearly inseparable data, kernel functions such as the Radial Basis Function (RBF) allow SVMs to project data into higher-dimensional spaces where separation becomes possible. This ability to model nonlinearity without explicitly transforming the input features makes SVM a theoretically attractive option for medical diagnosis [11]. However, SVMs can be computationally intensive and less interpretable, particularly when non-linear kernels are used, which may limit their usability in clinical settings where transparency and speed are critical [12].

The theoretical foundation of ML also includes generalization theory, which focuses on how well a trained model can perform on unseen data. In healthcare, where datasets are often small due to privacy restrictions or rare disease prevalence, overfitting is a major concern. Models like Random Forest and regularized versions of Logistic Regression address this through embedded mechanisms that promote generalizability [13]. Furthermore, interpretability—defined as the ability of a human to understand the reasoning behind a model's decision—is a critical consideration in medical AI. From a theoretical standpoint, there exists a trade-off between model complexity and interpretability, often referred to as the "accuracy-interpretability trade-off" [14]. Balancing this trade-off is key to building ML models that are not only accurate but also trusted by clinicians.

In conclusion, the theoretical foundation of machine learning in healthcare blends mathematical rigor, algorithmic innovation, and domain-specific constraints to create predictive systems capable of enhancing clinical decision-making. By understanding the core principles that govern models like Logistic Regression, Random Forest, and SVM, we

gain valuable insight into how these tools can be ethically, safely, and effectively applied to critical problems such as the early detection of coronary heart disease [15].

3. Comparative Theoretical Framework

From a theoretical standpoint, the three models offer distinct advantages and trade-offs when applied to CHD diagnosis (shown in table 1):

Table 1. Theoretical Comparison of Machine Learning Models for Coronary Heart Disease Diagnosis

Model	Theoretical Strengths	Theoretical Limitations	Clinical Applicability
Logistic Regression	Simplicity, interpretability, statistical grounding	Limited to linear decision boundaries	High (for transparency needs)
Random Forest	Handles nonlinearity, resists overfitting, scalable	Less interpretable, may require tuning	Very High (for complex datasets)
SVM	Excellent for small, high-dimensional data; flexible	Kernel choice critical, black-box nature	Moderate (needs expert handling)

This comparative framework allows stakeholders to theoretically assess model suitability before implementation. Logistic Regression remains a strong baseline for interpretable clinical decisions. Random Forest is ideal when the goal is accuracy and robustness, even at the expense of some transparency. SVM, while powerful, demands computational resources and expert tuning, limiting its general use in smaller clinical setups.

4. Conclusion

This theoretical paper has explored how three fundamental machine learning models—Logistic Regression, Random Forest, and Support Vector Machine—conceptually apply to the problem of Coronary Heart Disease diagnosis. While empirical studies are essential for validation, a theoretical understanding of each model’s strengths, limitations, and applicability provides a foundational basis for future research and development. Machine learning, when properly integrated into healthcare systems, holds transformative potential in improving early diagnosis, enhancing clinical workflows, and ultimately saving lives. Future research should explore hybrid models, explainable AI techniques, and the integration of real-time physiological data streams to further enhance CHD prediction models within intelligent clinical environments.

References

- [1] Obermeyer, Ziad, and Ezekiel J. Emanuel. “Predicting the Future — Big Data, Machine Learning, and Clinical Medicine.” *New England Journal of Medicine*, vol. 375, no. 13, 2016, pp. 1216–1219.
- [2] Detrano, Robert, et al. “International Application of a New Probability Algorithm for the Diagnosis of Coronary Artery Disease.” *The American Journal of Cardiology*, vol. 64, no. 5, 1989, pp. 304–310.
- [3] Topol, Eric J. “High-Performance Medicine: The Convergence of Human and Artificial Intelligence.” *Nature Medicine*, vol. 25, 2019, pp. 44–56.
- [4] Jordan, Michael I., and Tom M. Mitchell. “Machine Learning: Trends, Perspectives, and Prospects.” *Science*, vol. 349, no. 6245, 2015, pp. 255–260.
- [5] Hosmer, David W., et al. *Applied Logistic Regression*. 3rd ed., Wiley, 2013.
- [6] Kuhn, Max, and Kjell Johnson. *Applied Predictive Modeling*. Springer, 2013.
- [7] Breiman, Leo. “Random Forests.” *Machine Learning*, vol. 45, no. 1, 2001, pp. 5–32.
- [8] Deo, Rahul C. “Machine Learning in Medicine.” *Circulation*, vol. 132, no. 20, 2015, pp. 1920–1930.
- [9] Lundberg, Scott M., and Su-In Lee. “A Unified Approach to Interpreting Model Predictions.” *Advances in Neural Information Processing Systems*, 2017.
- [10] Cortes, Corinna, and Vladimir Vapnik. “Support-Vector Networks.” *Machine Learning*, vol. 20, 1995, pp. 273–297.
- [11] Scholkopf, Bernhard, and Alexander J. Smola. *Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond*. MIT Press, 2002.
- [12] Liang, Yulan, et al. “Evaluation and Interpretation of Machine Learning Models for Predicting Type 2 Diabetes: A Clinician’s Perspective.” *NPJ Digital Medicine*, vol. 2, 2019, p. 38.
- [13] Zhang, Yujin, et al. “Model Generalization and Overfitting in Predictive Healthcare Analytics.” *IEEE Transactions on Biomedical Engineering*, vol. 68, no. 1, 2021, pp. 49–60.
- [14] Lipton, Zachary C. “The Mythos of Model Interpretability.” *Communications of the ACM*, vol. 61, no. 10, 2018, pp. 36–43.
- [15] Esteva, Andre, et al. “A Guide to Deep Learning in Healthcare.” *Nature Medicine*, vol. 25, no. 1, 2019, pp. 24–29.
- [16] Munmun, Zakia Sultana, Salma Akter, and Chowdhury Raihan Parvez. “Machine Learning-Based Classification of Coronary Heart Disease: A Comparative Analysis of Logistic Regression,

Random Forest, and Support Vector Machine Models." *Open Access Library Journal* 12.3 (2025): 1-12.

[17] Hasan, Sakib, et al. "Analysis of Machine Learning Models for Stroke Prediction with Emphasis on Hyperparameter Tuning Techniques." *International Symposium on Computational Intelligence and Industrial Applications*. Singapore: Springer Nature Singapore, 2024.