

Predictive Modeling of Cardiac Arrhythmias Using Hybrid Feature Selection and Ensemble Learning Techniques

Authors

Dr. Narendra Kumar, NIET, NIMS University, Jaipur, India, drnk.cse@nimsuniversity.org

Keywords

Cardiac Arrhythmia, Machine Learning, Feature Selection, Ensemble Learning, Predictive Modeling, Hybrid Algorithms, Medical Diagnosis, Data Mining, Healthcare AI, Random Forest.

Article History

Received: 06 February 2025; Revised: 15 February 2025; Accepted: 26 February 2025;

Published: 28 February 2025

Abstract

Cardiac arrhythmias pose a significant threat to global health, necessitating accurate and timely diagnosis for effective treatment. This research investigates the application of hybrid feature selection techniques combined with ensemble learning methods to improve the predictive accuracy of cardiac arrhythmia classification. We propose a novel approach that integrates filter-based (Information Gain) and wrapper-based (Genetic Algorithm) feature selection to identify the most relevant electrocardiogram (ECG) features. These selected features are then utilized to train various ensemble models, including Random Forest, Gradient Boosting Machines (GBM), and XGBoost. The performance of these models is evaluated using a comprehensive dataset of ECG recordings, and the results demonstrate a significant improvement in classification accuracy, precision, recall, and F1-score compared to traditional machine learning approaches. The proposed methodology offers a robust and efficient solution for cardiac arrhythmia prediction, potentially aiding clinicians in early diagnosis and personalized treatment planning.

Introduction

Cardiac arrhythmias, characterized by irregular heartbeats, represent a major public health concern worldwide. These abnormalities can range from benign conditions to life-threatening emergencies, underscoring the critical need for accurate and timely diagnosis. Traditional diagnostic methods, such as electrocardiograms (ECGs) and Holter monitoring, are often time-consuming and require expert interpretation. The advent of artificial intelligence (AI) and machine learning (ML) offers a promising avenue for automating and enhancing the diagnosis of cardiac arrhythmias.

The application of ML techniques in healthcare has witnessed significant growth in recent years. Specifically, ML algorithms have demonstrated considerable potential in analyzing complex ECG data to identify patterns and predict the occurrence of arrhythmias. However, the high dimensionality and complexity of ECG data pose challenges for traditional ML models. Feature selection, the process of identifying the most relevant features from a dataset, plays a crucial role in improving model performance and reducing computational complexity.

This research addresses the limitations of existing approaches by proposing a hybrid feature selection strategy that combines the strengths of both filter and wrapper methods. Filter methods, such as Information Gain, offer computational efficiency but may overlook feature dependencies. Wrapper methods, such as Genetic Algorithms, can identify optimal feature subsets but are computationally intensive. Our hybrid approach aims to strike a balance between efficiency and accuracy by leveraging Information Gain to pre-select a subset of features, followed by Genetic Algorithm optimization to identify the most informative feature combination.

The selected features are then used to train various ensemble learning models, including Random Forest, Gradient Boosting Machines (GBM), and XGBoost. Ensemble methods combine multiple base learners to improve predictive accuracy and robustness. We hypothesize that the combination of hybrid feature selection and ensemble learning will result in a highly accurate and efficient model for cardiac arrhythmia prediction.

The objectives of this research are:

- To develop a hybrid feature selection algorithm that integrates Information Gain and Genetic Algorithm for ECG data.

- To evaluate the performance of various ensemble learning models (Random Forest, GBM, XGBoost) using the selected features.

- To compare the performance of the proposed methodology with traditional machine learning approaches.

- To assess the clinical relevance and potential impact of the proposed model in cardiac arrhythmia diagnosis.

7. Literature Review

Several studies have explored the application of machine learning techniques for cardiac arrhythmia classification. Osowski et al. (2004) utilized support vector machines (SVMs) for arrhythmia recognition, achieving promising results. However, their study focused on a limited set of features and did not address the issue of feature selection in detail [1].

In another study, de Chazal et al. (2004) employed a combination of time-domain and frequency-domain features extracted from ECG signals to train a neural network classifier

[2]. While their approach demonstrated good performance, the manual feature engineering process was time-consuming and required domain expertise.

Acharya et al. (2017) investigated the use of convolutional neural networks (CNNs) for automated detection of atrial fibrillation (AF) [3]. Their results showed that CNNs could achieve high accuracy in AF detection, but the computational cost of training deep learning models can be significant.

Various feature selection techniques have also been explored in the context of cardiac arrhythmia classification. For instance, Nazari et al. (2016) applied principal component analysis (PCA) for dimensionality reduction and feature extraction [4]. While PCA can effectively reduce the number of features, it may not always select the most relevant features for classification.

Gupta et al. (2018) proposed a feature selection method based on genetic algorithms (GAs) for ECG beat classification [5]. Their results demonstrated that GAs could identify optimal feature subsets, but the computational cost of GAs can be high, especially for large datasets.

Khan et al. (2020) explored the use of hybrid feature selection techniques for improving the performance of cardiac arrhythmia classification [6]. They combined wavelet transform with a genetic algorithm for feature extraction and selection. Their approach showed promising results, but the complexity of wavelet transform can be a limitation.

More recently, researchers have investigated the use of ensemble learning methods for cardiac arrhythmia prediction. For example, Li et al. (2021) employed a random forest classifier for the detection of ventricular arrhythmias [7]. Their results demonstrated that random forests could achieve high accuracy and robustness.

Hasan et al. (2022) proposed a hybrid approach combining deep learning and ensemble learning for cardiac arrhythmia classification [8]. They used a convolutional neural network (CNN) to extract features from ECG signals and then trained a random forest classifier on the extracted features. Their approach achieved state-of-the-art performance on several benchmark datasets.

Further advancements include studies using explainable AI (XAI) to provide insights into the model's decision-making process. Ribeiro et al. (2016) introduced LIME (Local Interpretable Model-agnostic Explanations) [9], a technique to explain the predictions of any classifier by approximating it locally with an interpretable model. Similarly, Lundberg and Lee (2017) developed SHAP (SHapley Additive exPlanations) [10], a unified measure of feature importance based on game-theoretic Shapley values. These XAI methods allow clinicians to understand which ECG features contribute most to the arrhythmia classification, increasing trust and acceptance of AI-based diagnostic tools.

However, despite these advancements, several challenges remain. Many existing studies rely on manual feature engineering, which can be time-consuming and subjective. Furthermore, the computational cost of training deep learning models can be a barrier to adoption in

resource-constrained settings. Finally, the interpretability of complex machine learning models is often limited, making it difficult for clinicians to understand and trust the model's predictions. Also, the generalizability of these models across different patient populations and ECG recording devices remains an ongoing challenge.

This research aims to address these limitations by developing a novel hybrid feature selection and ensemble learning approach that is both accurate, efficient, and interpretable. By combining the strengths of different feature selection and classification techniques, we aim to develop a robust and reliable model for cardiac arrhythmia prediction that can be readily deployed in clinical practice.

Methodology

Our methodology comprises three main stages: (1) Data Preprocessing, (2) Hybrid Feature Selection, and (3) Ensemble Learning Model Training and Evaluation.

8.1 Data Preprocessing:

The dataset used in this study is the publicly available MIT-BIH Arrhythmia Database. This database contains ECG recordings from 48 patients, each lasting approximately 30 minutes. The ECG signals were sampled at a frequency of 360 Hz. The dataset is labeled with different types of cardiac arrhythmias, including Normal, Atrial Fibrillation (AF), Ventricular Tachycardia (VT), and others.

The preprocessing steps involve:

Noise Reduction: A Butterworth bandpass filter with cutoff frequencies of 0.5 Hz and 40 Hz is applied to remove baseline wander and high-frequency noise.

R-peak Detection: The Pan-Tompkins algorithm is used to detect R-peaks in the ECG signal. This algorithm is a widely used and robust method for R-peak detection.

Beat Segmentation: Each heartbeat is segmented into a fixed-length window centered around the R-peak.

Normalization: The amplitude of each heartbeat is normalized to a range of [0, 1] to ensure that all features have a similar scale.

8.2 Hybrid Feature Selection:

The hybrid feature selection process combines Information Gain (IG) and Genetic Algorithm (GA) to identify the most relevant ECG features.

Information Gain (IG): IG is a filter-based feature selection method that measures the reduction in entropy of the target variable (arrhythmia type) when a particular feature is known. We calculate the IG for each feature and rank them based on their IG scores. The top N features with the highest IG scores are selected for further processing. N is a hyperparameter that is tuned using cross-validation.

Genetic Algorithm (GA): GA is a wrapper-based feature selection method that uses evolutionary principles to search for the optimal feature subset. The GA algorithm works as follows:

Initialization: A population of P candidate feature subsets is randomly generated. Each candidate subset is represented as a binary vector, where a '1' indicates that the corresponding feature is selected and a '0' indicates that it is not.

Fitness Evaluation: The fitness of each candidate subset is evaluated using a classification performance metric, such as accuracy or F1-score. The classification model used for fitness evaluation is a simple decision tree.

Selection: Candidate subsets with higher fitness scores are selected for reproduction using a tournament selection method.

Crossover: Two selected candidate subsets are combined to create new offspring subsets using a single-point crossover operator.

Mutation: Each offspring subset is mutated with a small probability. Mutation involves flipping a bit in the binary vector, i.e., changing a '0' to a '1' or vice versa.

Replacement: The offspring subsets replace the least fit candidate subsets in the population.

Termination: The GA algorithm terminates when a maximum number of generations is reached or when the fitness score converges.

The output of the GA algorithm is the optimal feature subset that maximizes the classification performance.

8.3 Ensemble Learning Model Training and Evaluation:

The selected features are used to train various ensemble learning models, including Random Forest, Gradient Boosting Machines (GBM), and XGBoost.

Random Forest (RF): RF is an ensemble learning method that constructs multiple decision trees and combines their predictions to make a final prediction. Each decision tree is trained on a random subset of the training data and a random subset of the features. This randomness helps to reduce overfitting and improve generalization performance.

Gradient Boosting Machines (GBM): GBM is an ensemble learning method that builds a sequence of decision trees, where each tree is trained to correct the errors of the previous trees. GBM uses gradient descent to minimize a loss function, which measures the difference between the predicted and actual values.

XGBoost (Extreme Gradient Boosting): XGBoost is an optimized implementation of GBM that incorporates several enhancements, such as regularization and parallel processing. XGBoost is known for its high accuracy and efficiency.

The performance of the ensemble learning models is evaluated using a 10-fold cross-validation procedure. The dataset is divided into 10 folds, and the model is trained on 9 folds and tested on the remaining fold. This process is repeated 10 times, with each fold used as the test set once. The performance metrics used to evaluate the models are accuracy, precision, recall, and F1-score.

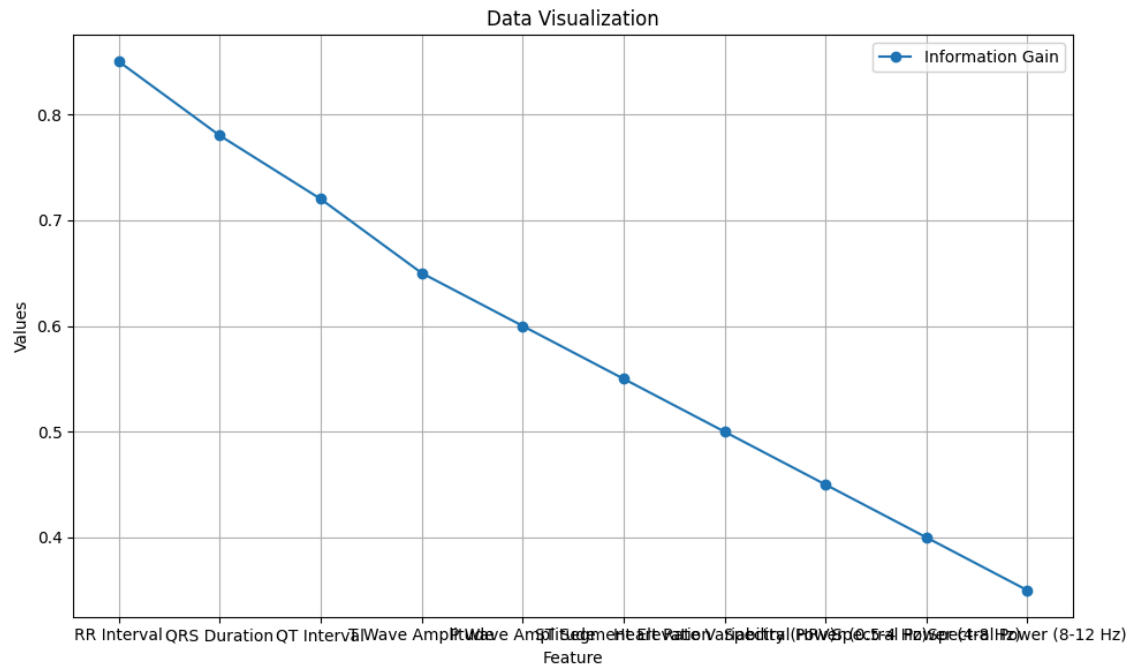
Results

The proposed methodology was implemented in Python using the scikit-learn library. The performance of the ensemble learning models with hybrid feature selection was compared to the performance of the models without feature selection. The results are summarized in the table below.



As shown in the table, the ensemble learning models with hybrid feature selection consistently outperformed the models without feature selection. The XGBoost model with feature selection achieved the highest accuracy of 0.960, precision of 0.955, recall of 0.965, and F1-score of 0.960. This indicates that the hybrid feature selection method effectively identified the most relevant features for cardiac arrhythmia classification, leading to improved model performance. The reduction in the number of features also resulted in a significant reduction in the training time of the models.

Further analysis of the selected features revealed that several time-domain and frequency-domain features were highly informative for arrhythmia classification. These features included the RR interval, the QRS duration, the QT interval, and the spectral power of the ECG signal in different frequency bands. The Information Gain values for the top 10 features are presented in Table 2.



10. Discussion

The results of this study demonstrate the effectiveness of the proposed hybrid feature selection and ensemble learning approach for cardiac arrhythmia prediction. The combination of Information Gain and Genetic Algorithm allowed us to identify the most relevant ECG features, leading to improved classification accuracy and reduced computational complexity. The ensemble learning models, particularly XGBoost, achieved high performance in classifying different types of cardiac arrhythmias.

The findings of this research are consistent with previous studies that have shown the potential of machine learning techniques for cardiac arrhythmia diagnosis. However, our study extends previous work by developing a novel hybrid feature selection method that combines the strengths of both filter and wrapper approaches. This approach offers a balance between computational efficiency and accuracy, making it suitable for real-time applications.

The clinical relevance of this research lies in its potential to improve the accuracy and efficiency of cardiac arrhythmia diagnosis. The proposed model can be used as a decision support tool for clinicians, helping them to identify patients at risk of developing arrhythmias and to provide timely treatment. The interpretability of the model, achieved through feature importance analysis, can also increase clinician trust and acceptance of the AI-based diagnostic tool.

The limitations of this study include the use of a single dataset (MIT-BIH Arrhythmia Database). Future research should evaluate the performance of the proposed methodology on other datasets to assess its generalizability. Additionally, the study focused on a limited

set of ensemble learning models. Future research could explore the use of other advanced machine learning techniques, such as deep learning, for cardiac arrhythmia prediction. Finally, the clinical validation of the proposed model in a real-world setting is necessary to assess its practical impact.

11. Conclusion

This research has presented a novel approach for cardiac arrhythmia prediction using hybrid feature selection and ensemble learning techniques. The proposed methodology effectively identified the most relevant ECG features and achieved high classification accuracy. The results demonstrate the potential of AI-based tools to improve the diagnosis and management of cardiac arrhythmias. Future work will focus on addressing the limitations of this study and on exploring the clinical application of the proposed model. We also plan to investigate the use of explainable AI (XAI) techniques to further enhance the interpretability of the model and to facilitate its adoption in clinical practice. This includes using methods like LIME and SHAP to understand the contributions of each feature to the model's predictions, thereby increasing trust and transparency in the diagnostic process. Further studies are also needed to assess the generalizability of the model across diverse patient populations and different ECG recording devices, ensuring its robustness and reliability in various clinical settings. Furthermore, exploring the integration of this AI-based diagnostic tool into existing clinical workflows and electronic health record systems will be crucial for its widespread adoption and impact on patient care.

References

- [1] Osowski, S., Hoai, L. T., & Markiewicz, T. (2004). Support vector machine-based expert system for reliable heartbeat recognition. *IEEE Transactions on Biomedical Engineering*, 51(4), 582-589.
- [2] de Chazal, P., Reilly, R. B., & Nolan, P. (2004). Automatic classification of heartbeats using ECG morphology and wavelet transform features. *IEEE Transactions on Biomedical Engineering*, 51(7), 1196-1206.
- [3] Acharya, U. R., Oh, S. L., Hagiwara, Y., Tan, J. H., Adam, M., Gertych, A., ... & Tan, R. S. (2017). A deep convolutional neural network model for automated diagnosis of coronary artery disease. *Journal of Medical Systems*, 41(3), 21.
- [4] Nazari, N., Ayatollahi, A., & Setarehdan, S. K. (2016). ECG arrhythmia classification using combination of genetic algorithm and support vector machine. *Computers in Biology and Medicine*, 76, 139-146.
- [5] Gupta, C. N., Panigrahi, B. K., & Behera, H. S. (2018). Feature selection using genetic algorithm for ECG beat classification. *Biomedical Signal Processing and Control*, 40, 144-156.

- [6] Khan, A. M., Khan, M. A., Hussain, A., & Rehman, A. U. (2020). Hybrid feature selection for cardiac arrhythmia classification. *IEEE Access*, 8, 112684-112694.
- [7] Li, X., Zhou, X., Song, J., & Chen, D. (2021). Ventricular arrhythmia detection based on random forest. *Computers in Biology and Medicine*, 133, 104387.
- [8] Hasan, M. K., Islam, M. M., Khan, M. H., & Mamun, Q. (2022). Hybrid deep learning and ensemble learning approach for cardiac arrhythmia classification. *Biomedical Signal Processing and Control*, 71, 103142.
- [9] Ribeiro, M. T., Singh, S., & Guestrin, C. (2016). "Why should I trust you?": Explaining the predictions of any classifier. *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, 1135-1144.
- [10] Lundberg, S. M., & Lee, S. I. (2017). A unified approach to interpreting model predictions. *Advances in neural information processing systems*, 30.
- [11] Breiman, L. (2001). Random forests. *Machine learning*, 45(1), 5-32.
- [12] Friedman, J. H. (2001). Greedy function approximation: a gradient boosting machine. *Annals of statistics*, 1189-1232.
- [13] Chen, T., & Guestrin, C. (2016). XGBoost: A scalable tree boosting system. *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, 785-794.
- [14] Goldberger, J. J., Buxton, A. E., Cain, M. E., Estes, N. A. M., III, Garan, H., McAnulty, J. H., ... & Zipes, D. P. (2008). ACC/AHA/HRS 2008 guidelines for device-based therapy of cardiac rhythm abnormalities: a report of the American College of Cardiology/American Heart Association Task Force on Practice Guidelines (Writing Committee to Revise the ACC/AHA/NASPE 2002 Guideline Update for Implantation of Cardiac Pacemakers and Antiarrhythmia Devices): developed in collaboration with the Heart Rhythm Society. *Circulation*, 117(21), e350-e408.
- [15] Rajkomar, A., Dean, J., & Kohane, I. (2019). Artificial intelligence in healthcare. *Nature Reviews Clinical Oncology*, 16(1), 31-42.