

A Comprehensive Review on Arrhythmia Disease Prediction Using Data Mining Classification Algorithms and Feature Selection Techniques

Sunil Kumar Saini¹, Ms. Pragya Bharti²

¹M. Tech Scholar, ²Assistant Professor

Department of Computer Science & Engineering

Rajasthan Institute of Engineering & Technology, Jaipur, Rajasthan

Abstract: The state arrhythmia, a state of irregular heartbeat, is a significant health problem that can lead to serious complications including sudden cardiac arrest. The exact and initial prediction of arrhythmia is important for effective medical intervention and patient care. With progress in machine learning and data mining, the classification algorithm has become an essential tool for detecting arrhythmia from the electrocardiogram signals (ECG). However, the high height of the ECG dataset creates a challenge, making the facility's electoral technique to improve model performance, reduce calculation complexity and increase the interpretation. This article offers various data mining classification algorithms such as decision trees, Support Vector Machines (SVM), k-Nearest Neighbors (k-NN), Naïve Bayes, Random Forest, and Deep Learning models employed for arrhythmia prediction. In addition, we discover different functional choice techniques, including filters, cover and built-in methods, to identify the most relevant properties that

contribute to better classification accuracy. We also discuss hybrid approaches that integrate many classify and facilitate selection strategies to achieve increased prediction benefit. In addition, this review addresses challenges such as unbalanced data sets, noise ECG signals and model lecturers, highlighting recent progress and potential future research directions. The findings from this study serve as a valuable resource for researchers and health professionals, who aim to develop a more strong and effective arrhythmia model.

Keywords: Arrhythmia Prediction, Data Mining, Classification Algorithms, Feature Selection, Machine Learning, Electrocardiogram (ECG).

I. INTRODUCTION

In essence, data mining aims to uncover hidden patterns and trends in large databases and facilitates automatic data exploration. From these patterns, some rules are derived that allow users to review and examine their business or scientific decisions, leading to more interaction with d

atabases and data warehouses. The term "data mining" refers to the process of extracting or "mining" knowledge patterns from vast amounts of data; the term "knowledge mining," which is shorter, may not reflect the emphasis on mining from large amounts of data.

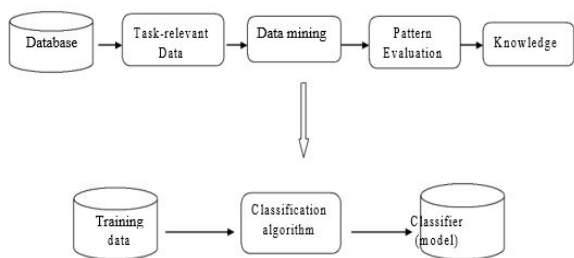


Figure 1: Knowledge Discovery

Large amounts of data are accessible in practically every sector in the modern world. In order to address the demands, researchers are trying to manage, evaluate, and make decisions. The term "data mining" refers to the process of extracting useful information from similarities. It forecasts the patterns and actions of the variables found in the data. One such vast subject that is best suited for examining the extent of mining is medical data. One such condition brought on by an irregular heartbeat is arrhythmia. Dealing with the vast amounts of data that are available presents numerous challenges. The step that cuts down on redundancy and even complexity is preprocessing.

High-dimensional data is reduced by feature selection, which keeps the majority of the information or traits required for categorization. It is used on an arrhythmia

dataset to eliminate characteristics that are not important for analysis and identify the variables that are. One way to describe the scope of feature selection is either local or global. Local uses a different collection of characteristics every time, but global uses the same set for all categories. Examining different feature selection methods for arrhythmia datasets is the aim of this thesis. To obtain the most suitable outcomes, it will also combine feature selection and classification techniques.

II. ARRHYTHMIA DISEASE AND ECG DATA

Arrhythmias encompass a range of irregular heartbeats, such as bradycardia (slow heartbeat), atrial fibrillation (irregular heartbeat), ventricular fibrillation (chaotic heartbeat), and tachycardia (fast heart rate). Severe heart failure, including stroke, heart failure, and sudden cardiac arrest, can result from these conditions. An electrocardiogram (ECG), a non-invasive test that tracks the electrical activity of the heart over time, is the main tool used to identify arrhythmia. By recording the wave form pattern, the ECG signal gives doctors important information about heart disease and aids in the diagnosis and categorization of various arrhythmias. Analyzing various signal characteristics, such as P waves, QRS complexes, and T-waves, which correlate to distinct cardiac cycle phases, is part of the process of interpreting

ECG data. Manual ECG analysis, however, can be laborious, subjective, and lead to inaccurate diagnoses. The use of data mining and machine learning techniques to identify automated arrhythmias has garnered a lot of attention as a solution to these problems. Researchers have increased the precision and efficacy of arrhythmia prediction by utilizing sophisticated classification algorithms and comfort selection techniques. Combining computation methods enables the extraction of significant features from ECG data, improving classification performance and lowering noise. Furthermore, real-time portable ECG equipment and surveillance have made it easier to identify patients early and follow them continuously, which has improved clinical outcomes.

III. DECISION TREE AND DATA PROCESSING

A. Decision Tree

The classification data is a basic component of mining, which helps identify general properties in a dataset, classify them and highlight meaningful patterns. Decision trees act as a powerful decision -making apparatus, offering a structured approach to achieve goals. A decision represents a characteristic under each internal node assessment in the tree, while the branches reflect the results of these evaluation. This method naturally evaluates convenience choice and characteristic. To maintain the balance of the

tree, dynamic pruning technology uses breeding trees, and ensures optimal height on the basis of priority control on each node, which includes the concept of node merger. Decision is one of the benefits of trees that they require minimal preparatory efforts for data preparation. To increase accuracy and reduce the depth of the tree, a limited number of properties - usually one or two - are selected as a divided criteria, and maximizes the information result. This approach also improves the possibility of finding optimal solutions globally. The decision trees are highly adaptable, which is able to handle different data types including nominal, numerical and text data. They can also manage the dataset with errors or missing values. His application in the health care system has received significant traction due to large amounts of medical data available. Decisions help trips to come up with predictions, diagnosis and treat treatment, contribute to more effective patient care.

B. Data Preprocessing

Because of its vast size and many sources, the data found in today's databases is especially susceptible to issues including missing values, inconsistent data, and noise. The quality of the data acquired must be improved in order to raise the caliber of the mining results. There are numerous methods that can be used to solve these issues. They are as follows:

- **Data Cleaning:** This technique is used to eliminate noise and inconsistencies from the data.
- **Data integration:** Since the data is gathered from many sources, it must be combined and stored in a single location known as a data warehouse.
- **Data transformation:** Remove any redundant information from the data and normalize it in accordance with analysis requirements to improve performance analysis accuracy.
- **Data reduction:** By using aggregation and removing characteristics that exhibit duplicate behavior, it is possible to clearly minimize the data.

IV. FEATURE SELECTION AND FILTERS METHODS

A. Feature Selection

Prior to using any mining approach, preprocessing is a crucial step. The process of choosing a subset of features is known as feature selection. This is not the same as feature extraction, which creates new features based on the functions of existing features. In order to identify feature subsets, the assessment method first applies training data. Variance, correlations, standard deviation, eigenvalues, and other statistical measures are used to evaluate the variables. Next, as seen in fig. 3.1, the information content is passed to the algorithm to be applied to the values that are acquired. Until the appropriate subset of

variables is chosen from the large dataset, this cycle is repeated. The chosen characteristics are now prepared for additional data processing and analysis.

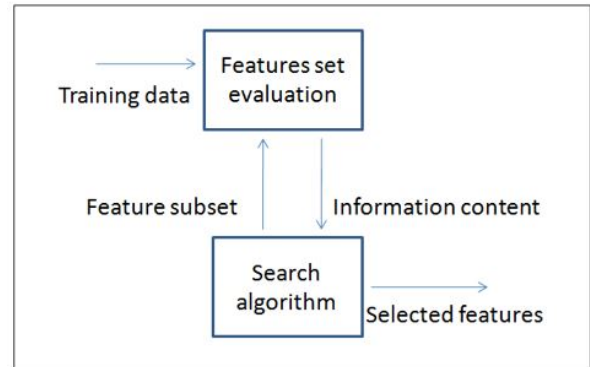


Figure 2: General Approach for Feature Selection

The following benefits are achieved through feature selection, which plays a very effective role in dimensional reduction in retaining the information content of data:

- Better model understanding and visualization—visualizing with fewer features will result in a more clear understanding of the model;
- Generalization of the model—overfitting reduction leads to more accurate learning;
- Efficiency in terms of time and space complexity—both gets reduced during execution;
- Avoid the curse of dimensionality;
- Shorten training time;
- Improved generalization by reducing overfitting;

- Avoid the curse of dimensionality;
- Provides good accuracy results.

B. Filter methods



Figure 3: General Procedure for Filter Method

Typically, filtering techniques are applied as a preprocessing step. No machine learning methods influence the feature selection process. Rather, characteristics are chosen based on how well they correlate with the end variable as determined by scores in a variety of statistical tests. Here, the term "correlation" is arbitrary. LDA, ANOVA, chi-square, and PCA are among them.

A target attribute is selected in order to employ filter-based feature selection. It predicts the subset's maximum power using statistical metrics. To obtain the needed or pertinent qualities, a module calculates a score value. The property with the highest relevance is chosen as the significant subset.

V. RESULTS AND CONCLUSION

By using data mining classification algorithms and techniques for the selection of comforts, the prediction of arrhythmia has improved significantly this year. Unbalanced data and model lecturers have paved the way for more accurate and effective arrhythmia diagnosis despite challenges such as progression, deep

education and surveillance of real -time in hybrid models. This review serves as a comprehensive resource for researchers and doctors for the health care system, who aims to detect arrhythmia using data mining and machine learning.

VI. FUTURE SCOPE

Future work involves many more experiments on analysis of data from other areas. Research based on their impact on the hybrid approach and classification of functional choice techniques. Handling lack of values using different techniques. Search strategies using built -in functional choice techniques. Integration of other sources of information, different combinations of functional choice methods and fuzzy set principles to improve model development. Analysis of results with unsafe technology. Hybrid models can also be used to improve and achieve better results. Some other future prospectus are given below;

- The hybrid model combination of multiple classifies and functional choice techniques can improve predicting accuracy.
- Explanation of AI- Development of interpretable models ensures reliability and clinical acceptance.
- ECG monitoring- Integrated portable laptop equipment with machine learning models provides for detecting real-time arrhythmia.

REFERENCES

- [1] Bhukya D. and Ramachandram S., “Decision Tree Induction: An Approach for Data Classification Using AVL-Tree”, International Journal of Computer and Electrical Engineering, Vol. 2, No. 4, August, 2010 1793-8163.
- [2] Madadipouya K., “A New Decision Tree Method For Data Mining In Medicine”, Advanced Computational Intelligence: An International Journal (ACII), Vol.2, No.3, July 2015.
- [3] Teli S., Kanikar P., “ A Survey on Decision Tree Based Approaches in Data Mining”, International Journal of Advanced Research in Computer Science and Software Engineering, Volume 5, Issue 4, 2015.
- [4] Țăranu I., “Data mining in healthcare: decision making and precision”, Database Systems Journal vol. VI, no. 4/2015.
- [5] Dey M., Rautaray S., “Study and Analysis of Data mining Algorithms for Healthcare Decision Support System”, International Journal of Computer Science and Information Technologies, Vol. 5 (1) , 2014, 470-477.
- [6] Goel V., “Decision analysis: applications and limitations”, CAN MED ASSOC J1992.
- [7] Werner E., Wheeler S. and Burd I., “Creating Decision Trees to Assess Cost-Effectiveness in Clinical Research”, J Biomet Biostat 2012.
- [8] Elwyn G., Edwards A., Eccles M., Rovner D., “Decision Analysis In Patient Care”, The Lancet , Vol 358 , August 18, 2001.
- [9] Soleimanian F., Mohammadi P., Hakimi P., “Application of Decision Tree Algorithm for Data Mining in Healthcare Operations: A Case Study”, International Journal of Computer Applications (0975 – 8887) Volume 52 – No. 6, August 2012.
- [10] Quinlan J., “Induction Of Decision Trees”, Machine Learning 1: 81-106, 1986.
- [11] Szolovits P., “Uncertainty and Decisions in Medical Informatics”, Methods of Information in Medicine, 34:111–21, 1995.
- [12] Teli S., KanikarP. ,“ A Survey on Decision Tree Based Approaches in Data Mining”, International Journal of Advanced Research in Computer Science and Software Engineering, Volume 5, Issue 4, 2015.
- [13] Srimani P.K. And Koti M.S, “Evaluation Of Principal Components Analysis (Pca) And Data Clustering Techniques (DCT) On Medical Data”, International Journal Of Knowledge Engineering, ISSN: 0976-5816, Volume 3, Issue 2, 2012. 2016.
- [14] Wimmer H., Powell L., “Principle Component Analysis for Feature Reduction and Data Preprocessing in Data Science”, 2016 Proceedings of the

- Conference on Information Systems Applied Research Las Vegas, Nevada USA.
- [15] Sjostrand K., Stegmann M. and Larsen R., “Sparse Principal Component Analysis in Medical Shape Modeling”.
- [16] Jayaprada S., “Enhanced C-Means Clustering with PCA for medical dataset”, IJDCST, April-May-2016.
- [17] Luukka P., “A New Nonlinear Fuzzy Robust PCA Algorithm and Similarity Classifier in Classification of Medical Data Sets”, International Journal of Fuzzy Systems, Vol. 13, No. 3, September 2011.
- [18] Kalaiselvi.R, Premadevi.P and Hamsathvani.M, “Weighted Principle Component Analysis for Dimensionality Reduction in Medical Dataset”, IJISSET - International Journal of Innovative Science, Engineering & Technology, Vol. 2 Issue 5, May 2015.
- [19] Salama G., Abdelhalim M., and Zeid M., “Breast Cancer Diagnosis on Three Different Datasets Using Multi-Classifiers”, International Journal of Computer and Information Technology (2277 – 0764) Volume 01– Issue 01, September 2012.
- [20] Naik G., Selvan S., Gobbo M., Acharyya A., Nguyen H., “Principal Component Analysis Applied To Surface Electromyography: A Comprehensive Review”, IEEE Access, VOLUME 4, 2016.
- [21] Hu B., Dai Y., Su Y., Moore P., Zhang X., Mao C., Chen J., Xu L., “Feature Selection for Optimized High dimensional Biomedical Data Using An Improved Shuffled Frog Leaping Algorithm”, 2016 Ieee.
- [22] Williams B., Onsman A., Brown T., “Exploratory factor analysis: A five-step guide for novices”, Journal of Emergency Primary Health Care (JEPHC), Vol. 8, Issue 3, 2010.
- [23] Pinto C., “Data Reduction I: PCA and Factor Analysis”, Data Analysis Seminars 11 November 2009.
- [24] Anna B. and Jason W., “Best Practices in Exploratory Factor Analysis: Four Recommendations for Getting the Most From Your Analysis”, Practical Assessment, Research & Evaluation, Volume 10 Number 7, July 2005.