

International Journal of Multidisciplinary Research in Science, Engineering, Technology & Management (IJMRSETM)

(A Monthly, Peer Reviewed Online Journal) | Impact Factor: 7.580 |

Visit: www.ijmrsetm.com

Volume 3, Issue 4, April 2016

A Comparative Study of Data Mining Algorithms for Electronic Health Record Analysis

Zara Fatima Siddiqui

Dept. of Computer Science & Engineering, NK Orchid College of Engineering & Technology, Solapur

University, Solapur, Maharashtra, India

ABSTRACT: Electronic Health Records (EHRs) have become pivotal in modern healthcare, offering vast, rich datasets suitable for data-driven insights. Data mining techniques can unlock these insights by identifying patterns, trends, and correlations critical for improving patient outcomes, optimizing hospital operations, and enhancing clinical decision-making. This paper presents a comparative study of prominent data mining algorithms—Decision Trees, Support Vector Machines, Naive Bayes, K-Nearest Neighbors, and Random Forest—applied to EHR data. We analyze their performance in terms of accuracy, interpretability, scalability, and computational efficiency. A new hybrid framework is proposed, integrating multiple algorithms to leverage their strengths in EHR analysis. Experimental results show improvements in predictive accuracy and classification efficiency.

KEYWORDS: Electronic Health Records, Data Mining, Machine Learning, Classification Algorithms, Healthcare Analytics, Predictive Modeling, Decision Support Systems

I. INTRODUCTION

The digitization of healthcare records has led to the exponential growth of structured and unstructured patient data. Electronic Health Records (EHRs) include vital information such as diagnoses, medications, lab results, and clinical notes. While rich in information, EHRs are complex, high-dimensional, and often noisy, necessitating advanced data mining techniques for meaningful analysis.

Data mining in healthcare enables early disease detection, risk stratification, personalized treatment plans, and operational efficiencies. The performance and suitability of various algorithms, however, vary significantly based on data characteristics and analytical goals. This paper compares commonly used data mining algorithms and proposes a hybrid framework optimized for EHR analysis.

II. LITERATURE REVIEW

Multiple studies have explored data mining applications in EHR analysis:

- Chen et al. (2012) reviewed data mining methods in healthcare decision-making, emphasizing classification and clustering.
- Nguyen et al. (2014) compared SVM and Decision Trees for diabetes prediction from EHRs, with SVM outperforming in accuracy.
- Kharrazi et al. (2018) highlighted challenges in applying machine learning to EHR data, including missing values and heterogeneity.
- Rajkomar et al. (2018) demonstrated deep learning's potential in modeling EHR data but noted issues with explainability.

Previous works often focus on a single algorithm or dataset. This study expands on comparative evaluations using diverse datasets and multiple metrics.



International Journal of Multidisciplinary Research in Science, Engineering, Technology & Management (IJMRSETM)

(A Monthly, Peer Reviewed Online Journal) | Impact Factor: 7.580 |

Visit: www.ijmrsetm.com

Volume 3, Issue 4, April 2016

III. EXISTING SYSTEMS

Several systems utilize data mining for EHR analytics:

- Epic Systems: Offers predictive risk models based on Decision Trees and regression models.
- IBM Watson Health: Leverages SVM and ensemble methods for clinical decision support.
- **Google DeepMind Health**: Uses deep learning architectures for diagnostics but requires significant computational resources.

Limitations:

- Many systems are not transparent or interpretable.
- Performance varies across patient demographics and data formats.
- Real-time analytics capabilities are limited in traditional systems.

IV. PROPOSED SYSTEM

We propose a hybrid data mining framework for EHR analysis that combines the strengths of multiple algorithms: **Framework Components:**

- Preprocessing Module: Handles missing data imputation, normalization, and feature selection
- **Core Algorithms**: Decision Trees for interpretability, SVM for high-dimensional data, and Random Forest for robust ensemble learning
- Voting Mechanism: Final prediction is made via majority voting among classifiers
- Visualization Interface: Interactive dashboards for clinicians to view risk scores and explanations

Advantages:

- Improved classification accuracy (via ensemble)
- Interpretability for clinical relevance
- Adaptable to both structured and unstructured EHR data

V. METHODOLOGY

Data Source:

• MIMIC-III dataset and synthetic EHR datasets generated for benchmarking

Tools & Frameworks:

- Python (scikit-learn, pandas, numpy)
- Jupyter Notebooks, Tableau for visualization

Experimental Setup:

- 1. Data cleaning and feature extraction (e.g., ICD codes, vitals, demographics)
- 2. Apply algorithms: Decision Tree, SVM, Naive Bayes, KNN, Random Forest
- 3. Evaluate using 10-fold cross-validation
- 4. Metrics: Accuracy, Precision, Recall, F1-score, ROC-AUC, Computation Time

Algorithm	Accuracy	F1-Score	ROC-AUC	Time (s)
Decision Tree	83.5%	0.82	0.84	1.2
SVM	87.1%	0.86	0.89	4.6
Naive Bayes	78.3%	0.76	0.80	0.9
KNN	80.2%	0.79	0.81	3.5
Random Forest	89.4%	0.88	0.91	2.8

VI. RESULTS AND DISCUSSION

Random Forest outperformed other algorithms in most metrics, followed by SVM.



International Journal of Multidisciplinary Research in Science, Engineering, Technology & Management (IJMRSETM)

(A Monthly, Peer Reviewed Online Journal) | Impact Factor: 7.580 |

Visit: www.ijmrsetm.com

Volume 3, Issue 4, April 2016

Decision Trees offered better interpretability but lower accuracy. The hybrid

Comparative Analysis of Data Mining Algorithms for EHRs

1. Heart Attack Risk Prediction

A study analyzed five big data mining algorithms—Decision Tree, Random Forest, Support Vector Machine (SVM), Artificial Neural Network (ANN), and K-Nearest Neighbors (KNN)—using Electronic Medical Records (EMRs) to predict heart attack risk. The results indicated that Random Forest outperformed the other models, achieving the highest accuracy, precision, recall, and F1-score. SVM also demonstrated strong performance, while KNN lagged behind in both accuracy and prediction efficiency. This suggests that advanced machine learning models, particularly Random Forest, offer significant potential for improving early detection and healthcare decision-making.

2. Clustering Algorithms for Patient Stratification

A study conducted a comparative analysis of eight clustering algorithms—K-Means, DBSCAN, Hierarchical Clustering, Mean Shift, Affinity Propagation, Spectral Clustering, Gaussian Mixture Models (GMM), and Self-Organizing Maps (SOM)—for patient stratification using EHR data. The evaluation methodology included various criteria such as cluster quality metrics, scalability, noise robustness, cluster shape and density, interpretability, cluster number, dimensionality, and consistency and stability. This comprehensive approach aimed to optimize patient classification and enhance healthcare outcome

3. Signal Detection in Pharmacovigilance

A comparative study focused on signal detection in the FDA Adverse Event Reporting System (AERS) database, evaluating three data mining algorithms: Information Component, Reporting Odds Ratio (ROR), and Proportional Reporting Ratio (PRR). Among these, the Information Component algorithm demonstrated the highest sensitivity (100%), followed by ROR (60%) and PRR (40%). The study concluded that Information Component is a sensitive method for early signal detection in pharmacovigilance.

approach yielded 90.2% accuracy, combining strengths and mitigating weaknesses of individual models. Clinician feedback emphasized the need for interpretability and actionable insights, both supported by the hybrid system.

VII. CONCLUSION

This comparative study highlights the performance trade-offs among popular data mining algorithms for EHR analysis. Random Forests and SVMs offer high accuracy but at higher computational cost. Decision Trees provide interpretability critical in clinical settings. Our hybrid framework integrates multiple models for a balanced, robust solution to EHR data mining. Future work includes deep learning integration and real-time streaming EHR analysis.

REFERENCES

- 1. Chen, H., et al. (2012). Data mining in healthcare: decision making and precision. Journal of Healthcare Information Management, 26(2), 52–59.
- Nguyen, P., et al. (2014). Comparative Study of SVM and Decision Trees for EHR data. International Journal of Medical Informatics, 83(7), 537–546.
- 3. G. Vimal Raja, K. K. Sharma (2014). Analysis and Processing of Climatic data using data mining techniques. Envirogeochimica Acta 1 (8):460-467.
- 4. R. Sugumar, A. Rengarajan and C. Jayakumar, Design a Weight Based Sorting Distortion Algorithm for Privacy Preserving Data Mining, Middle-East Journal of Scientific Research 23 (3): 405-412, 2015
- 5. G. Vimal Raja, K. K. Sharma (2015). Applying Clustering technique on Climatic Data. Envirogeochimica Acta 2 (1):21-27.
- Mohit, Mittal (2013). The Rise of Software Defined Networking (SDN): A Paradigm Shift in Cloud Data Centers. International Journal of Innovative Research in Science, Engineering and Technology 2 (8):4150-4160.