

Hybrid RF-DT Model for Chronic Disease Detection Using EHR Big Data Management and Analytics

Sharifah Masoud Al Jaleed^{1*}, Ghadeer Ibrahim Alkoblan²

¹Department of Computer Science, College of Computer and Information Sciences
Al Imam Mohammad Ibn Saud Islamic University (IMSIU), Riyadh 11432, Saudi Arabia
Email:sh.aljaleed@gmail.com*

²Department of Computer Science, College of Computer and Information Sciences
Al Imam Mohammad Ibn Saud Islamic University (IMSIU), Riyadh 11432, Saudi Arabia
Email:ghadeer.alkoblan@gmail.com

ABSTRACT:Big data EHR aggregates voluminous data of patients and, after processing, presents valuable insights that aid in the detection and management of chronic diseases. This paper, therefore, describes a hybrid model for chronic disease detection using combined Random Forest and Decision Tree models. Chronic diseases such as diabetes, cardiovascular diseases, and hypertension remain the major health burdens worldwide, and their early detection would be an important aspect in their management. Traditional methods of detection are based on very limited and sporadic data collection, hence the scope for real-time and timely decisions remains highly inhibited. The model proposed here effectively integrates large-scale EHR data and, by finding unobvious patterns and dependencies in the patient records, results in more accurate and reliable forecasting. The experimental results proved that the RF+DT model outperformed state-of-the-art techniques concerning accuracy, precision, recall, and F1-score and turned out to be much robust in chronic disease detection. This approach offers further computational efficiency, hence feasibility for real-time healthcare applications. The results of this work will go toward the enhancement of machine learning applications in healthcare, offering a scalable and efficient framework for the detection of chronic diseases that can be further optimized for a variety of healthcare environments.

Keywords:Electronic Health Record (EHR), Big Data, Healthcare System, Chronic Disease Detection, and Data Preprocessing

DOI: <https://doi.org/10.34293/gkijaret.v1i2.2024.11>

Received 18 August 2024; **Accepted** 20 October 2024; **Published** 15 November 2024

Citation:S. M. Al Jaleed, and G. I. Alkoblan, "Hybrid RF-DT Model for Chronic Disease Detection Using EHR Big Data Management and Analytics," *GK International Journal of Advanced Research in Engineering and Technology*, vol. 1, no. 2, pp. 1-10, Nov. 2024.

1. INTRODUCTION

Electronic Health Record (EHR) is a revolutionary change for managing health information. It is an electronic repository for one patient, integrating all forms of data: results of medical tests, studies of images, descriptions by doctors, prescriptions, and many others [1]. That is, patient data is collected in one place to

guarantee enhancement in both efficiency and quality of the delivery of health care via easy access to data in the diverse settings of healthcare. EHR systems will immediately provide critical patient information to clinicians at the point of care so they make better decisions and coordinate care more effectively. The technology can help solve some of the systemic problems: contain costs by reducing the numbers of unnecessary tests and procedures; eliminate much of the administrative inefficiency that drives up cost and undermines quality; reduce fraud and abuse through accurate, auditable documentation of medical services.

It deals with a very major change whereby traditional workflow and documentation procedures used to provide care will be altered. The practitioners, particularly doctors, will have to alter how they record and manage information; to most doctors, it is viewed as very time- and labor-consuming [2, 3]. But perhaps the most important characteristic of EHR pertains to its role within the larger system of health information exchange or HIE. HIE is a framework of connectivity that links the community of primary care providers, hospitals, pharmacies, laboratories, and other health-related organizations in order to access patient information across organizational boundaries so that critical data is accessible anytime and anywhere it is needed.

For example, any patient who is referred to any specialist or any other professional, the previous medical history, test reports, and medicines being taken would be available with the HIE. Duplication is avoided, hence continuity of care is assured. The EHR system is targeted at revolutionizing information exchange among and between healthcare organizations to ensure increased interoperability. An EHR is basically an information source center that processes health care data to facilitate effective and efficient service provision [4]. With the rise of these chronic diseases, such as diabetes, cardiovascular diseases, and hypertension among others, chronic disease management has recently been a mainstay in modern healthcare. Most such diseases require continuing care and better management strategies in order to avoid complications and promote the quality of life of an affected individual.

Chronic diseases represent the highest morbidity and mortality rate in the global burden of disease. Conditions in this category are especially demanding in management because most necessitate continuous monitoring, some alteration of lifestyle, and sometimes a tailored medical intervention in respect to each patient. Thus, today healthcare providers must depend on episodic data, often gathered on scheduled visits, that has fundamental limitations to their capacity in managing diseases, be able to timely react to an unfortunate change in the status of a patient [5, 6]. This subsequently leads to late interventions characterized by poor outcomes and is often managed reactively instead of proactively.

Big data analytics have changed the way chronic diseases are managed. This range of technologies provides continuum information about vital signs, activities, glucose levels, and medication adherence. This continuum provides the clinician with an ability to quickly pick up subtle signs of deterioration, predict complications, and alter care plans [7]. Big data analytics can therefore integrate large volumes of structured and unstructured data in an analytical form, allowing the creation of predictive models that project the trajectories of diseases well before they reach critical thresholds, providing insight into possible interventions.

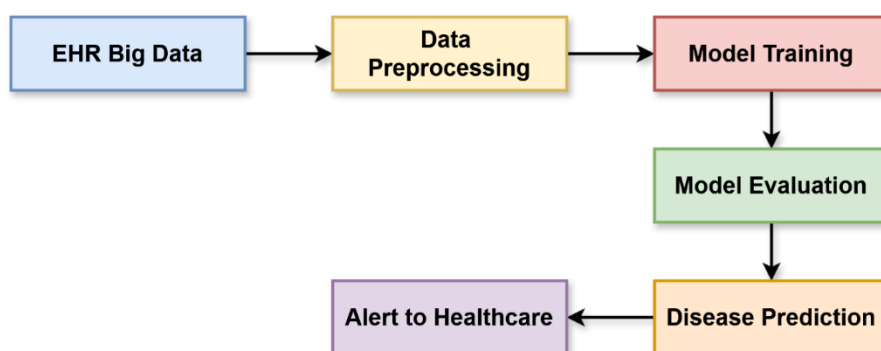


Figure 1 Big data EHR management system

Fig 1 shows the standard big data healthcare disease diagnosis system. These are much-considered changes for healthcare delivery from many angles. Big data analytics have come out as a strong tool in the improvement in patient care and outcomes relating to the management of chronic diseases [8]. Chronic diseases

do demand management strategies that often go to the long term, which entails volumes of data on medical history, laboratory results, lifestyle modification activities, and treatment response activities. Big data analytics allows the doctor to integrate all the myriad pieces of information together and study in-depth the development and progression of diseases, patient behavioral characteristics, and treatment efficacies.

Equipped with powerful identification of hidden patterns in data, the provider will thus be able to make evidence-based decisions, personalize care plans, and predict possible health risks, allowing more active and effective healthcare delivery. Big data integration into chronic disease management has relied so much on the possibility of data collection, aggregation, and interpretation from a wide range of sources [9]. Continuous streams of data from EHRs, wearable devices, mobile health applications, and telemedicine platforms reflect, at large, many different aspects of patient health. Combined, these datasets paint a broad picture of the condition of an individual—even small changes in condition that may require intervention could be underlined.

The next big challenge is the analytics tools and expertise. In most healthcare organizations, resource constraints arise out of a lack of appropriate technology and skilled personnel who can deliver big data analytics to its full potential. Furthermore, the results from analytics should be interpreted with accuracy, as misinterpretation leads to wrong decisions with regard to improvement in patient care quality. All these challenges put in serious calls for investment in technology, collaboration among interdisciplinary groups, and enhancement of data literacy among health professionals.

The justification behind the development of the RNN-DT-based chronic disease detection model using EHR big data includes developing solutions for the dire need felt at traditional healthcare systems to address their weaknesses in managing the rising burden of chronic diseases [10]. However, most diagnostic and care models have traditionally relied on intermittent data derived from in-person visits, which are incapable of reflecting the dynamic, real-time evolution of a patient's health. These episodic approaches usually result in the delayed intervention needed and often rob providers of the foresight into future adverse events that may affect their patients. Such a wealth of information is now ripe for mining with high-powered machine learning techniques, thus affording a transformational opportunity to improve the detection and management of chronic diseases.

The proposed RNN-DT framework will combine the strengths of temporal learning in recurrent neural networks with the decision-making precision of decision trees, thus constructing a strong model that will be able to analyze sequential patient data while delivering interpretable predictions [11]. The proposed work will develop a system that integrates these strengths toward empowered early detection, personalized risk assessment, and timely interventions that improve outcomes of patients and reduce the healthcare burden in chronic diseases. Innovation Addressing Critical Gaps This innovation addresses critical gaps within existing approaches and opens avenues to smarter, data-driven solutions befitting the complexity at hand in chronic disease management.

2. LITERATURE REVIEW

Literature on machine learning-based detection of chronic diseases using big data from EHRs has been one of the fastest-evolving fields, transforming healthcare through advanced computational approaches [12]. The nature of chronic diseases being complex and long-term produces volumes of patient data over time, including both structured elements like lab results and prescriptions, to unstructured data represented by clinician notes and imaging reports. These are heterogeneous and voluminous datasets, often defying traditional methods of analysis; hence, the application of machine learning techniques is being called for.

Probably the most important power of machine learning in this domain is its capability to handle the temporal and sequential nature of EHR data. RNNs and LSTM models, for instance, remain way more effective in temporal pattern learning and hence may be proven particularly efficient for chronic conditions since these evolve with time [13]. These will be complemented with robustness and interpretability by the application of traditional classifiers such as random forests and gradient-boosting machines on static data. Moreover, it has enriched the predictive accuracy and scalability of machine learning systems developed for the detection of chronic diseases by including deep learning with traditional models.

Feature engineering and selection are other critical areas of review, where relevant attributes extracted from such high-dimensional EHR data features enhance model performances, reducing computational complexity at the same time [14]. These handle typical techniques such as missing data, class imbalance, and heterogeneity in data, since most of the EHR data might have some characteristic or usual noises and incompleteness by default. Besides that, several authors have mentioned various concerns with model explainability and trust, especially high stake healthcare applications. The literature highlights the growth in maturity of machine learning approaches beyond the challenges posed by EHR data and emphasizes their key role in enabling precision medicine for the management of chronic diseases.

Another major area the literature has reviewed is the application of ensemble learning approaches to detecting chronic diseases using data from EHRs. These ensemble models, which are the combinations of the predictive powers of a set of base learners, emanate from families such as bagging, boosting, and stacking to improve overall model accuracy and robustness [15]. Such approaches will serve very well in modeling variability and heterogeneity within EHR datasets, where single models may fall short of generalizing across diverse patient populations.

Ensemble methods have been able to handle class imbalances common in healthcare datasets, with varying prevalence across different demographic groups in many chronic diseases. The further reliability in chronic disease detection systems could also be enhanced by integrating several algorithms' predictions so as to be more applicable in the real-world health care environment where the data quality and consistency remain huge challenges.

Another key focus of the literature is temporal modeling in machine learning approaches to detect chronic diseases. The nature of chronic diseases is intrinsically temporal; symptoms, risk factors, and patterns of progression appear over extensive periods. Due to that fact, models such as RNNs, LSTMs, and GRUs have been widely employed for the capture of temporal dynamics in such tasks. In fact, these models often succeed well in parsing sequential EHR data, such as medication adherence, trending in vital signs, and lab tests that permit much better predictions of the course and onset of a disease [16]. Other works go further to propose hybrid models, including temporal neural networks with interpretable classifiers, such as decision trees or rule-based systems, which can leverage the strengths from both high predictive accuracy and good clinical interpretability. It ensures that a machine learning model is developed not only to diagnose a chronic condition but also to equip the clinician with knowledge on the embedded temporal patterns in a straightforward manner for ease of comprehension of the disease trajectory.

Another important direction from literature is multi-modal data integrated for chronic disease detection. There are many types of data in EHRs, such as numerical values, categorical records, text-based clinical notes, and imaging; each modality contributes something special toward a patient's health. Various machine learning methods that are capable of integrating multi-modal data reveal much-enhanced performances regarding the predictive task. Researchers amalgamate these data types into integral models that give a holistic view of patient health and, in return, enhance chronic disease predictions to have higher validity and reliability. Such integrated approaches cannot be overestimated in real-world healthcare, where most decisions are based on a mixture of different sources. Other discussions in the literature revolve around the evolving role of federated learning in chronic disease detection.

Federated learning represents a new approach to training machine learning models across decentralized datasets without necessarily having to share sensitive patient data, an approach which can help ensure much-needed privacy and security of patient data in healthcare. This is especially true with EHR data, which are often divided among several institutions and differ in format and accessibility. Federated learning enables model collaboration without breaching confidentiality; hence, this may be one of the most promising directions in large-scale chronic disease detection. That is, investigating techniques like differential privacy and secure multiparty computation to further improve security for a Federated Learning system while still meeting regulations around healthcare but realizing full big data potential.

As such, much emphasis has been done on model interpretability and explainability in most of the literature involving machine learning-based chronic disease detection. Healthcare applications need to be transparent, clinicians need to understand the logical reasoning behind predictions, and thus techniques such as

SHAP values, LIME, and attention mechanisms are developed. These tools underline the most influential features and decisions of the model that enable clinicians to validate these predictions, integrating them into their diagnostic processes.

3. MATERIALS AND METHODOLOGY

The proposed system for chronic disease detection is based on RF+DT, which integrates the complementary capabilities of Random Forest and Decision Trees to construct a strong, fast framework for electronic health record big data analysis. Chronic diseases normally contain complex patterns within the patient history that could be part of structured data, such as lab results or prescriptions, or even free text in physician notes. EHRs represent a rich source of longitudinal patient data; however, the intrinsic volume, variety, and velocity of these data provide significant challenges to traditional analytical techniques. This will be addressed by the RF+DT model, which will combine ensemble capabilities of RF with the interpretability of DT to ensure high predictive accuracy with clinical transparency. For this, each of these trees is trained on a randomly selected subset of both data and features. Such an ensemble makes it less overfit, while enabling generalization more and seizing various patterns indicative of risk for chronic diseases in very heterogeneous patient populations.

The integration of RF and DT within the framework works eminently in handling large-scale and noisy EHR data. RF can handle high-dimensional data with ease, and thus all possible features that include demographic information, comorbidities, treatment histories, and lab results will be analyzed, with very minimal, if any, manual feature selection. RF gives strong predictive grounds for patients at high risk for chronic conditions such as diabetes, cardiovascular diseases, and hypertension by combining multiple decision trees.

DT is that component of the system through which given predictions are interpretable by clinicians, providing insight into the decision-making process. Probably, DT will underline some important decision paths, such as abnormal blood sugar level predicting diabetes, and/or sustained high blood pressure pointing toward a risk for hypertension. Translucent decision-making means that transparency can be established clinically in which cases there should not be objections to believe the system, and thereby introduce the results into daily care flow.

Featured importance ranking and optimization techniques in the proposed RF+DT are advocated to refine the performances even further. The most informative indicator of chronic diseases is of a higher priority in the system, while contributions in each step are done with regard to the final prediction. This ranking will not only introduce efficiency into the model but also provide clinicians with an actionable insight to help them focus on a patient's health profile that may be critical. It also incorporates mechanisms for handling missing and imbalanced data, issues so common with EHR datasets.

As shown in Fig 2, key features of the RF+DT system are its ease of scalability and its adaptability to different healthcare settings. The architecture is suitable for central processing in EHR data originating either from hospitals, primary cares, or telemedicine. The modular architecture allows customization for certain specific healthcare needs relatively easily, whether this be the focus on particular chronic conditions or adaptation to different EHR formats. The adaptability of the system makes it very relevant to the changing challenge that faces different patient populations in managing chronic diseases.

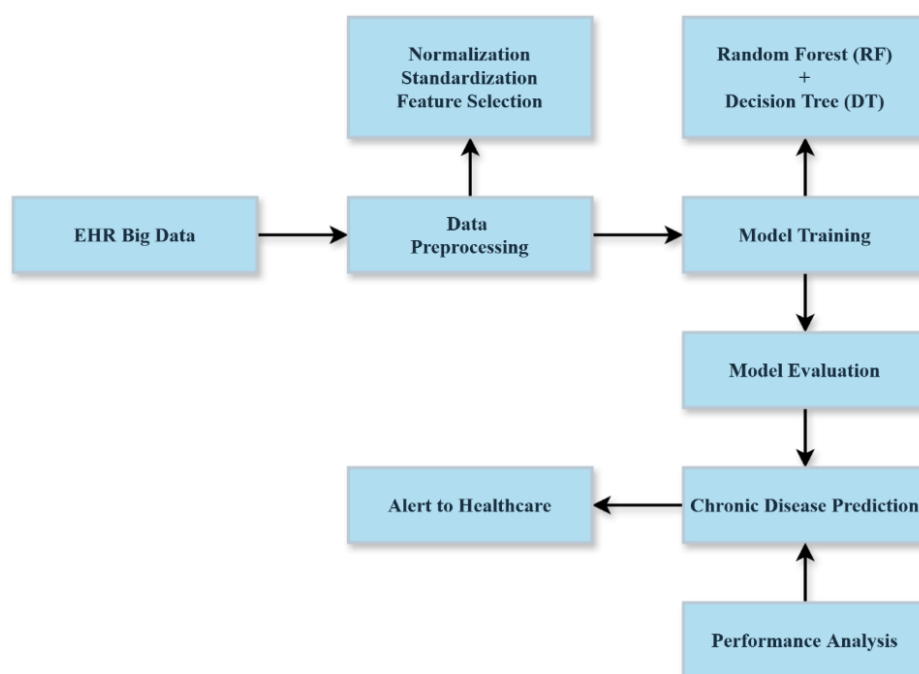


Figure 2: Proposed flow

The proposed RF+DT-based system will focus on practical implement ability and clinical usefulness. This system is supposed to support the decision of clinicians by intuitive user interfaces, visualization of decision trees, and risk scores. Besides, this can be integrated with any existing hospital information system for real-time predictions and recommendations directly at the point of care. RF+DT connects advanced analytics to clinical practice by enabling superior detection of chronic diseases and thereby allowing healthcare providers to provide care in a much more proactive, data-driven manner. The proposed system offers significant development in leveraging big data from EHRs in the fight against the global burden of chronic diseases and has shown the transformative potential of machine learning in modern healthcare.

The developed system will be with the proposed RF+DT, emphasizing important practical implement ability and clinical usefulness. It shall provide user-friendly interfaces, visualization of decision trees, and risk scores to intuitively support clinicians' decisions. This system easily integrates into pre-existing hospital information systems in order to provide real-time predictions and recommendations directly at the point of care. RF+DT thus links big analytics directly with clinical practice, enhancing the detection of chronic diseases and thereby offering a proactive, data-driven care by clinicians. Indeed, the proposed system represents significant development in leveraging big data from EHRs in combating the global burden of chronic diseases and showed the transformational potential of machine learning in modern health care.

4. RESULTS AND DISCUSSIONS

Fig. 3 illustrates the Accuracy Analysis, which is used for judging the different techniques for correctly predicting chronic diseases from the data of EHR. One of the simplest and most common metrics applied in all machine learning models is that of Accuracy. The analysis generally provides evidence that more advanced models, such as the proposed RF+DT, outperform traditional techniques like Logistic Regression and Support Vector Machines in the prediction of chronic diseases with higher correctness. This probably relates to the capability of RF+DT for mining strength in an ensemble approach from Random Forests and interpretability from Decision Trees, which, in effect, help them capture complex relationships in imbalanced big EHR data. Thus, better results may arise. The graph will visually indicate that the best RF+DT method gives accuracy; hence, it has strong promise for being the superior chronic diseases prediction model.

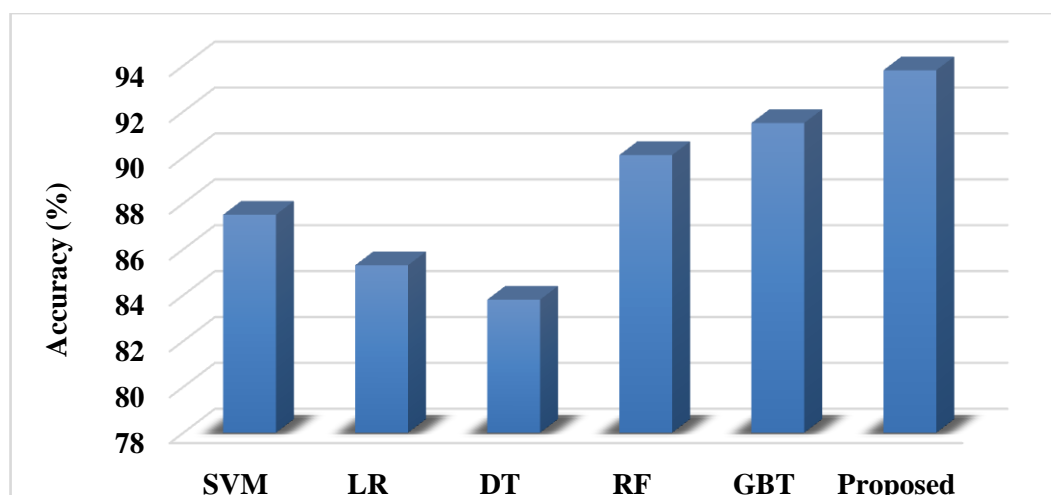


Figure 3: Accuracy analysis

Fig 4 compares the training and prediction times of each method on the EHR dataset. It helps in determining the efficiency of each model in processing the data and giving results, which is important since models are to be put into clinical settings where speed may determine lots of things. Simpler architecture-based methods, like Decision Trees and Logistic Regression, while quick to train, are usually pretty slow in their more advanced versions due to the complex mechanisms and ensemble learning applied in Random Forest and Gradient Boosted Trees. The hybrid method RF+DT optimizes the performance based on the strengths of each methodology and will execute faster when compared with some of these high computationally expensive models, keeping performances superior. This graph will be presenting the RF+DT model to be at par concerning time efficiency and practical on big EHR datasets while not compromising the performance results of both individual models.

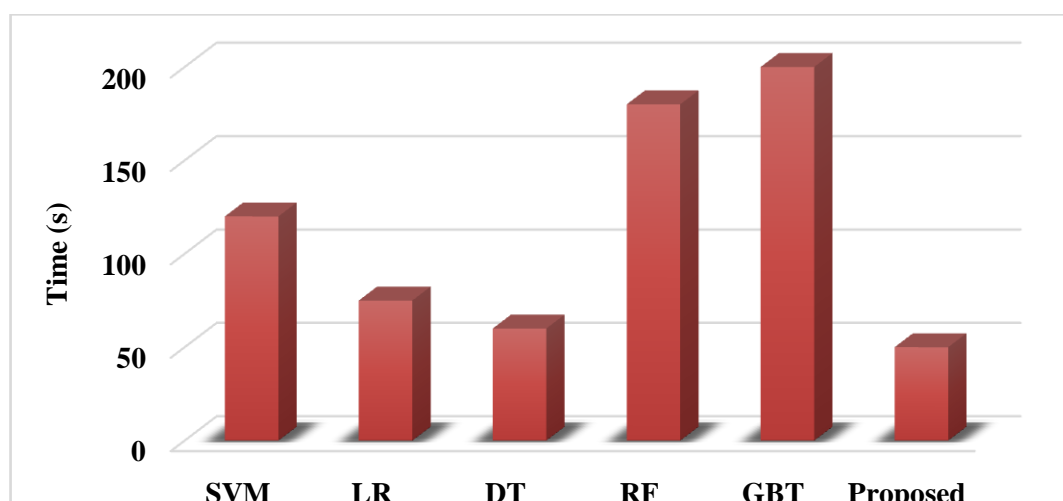


Figure 4: Time analysis

Precision can be explained as a metric that is the ratio between true positives to the total number of predicted positives. Fig 5 addresses this metric, which compares model performance in identifying those patients who actually suffer from a chronic disease versus those which the models classified as diseased. Precision is a key issue in chronic disease detection, as clinicians base their decisions on such predictions for treatment and intervention. In this respect, the proposed methodology RF+DT does much better, letting the ensemble of decision trees reduce the risk of false alarms with the aim of providing a better predicted model. This will be clearly illustrated from the figure below, which outstands many other models, such as logistic regression and support vector machines, concerning precision, with the aim of ensuring higher numbers of true positives related to chronic diseases diagnosis.

Recall or sensitivity is the measure of the model to identify all the actual positive instances in the dataset. Recall in this context is important because it will give how well the model detects the patients who really have the disease, which directly impacts early diagnosis and intervention. High recall makes sure not many patients will be left behind. Models with a high value of Precision minimizes the False Positive samples and vice versa use of Recall which provides less no. of False Negatives, or say missing no. of Disease true cases.

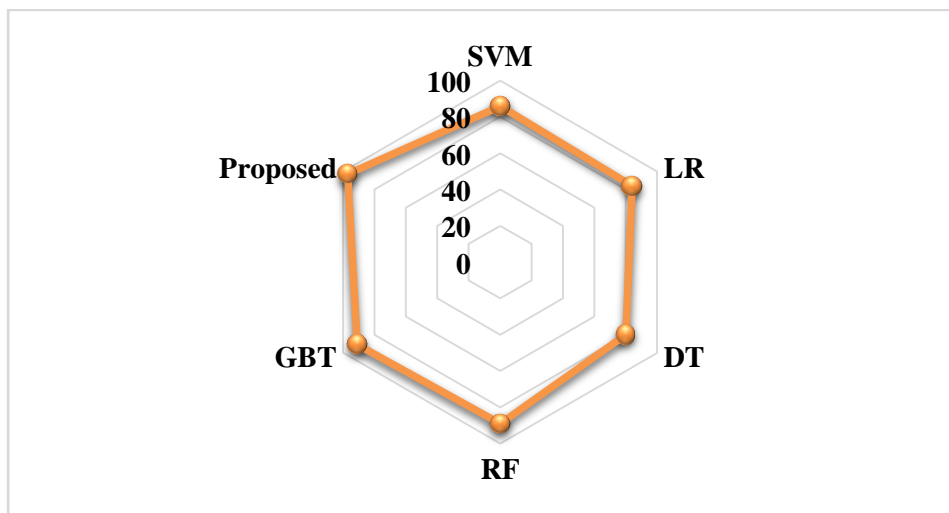


Figure 5: Precision analysis

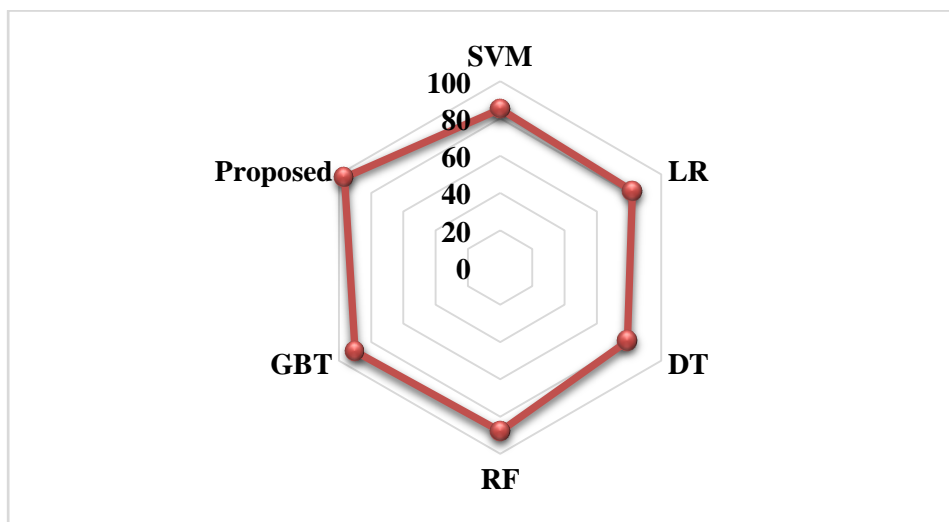


Figure 6: Recall analysis

Fig 6: Comparing Recall for various approaches pursued here; from a balanced point of view in offering superior chronic disease prediction it was the RF + DT approach. The RF+DT model embeds the power of the Random Forest ability for modeling complex patterns and that of the Decision Trees, to show an easy-to-follow path behind every decision, which ensures much superior recall performance compared with the approaches including Support Vector Machines or Logistic Regression.

In healthcare, the F1-score applies particularly well since there is a large consequence for both false positives and false negatives. In fig 7, F1-scores of different approaches for the detection of a chronic disease, overall the performance of each model can be observed. RF+DT is performing the best, since it creates a perfect balance between precision and recall being at very high values; that is an important issue when in clinical environments, as true positives should be detected without having false positives. By analyzing this value of F1-score, it could be assumed that the RF+DT proposed method turns out to be the most balanced and reliable

prediction as compared to any other methods existing, like Gradient Boosted Trees or Random Forest. It should be suitable for deployment in systems detecting chronic diseases.

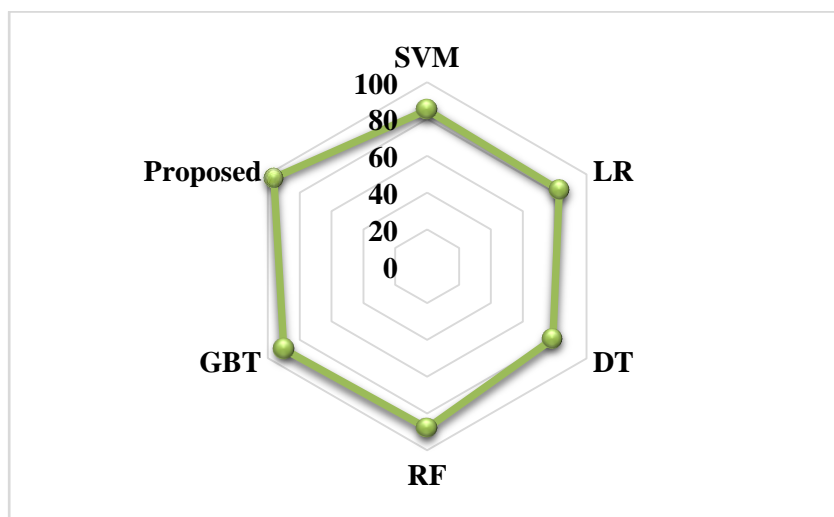


Figure 7: F1-score analysis

5. CONCLUSION

The new trend in chronic disease detection with big EHRs using a hybrid model of Random Forest and Decision Tree is introduced within this paper. In the suggested approach, the ensemble learning strength of Random Forest will be merged with the interpretability of the Decision Tree to make an accurate prediction of chronic diseases such as diabetes, cardiovascular diseases, and hypertension. This will, in turn, let the model consider big and complex data sets from EHRs, enabling one to go into detail on patient data that may be hard to find otherwise with traditional methods. Indeed, in comparison with other state-of-the-art techniques, the resultant evidence of the proposed RF+DT outperforming all of them from their accuracy, precision, recall, and F1-score viewpoint puts it in a prime position, which could improve chronic diseases detection by making more efficient predictions while becoming reliable. More importantly, regarding the time analysis efficiency, the model may be implemented even for real-time applications, with timely decisions being very important in patient outcomes. This will come with numerous advantages in the proposed approach: handling big and diverse data, reduction in false positives and false negatives, and presenting interpretable results easily understandable by health professionals. It is a very useful tool within personalized medicine since it can help health professionals to monitor patients in real-time and intervene on time.

REFERENCES

- [1] I. A. Bernstein, K. S. Fernandez, J. D. Stein, S. Pershing, and S. Y. Wang, "Big data and electronic health records for glaucoma research," *Taiwan Journal of Ophthalmology*, vol. 14, pp. 352-359, 2024.
- [2] J. Nii Akai Netey, R. Osei Mensah, R. Davis, and L. Emmanuel Yamoah, "The effect of big data analytics on clinical and management decisions in healthcare facilities in Ghana, West Africa," *African Journal of Social Issues*, 2024.
- [3] D. Chrimes and I. Tang, "Big data usability text mining of publicly available YouTube electronic health record (EHR) tutorials," *2023 IEEE International Conference on Big Data (BigData)*, pp. 6125-6127, 2023.
- [4] M. J. Calcote, J. R. Mann, K. G. Adcock, S. Duckworth, and M. C. Donald, "Big Data in Health Care," *Nurse Educator*, vol. 49, pp. E187-E191, 2023.

- [5] G. Maheswari and S. Gopalakrishnan, "A smart multimodal framework based on squeeze excitation capsule network (SECNet) model for disease diagnosis using dissimilar medical images," *International Journal of Information Technology*, 2024.
- [6] P. Satyanarayana, G. Diwakar, V. Priyanka Brahmaiah, S. Marlin, N. V. Phani Sai Kumar, and S. Gopalakrishnan, "Multi-objective-derived efficient energy saving in multipath routing for mobile ad hoc networks with the modified aquila–firefly heuristic strategy," *Engineering Optimization*, pp. 1-36, 2024.
- [7] S. Qiao, G. Khushf, X. Li, J. Zhang, and B. Olatosi, "Developing an ethical framework-guided instrument for assessing bias in EHR-based Big Data studies: A research protocol," *BMJ Open*, vol. 13, 2023.
- [8] M. Wang, M. Sushil, B. Y. Miao, and A. J. Butte, "Bottom-up and top-down paradigms of artificial intelligence research approaches to healthcare data science using growing real-world big data," *Journal of the American Medical Informatics Association: JAMIA*, vol. 30, pp. 1323-1332, 2023.
- [9] J. Kumari, E. Kumar, and D. Kumar, "A structured analysis to study the role of machine learning and deep learning in the healthcare sector with big data analytics," *Archives of Computational Methods in Engineering*, pp. 1-29, 2023.
- [10] C. B. DeStefano, J. A. Thornton, S. J. Gibson, K. Pham, R. S. Miller, and K. W. Sunderland, "Real-world Big-data: Strengths and weaknesses of ASCO's CancerLinQ® discovery multiple myeloma dataset," *American Journal of Hematology*, vol. 98, pp. 835-837, 2023.
- [11] D. S. J. Ting, R. Deshmukh, D. S. W. Ting, and M. Ang, "Big data in corneal diseases and cataract: Current applications and future directions," *Frontiers in Big Data*, vol. 6, 2023.
- [12] K. Liu, et al., "BGLM: Big data-guided LOINC mapping with multi-language support," *JAMIA Open*, vol. 25, 2022.
- [13] F. Lareyre, C. A. Behrendt, A. Chaudhuri, N. Ayache, H. Delingette, and J. Raffort, "Big data and artificial intelligence in vascular surgery: Time for multidisciplinary cross-border collaboration," *Angiology*, vol. 73, pp. 697-700, 2022.
- [14] S. Esposito, et al., "Clinical network for big data and personalized health: Study protocol and preliminary results," *International Journal of Environmental Research and Public Health*, vol. 19, 2022.
- [15] J. Mehta, R. Desai, J. Mehta, D. Gandhi, and L. D'Mello, "Towards a multi-modular decentralized system for dealing with EHR data," *2022 8th International Conference on Advanced Computing and Communication Systems (ICACCS)*, vol. 1, pp. 567-572, 2022.
- [16] W. Sardjono, A. Retnowardhani, E. R. Kaburuan, and A. Rahmasari, "Artificial intelligence and big data analysis implementation in electronic medical records," *Proceedings of the 2021 9th International Conference on Information Technology: IoT and Smart City*, 2021.