

# A Comprehensive Review of Machine Learning and Feature Reduction Techniques for Effective Disease Risk Prediction in Healthcare Systems

<sup>1</sup>\*Yashwani Singh, <sup>2</sup>Deepshikha Patel

<sup>1</sup>\* Department of Computer Science & Engineering Oriental Institute of Science & Technology Bhopal, India

<sup>1</sup>\*[yashwanisingh61@gmail.com](mailto:yashwanisingh61@gmail.com), <sup>2</sup> [deepshikhapatel@oriental.ac.in](mailto:deepshikhapatel@oriental.ac.in)

\* Corresponding Author: Yashwant Singh

## Abstract:

*A rapid surge in health-care budgeting and occurrence of chronic diseases has now trended the utilization and/or advancement of machine learning-driven algebra for enhanced examination and registry of accurate "disease-risks prediction". Machine Learning (ML) and Data Mining (DM) have served as highly-efficient tools for processing medical datasets, establishing expert systems needed for certain clinical outcomes. In the forthcoming coupon, critical overview and direction are established for disease-risk predictive and existing patterns of adaptations in operationalizing health systems' context in relation to feature filtering and optimal feature selection. Moreover, many supervised learning algorithms, from Logistic Regression (LR) and SVM to RF and SGB; from BRT to DT; and from NB to NN, have gracefully trodden the path towards cardio- and chronic disease forecasting. In depth, the paper emphasizes an important role in the field, i.e., feature selection, which helps to improve the prediction capabilities while managing with the issues concerning the "efficiency," "redundancy," "convergence speed," and "speed of computational convergence." The different feature selection approaches, for example, the wrapper-based, filter-based, hybrid, evolutionary genetics, principal component analysis, probabilistic principal component analysis, Relief, fuzzy systems, and ensemble techniques, are carefully examined. In further discussion, particular stress has been laid on the importance of AUC and accuracy for purporting good performance optimization for strong medical diagnosis. Units also discuss handling missing values, skilled resolution of the imbalanced dataset by the SMOTE tool, and the cleaning of healthcare datasets purchased from repositories like the University of California Irvine (UCI) Machine Learning Repository. The other major research gaps in contemporary disease prediction systems were delineated by the review, which included the limited opportunity for effective feature-reduction frameworks, the need for more emphasis on AUC-based evaluation, the high computational time, and the challenges of imbalanced datasets. The integration of intelligent feature reduction models into machine-learning algorithms, the study concludes, could give a significant boost to the integrity of disease prediction, reduce medical expenses, and support early detection in modern health care settings.*

**Keywords:** Machine Learning, Disease Risk Prediction, Feature Reduction, Healthcare Analytics, Data Mining, Cardiovascular Disease Prediction

## I. INTRODUCTION

Disease risk prediction has gained importance due to the enormous increase in chronic diseases, increasing patient populations, and massive creation of medical data by electric-healthcare (e-health) facilities. Some common chronic diseases include cardiovascular diseases, diabetes, cancer, kidney ailments, cerebrovascular diseases, etc., causing a lot of mortality in the world-which in turn take drag resources of healthcare hard enough to bear immense pressure as well as casting clouds upon the economy. The early detection of the extent risk of disease aids provision for immediate treatment thus reducing death rates, minimizing healthcare costs, and, at the same time, improving the quest in the quality of care the patient receives. Disease diagnosis in the past was majorly dependent on the clinical acumen, laboratory tests, and the patients' history [1]. However, human intervention can be time wasting and could sometimes trigger errors due to situations that can be gleaned from the data set of the error from the complexities and spectral curse of some medical data. The need for a digital trust structure that supported developing learning ratios in healthcare at once meant that computational mechanisms stepped off with the task. Disease risk prediction systems study historical and real-time datasets of the patient's medical history to reveal the geometrical probability while its effects  $x$  in the established condition. These predictions correlate these past and current condition parameters, such as age, blood pressure, cholesterol level, glucose level, heart rate, lifestyle habits, medical history, and varied clinical measures with disease pattern and risk factors. The digitization of health records through Electronic Health Records (EHRs), wearable devices, smart sensors, and hospital management systems brings about accelerated growth in predictive healthcare technologies. Healthcare analytics is a revolutionary domain that combines statistical analysis, data mining, artificial intelligence, and machine learning techniques to extrapolate significant insights from medical data sets. Through healthcare analytics, hospitals and healthcare organizations enhance diagnosis accuracy, optimize resource utilization, manage treatment costs, and enrich their decision-making procedure [2]. The effective application of healthcare analytics helps identify hidden trends, relationships between clinical attributes, predict disease progress, and make way for personalized medicine. The implementation of the innovative

analytical techniques in healthcare enhances the provision of relevant patient metrics with large-scale monitoring and preventive healthcare strategies. Additionally, analytics can reduce hospital readmission rates, mitigate chronic illnesses, and instruct evidence-based clinical insights for reference. Machine learning receives attention from all directions in healthcare analysis primarily because of its ability to establish learning patterns and establish predictive models right from historical data without needing explicit programming. Meanwhile, there are quite a few illustrative models of prediction benchmarked on heart diseases, diabetes, liver diseases, kidney diseases, breast cancer, and so forth. Machine learning further contributes to automate the healthcare setup for quick diagnosis and decreases dependency on manual interpretation. The Lean-back artificial learning models are proving to be good models in image detection, forecasting, and monitoring patients' illnesses. Even with the breakthroughs made in machine learning-based healthcare systems, many challenges need to be addressed in predicting diseases robustly and reliably [3].-Further, nearly all medical databases find themselves being congregated with various numbers of irrelevant, redundant, and noisy attributes. In healthcare databases, high-dimensional features are usually contained in either the irrelevant or less importantly contributing features. For every unnecessary attribute included, the computational rate gets laborious; processing times go up with the increased burden, the model becomes difficult to interpret, and eventually it windows down the accuracy intended to be achieved. Thus, any healthcare data analyst is required to become vigilant to validate any potentialities to feature selection and feature reduction. Feature reduction is the process that focuses on detecting and eliminating the relatively unimportant attributes from a dataset while preserving a high level of accuracy, hence making it an essential component of machine-learning methods in healthcare. Reducing the number of features leads to reducing overfitting so that the learning algorithms can generalize over a dataset. The following methods fall under the umbrella of feature reduction: filter methods, wrapper-hybrid methods, heuristic methods, PCA, Genetic Algorithms (GAs), and others. These research methods are involved in extracting the most contributing relevant features for clinical or disease diagnosis. Overall, when applied to any healthcare application, number reduction of medical tests and diagnostic parameters may directly reduce healthcare costs as well as patient suffering. The appropriate set of evaluation metrics is likewise very critical for assessing the performance of disease prediction models. The predictive modelling for biomedical cases employed the accuracy, but herein lies the problem that it may give extreme feedbacks about the minority class in the model since medical datasets are imbalanced. Therefore, to quantify the reliability of predictive indicators and their engagement in any future model implementations, metrics emphasize the use of Area Under the Curve (AUC) as substitutes to identify diseases and their probability distribution in healthcare systems [4]. Thereon, for future prospects, among the fields of sensitivity, specificity, precision, recall, and Receiver Operating Characteristic (ROC) analysis. Combining feature selection with machine learning technique enhances the utility of disease prediction and develops a predictive system able progress is made to enhancing feature reduction with intelligent and sophisticated machine-learning algorithms that have been designated to coordinate outcomes with augmented precision and upon detecting emergent critical diseases. Lastly, the paper targets showing promising aids in likely outcomes for potent, scalable, and extremely reliable prediction systems, which could, in the long run, support most medical practitioners in advancing treatment schemas. Figure 1 illustrates the overall machine learning-based framework for disease detection, showing the flow from data collection and preprocessing to feature extraction, model training, and final disease prediction.

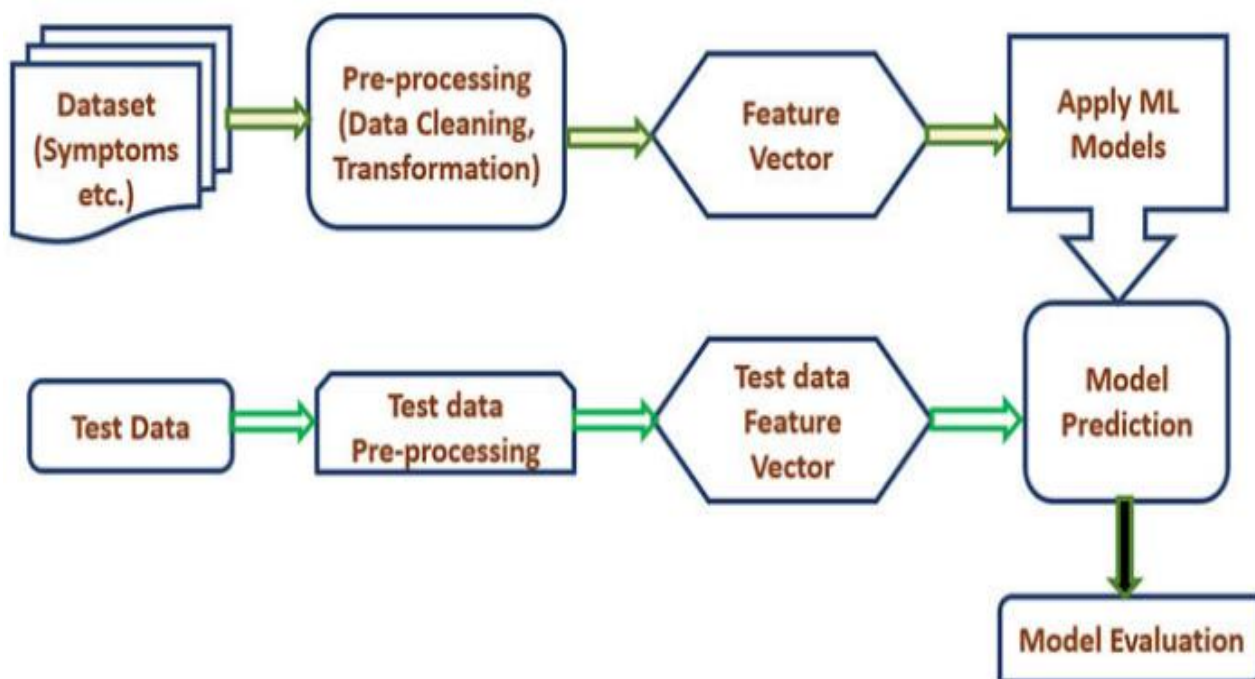


Figure 1 A framework for disease detection method using Machine Learning [22]

## II. OVERVIEW OF DISEASE RISK PREDICTION SYSTEMS

Disease predictions systems are the most intelligent healing apparatuses designed to assess the prediction of diseases by making use of the patient's pathologic events, indirect synonyms, and laboratory reports and past health care notes. These systems help healthcare providers make the right choice of medical vs prevention and health care management. Modern-day disease predictors need to incorporate data analytics, machine learning, or any artificial intelligence in overall process flow. This is to provide choices and make things quick concerning diagnostics, minimize medical errors, and promote effective clinical decision-making.

### a. Traditional Disease Prediction Methods

The traditional methods of disease prediction typically benefited manual medical examination, physician experience, and laboratory findings and data analysis for interpretation regarding patient care. Before the advent of computer technologies, healthcare professionals traditionally decreed diseases based upon clinical symptoms, obtaining histories by means of physical examinations, and relying upon medical expertise. Physicians used the typical diagnostic methodologies like blood pressure measurement, electrocardiograms, X-rays, laboratory blood tests, urine analysis, and imaging techniques to find out what diseases are, and conditions held by patients. These methods were the backbone for healthcare diagnosis over decades and played a crucial role in the identification of infectious diseases, cardiovascular events, diabetes, and other chronic diseases. Yet, traditional methods of prediction were very time-consuming, expensive, and required high level of expertise and reliance on the medical practitioner judgment [5]. In multiple cases, manual interpretation of the medical records and clinical reports resulted in discrepancies and diagnostic mistakes. According to surveys, interpretation of human disorder happens due to numerous reasons: poor vision--fatigue, tons of challenges overcome, considerably high incompleteness of information, lack of expertise, handling and conveying enormous data within minutes, or complexities in medical data. As the infrastructure became serviceable and hospital facilities increased under the purview of operating whole healthcare portfolios with a potentially greater population base, these systems had provided the turf for plenty of clinical health data learnings. Big data learning in the medical industry was riddled by myriad obstacles for pediatric applications across the medical spectrum. Figures, alias big clinical data, have been presented to include patient electronic medical records, laboratory information, physician's notes, and so on. Traditional statistical methods were generally perceived as capable strategies to observe disease patterns and risk factors. Logistic regression is one such most commonly used statistical model for prediction of disease occurrences because it estimates the probability of binary outcomes representing the presence or absence of diseases. Similarly, linear regression models were quite specialized in linear relationships between the selected clinical variables and the disease severity. In order to proceed with medical diagnosis by adhering to the predefined clinical rules along with probabilistic reasoning, Bayesian probability techniques and rule-based expert systems were established. These health decision-making tools were based on a series of "if-then" rules put together by the disease-domain experts to determine disease expression from patient symptoms. Though they offered some support to the healthcare provider, these expert systems suffered from the very core opportunity to dynamically adapt to new datasets or the complexity of the medical domain [6]. Traditional methods for disease prediction mainly made use of clinical guidelines and standard medical protocols researched for many years and produced purely diagnostic comparisons, based on patient symptoms against criteria for the diseases commonly defined. This method had proven to be highly efficient in diagnosing and treating common diseases with clearly recognizable symptoms, but effectiveness was compromised in the context of complex diseases based on multiple interacting risk factors. Chronic diseases like cardiovascular diseases, cancer, diabetes necessitate an analysis of numerous patient attributes like age, cholesterol, blood sugar, blood pressure, lifestyle habits, and genetic factors. Because of the increasing difficulty of manual processing of such types of multidimensional healthcare datasets, the uptake of traditional-method attributes in modern predictive healthcare became equally inefficient in these circumstances. Another challenge with the traditional methods was the handling of large datasets that the modern hospital and diagnostic systems typically generated. As hospital health records were digitized, enormous volumes of patient data became available in electronic files, which really rendered manual processing impractical. These approaches do not hold for automation and predictive features, operating naturally well to observe diseases after the symptoms had already struck, while not showing any promise for the prediction of risks at the preliminary stages. These methods are, for that matter, generally known to be riddled with shortcomings when used in practice dealing with incomplete, noisy, or imbalanced medical datasets. Delay in diagnosis, misinterpretations, and the difference in inappropriate management could all induce an increase in healthcare costs and even damage the patient. Despite the shortcomings, traditional disease prediction techniques provided the initial framework for medical diagnosis and healthcare decision-making [7]. They contributed greatly to the development of modern predictive healthcare systems to recognize patient information analysis and evidence-based diagnosis. Many statistical techniques have held on from the traditional healthcare system and serve as a backbone to the modern machine learning algorithms being used in disease prediction today. Moreover, conventional medical practices have still been of indispensable value within the clinical practice when physician expertise and clinical knowledge validate any machine prediction while ensuring patients' welfare. In further advancement, with the evolution of multiple computational, AI and data-driven approaches, investigators have injected these in the conventional systems to overcome the limitation and improve the performance of prediction. Owing to the implementation of such technologies, the era of dedicated health-care artificial intelligence thus commenced with smarter technologies analyzing big medical data more effectively and efficiently than traditional systems [8].

## **b. Intelligent Healthcare Systems**

The modern healthcare transformation in the healthcare sector, intelligent healthcare systems involve the use of artificial intelligence, machine learning, data mining, Internet of Things (IoT), cloud computing, and big data analytics. The healthcare facilities are now positioned to automate medical services, improve the diagnosis of diseases, expedite clinical decision-making, and enhance patient care through treatments suited for each individual. These intelligent healthcare systems help analyze great volumes of data from medical facilities like hospitals, diagnostic labs, sensors, electronic health records, and mobile healthcare applications. The proliferation of healthcare technologies of the digital age has widened the horizon for patient data introductions for developing intelligent models capable of predicting disease with greater accuracy and efficacy. Intelligent healthcare systems are wonderful for analysing high-dimensional and complex medical data on the fly. They are, in addition to the traditional healing practices, capable of carrying out an analysis on the data using some of the hidden patterns by use of the machine learning and artificial intelligence [9]. They are of great help to the physician for early detection of disease, risk assessment, treatment recommendation, and monitoring of patients. Intelligent healthcare employs machine learning algorithms such as Logistic Regression, Decision Tree, Random Forest, Support Vector Machine, Naive Bayes, and Neural Networks that have been used widely in various applications for disease classification and prediction. The deep learning process improves the performance of healthcare analytics to such an extent that it enables automated recognition of images, medical imaging analysis, and disease detection from radiology scans as well as pathological images. Intelligent healthcare systems seem to be a much-needed requirement in preventive healthcare. By identifying high-risk patients way before they even develop serious illness, predictive models actually look at patients' attributes, which include blood pressure, cholesterol levels, glucose levels, age, weight, genetic information, and lifestyle habits in order to determine probable disease, recommend preventive actions, and augment hospitalization rates, thereby reducing healthcare costs and improving patient outcomes. The use of wearable healthcare devices and IoT sensors have served as an enhancer for these systems with regards to continuous patient monitoring. The data collected by these devices are transmitted in real-time to healthcare systems to analyze and predict health conditions. Besides this, cloud computing technology provides scalable storage and processing abilities to high-powered big data analytics solutions to tackle valuable information from both structured and unstructured medical data. Telemedicine has provided the avenue to remote medical consultations even in rural and underserved regions, where access to medical specialists is restrained. They combine these platforms with intelligent diagnostic systems to screen patients at a distance and advise them in turn, thereby assisting remote monitoring of patients by doctors [10]. Meanwhile, intelligent healthcare technologies have been most beneficial during global health crises, especially in the areas of disease surveillance, outbreak prediction, patient tracking, and healthcare resource management during pandemics. Another fundamental area of intelligent healthcare systems is the use of natural language processing (NLP) in analyzing the natural language to understand and extract clinical notes, medical reports, and physician prescriptions. These natural language processing methods help in extracting useful information from healthcare documents to aid in more refined clinical decision-making. Intelligent healthcare systems are also indispensable in the development of personalized treatment based on individual patient characteristics and genetic information within the paradigm of precision medicine. Nevertheless, intelligent healthcare systems are fraught with various challenges, such as data security issues, cyber-attack threats, and interoperability issues, as well as incomplete standardization of healthcare datasets and ethical considerations [11]. Furthermore, healthcare datasets themselves include some issues, such as missing values, noisy data, and class imbalances, all of which are likely to affect predictive-accuracy. Journalists are really trying their best to provide various techniques for pre processing, strong techniques for feature selection, and an array of in the realm of making some more detailed explanations in their AI-based systems, simply to generate transparent healthcare analytics. In general, intelligent healthcare systems have been changes, which lead to accurate disease prediction, efficient patient management, and medical decision-making that is based on data. The integration of intelligent healthcare systems with strong computational technologies not only improves healthcare quality, reduces medical error, but also enables the evolution of the smart healthcare environment.

## **c. Clinical Decision Support Systems**

CDSS, or Clinical Decision Support Systems, use information technology to aid the doctors, medical staff, and medical establishments in making sensible medical decisions. Using patient-specific data, these systems facilitate the provision of guidance cues, alerts, diagnostic ideas, treatment plans, and risk assessments for improved quality of health care and patient safety. They are growing increasingly important in the current healthcare setting because they combine medical knowledge, patient data, and intelligent algorithms to compound the world of evidence-based medical practice. The CDSS came about in view of the necessity to diminish diagnostic errors, enhance the effectiveness of treatments, and manage the growing complexity in health information. In the past, health care systems depended heavily on the experience of medical doctors and manual interpretations of clinical data, thereby resulting in differences in opinions and tardiness in diagnosis. Exclusive features like calculated automated evaluation of medical information provided by the implementations of CDSS have been made available to physicians to make them better prepared in real-time while diagnosing sceneries. This technology can be employed from within an enterprise environment to enable the sharing or linking of data among the EHRs and LIS's and RIS's and HIS, all working in a real-time setup. A typical CDSS consists of a knowledge base, an inference engine, and a user interface. Knowledge bases contain medical rules, clinical guidelines, disease information, drug interactions, and treatment protocols consigned from healthcare experts and numerous research studies [12]. The inference engine processed patient data using available data in logical reasoning, statistical analysis, or machine learning algorithms and set forth recommendations and predictions. User interfaces facilitate the interaction between physicians or healthcare

professionals and the system, allowing for diagnostic input straightforwardness. CDSS can come in the form of either knowledge-based or non-knowledge systems. The knowledge-based CDSS system depends on previously set rules that assist in deducing the range of diseases, while machine learning methodologies, include those that plays on artificial intelligence, are just embedded into the non-knowledge-based system to provide a solution. The application of machine learning to CDSS is a trendsetter basically for its ability to adapt in dynamics and improve their predictive accuracies step by step. Decision Tree, Logistic Regression, Support Vector Machine, Random Forest and Neural Networks are varied machine learning algorithms widely used in the CDSS for disease diagnosis and patient risk assessment. CDSS works well in a number of areas such as disease diagnosis, medication management, treatment plans, medical record intelligence, continuous monitoring, and wellness development. The decision system for disease prediction in the cardiovascular field works based on patient attributes such as blood cholesterol, blood pressure, heart rate, status of diabetes, and smoking, so as to predict disease risk. On the same token, CDSS for cancer diagnosis may use imaging data and clinical data to identify the tumor and classify the stage of the disease. Medication-tied CDSS alert healthcare providers of drug interactions, allergies, inappropriate dosages, and adverse drug reactions, significantly reducing medication errors in one way and fostering patient safety in another. This technology further helps in cutting down healthcare costs and making operations more efficient in hospitals by ensuring early illness detection and correct treatment planning, as a result of which unnecessary laboratory tests are avoided, hospitalization rates cleaved down, and resources optimally allocated. Again, personalized medicine prospects are enhanced by CDSS in that the systems are capable of recommending bespoke treatment options in respect of patient-specific parameters and medical records. Integrating CDSS with smart healthcare systems, IoT devices, cloud computing, and big data analytics may take its capabilities further still in the biggest manner in today's healthcare setup [13]. On the other hand, several problems pose challenges for CDSSs' implementation and usefulness. They point to data privacy challenges affecting interoperability among healthcare systems, healthcare professionals being ignorant of alerts, and not having relevant data to feed into the algorithms. Also, diagnosing concerns are related to the process of interpreting machine learning predictions. Another factor involves alert churn. Usability decreases from excessive alerts and recommendations, resulting in alert fatigue for the practitioner. Consequently, efforts are directed at establishing a scientific approach for addressing the limitations described above. They are looking at ways of improving transparency, reliability of input, and acceptance by the clinical staff and patients. As a result, there is still a lot for CDSS that herald immense advancements in modern healthcare, such as the life-saving predictive accuracy for diagnostic services, enhancing clinical decisions, reducing medical errors, and ultimately bettering patient care through intelligent data-driven technology.

### III. MACHINE LEARNING TECHNIQUES FOR DISEASE PREDICTION

Machine Learning has made it possible to predict diseases with a high degree of accuracy by analyzing large volumes of medical data, identifying hidden patterns and relations in different data sets. Early prediction of a patient's disease aids healthcare givers, allowing them to plan the treatment in advance; assess diseases up-front; and even monitor several facets. The conventional healthcare system relied on physician expertise for interpretations and tested results related to diagnostics, sometimes with delay and inaccuracies. Contrarily, the machine learning modes can learn from historical healthcare systems and clinical data to develop predictive models for earlier identification of disease causations. The sudden availability of Electronic Health Records (EHR), sensed wearable data, medical imaging data, and clinical reports advanced the technology's applications in Healthcare Analytics.- These technologies were widely applied for the prediction of cardiovascular diseases, diabetes, cancer, kidney disorders, hepatic diseases, other chronic illnesses. Machine learning predictions involve patient variables like age, blood pressure, cholesterol, glucose, BMI, smoking habits, genetic information, and medical history, used to categorize them into positive patient and negative patient categories.

This prediction approach employs typical supervised machine techniques including Logistic Regression (LR), Support Vector Machine (SVM), Random Forest (RF), Decision Tree (DT), Naive Bayes (NB), Neural Networks (NN), and ensemble learning. Each algorithm has its set of advantages, limitations, and performance measures, depending on the way the datasets and prediction tasks are related to the medical field. Logistic Regression (LR) related stats and machine-learning techniques provide a very simple, interpretable, and useful method for making disease predictions. LR is commonly used for binary classification tasks where the output variable stands for the existence or absence of a disease. The algorithm involves estimating the probability of disease occurrence by mapping the values derived using logistic or sigmoid functions from the input feature variable between 0 or 1; thereby it makes a prediction. Logistic regression proves to be vital for the healthcare community, as this algorithm provides a rationale for the links between the clinical variables and the outcomes through disease prediction [14]. It is an obvious choice in the prediction of heart disease, diabetes, and cancer risk in public health management from the data commonly present in the domain. He, a rule-based classifier, performs very well on linearly separable data, way low in computational expense, but his effectiveness may wane against extremely complex, nonlinear medical data. Support Vector Machine (SVM) is yet another powerful supervised machine learning algorithm of great usage for predicting class label in disease classification and predictions.

SVM operates by discovering the optimal hyperplane separating the different classes in the dataset, creating the maximum margin between them. Thus doing, it is very effective at handling high-dimensional healthcare datasets, as it is applicable to linear as well as nonlinear classification-space using kernel functions, which include linear, polynomial, radial basis function (RBF), and sigmoid kernels; for its high accuracy in predicting cardiovascular diseases, cancer, and classifying

medical images for handling and learning from complex data. One of the most enduring advantages of SVM is the protection afforded against overfitting, particularly when the number of features to be considered exceeds the number of training points that exist. Nonetheless, one disadvantage of SVM would be its need for a considerable amount of parameter estimation and computational expenses that would set in when handling very large data sets from healthcare sciences' domain [15]. Random Forest (RF) builds upon Decision Trees by using an ensemble of trees to correct overfitting and boost the learning performance. The algorithm bootstraps the random selection of feature subsets and training data for building as many decision trees as specified. The Random Forest is ideally suited for disease prediction since it can handle missing values, noisy data, and high-dimensional medical datasets with sufficient feature importance analysis for healthcare research and professional use in identifying highly influential clinical attributes for disease prediction. Owing to its good performance in terms of high prediction accuracy, Random Forest has been used in predicting cardiovascular diseases, diabetes, and chronic kidney diseases. About its ability to maintain its predictive accuracy and robustness, RF is less sensitive to outliers and data imbalance compared with other traditional model learning techniques. Decision Trees (DTs) are widely known as a popular machine learning algorithm in healthcare applications for diagnosis diseases and aid clinical-decision-making systems [16]. Decision trees create a tree-like structure where internal nodes represent criteria based on feature values, branches represent outcomes, and leaf nodes represent final classifications or predictions. Decision trees produce easy-to-understand and interpretable outputs, which means they are very suitable for medical environments where transparency is an important aspect. Physicians can quickly interpret the decision built by trees and thus identify as to when specific patient attributes participate in the prediction of yet another ailment. Decision Tree algorithms have frequently been applied for issues to predict heart disease, diagnose diabetes, or classify cancer. One disadvantage of Decision Trees that is quite worrisome will be their vulnerability to overfitting when they are trained on very small or noisy data. Naive Bayes (NB) is a classification technique based on the Bayes probabilistic formula using the strong assumption of feature independence. Naive Bayes is simple but very effective in many applications due to its strength with the relatively low-dimensional input set of the identically distributed features, which represent wealth procurement for health systems applications. This theory generates disease probability based on both prior probabilities and input feature likelihood. Since NB is highly computationally efficient, requires somewhat less training data, and is competitive even in text applications such as clinical notes classification and symptom analysis, it makes it commonplace in disease prediction systems. Nevertheless, this strong assumption may not always hold in virtually any kind of disease because medical data may not show equal-feature independence. Neural Networks (NN) are machine learning methods capable of learning complex, nonlinear clinical data. They are inspired by the structure and function of the brain and consist of multiple interconnecting layers of neurons. They are most effective for disease prognosis since they can derive an array of hidden patterns and interactions among clinical factors. Artificial Neural Nets (ANNs) and Deep Neural Nets (DNNs) are well established in medical image analysis, predicting cancers, cardiology, and the monitoring of patients. Deep learning-type neural networks, such as Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs), have shown impressive results in image-based disease diagnosis, speech recognition, and the analysis of sequences of healthcare data. All forms of neural nets can boast of higher predictive accuracy and adaptability and are furthermore associated with wide-ranging tenability [17]. It goes, however, without saying that every good model naturally requires the processing of huge amounts of training data and provision of significant support for sophisticated computational resources. Additionally, they are characterized by their black-box controlled nature, which hinders widespread use for interpreting model output in clinical settings. Ensemble Learning Methods have garnered substantial attention, predominantly in disease prediction, where they amalgamate various machine learning models to produce superior model accuracy and robustness than singular classifiers. Ensemble methods enhance predictive performance by minimizing variance, bias, and influence their generalization capabilities. It encompasses widespread ensemble methods such as bagging, boosting, stacking, and voting classifiers. The Random Forest being an example of a bagging-based ensemble method; whereas boosting methods, including AdaBoost, Gradient Boosting, and XGBoost, train a series of weak learners in a sequential manner to improve classificatory. One of the major advantages of ensemble models is that they started in the health field, where we have complex and highly imbalanced datasets before ensembles." Ensemble models combining Logistic Regression, Decision Trees, Neural Networks, and Random Forests perform reasonably well in producing much better qualities as to reduce the generalization gap. More specifically, these methods significantly outperform in the areas of discrimination (accuracy), sensitivity, specificity, and good, if not strong, Area Under the Curve (AUC)—all qualifying them as optimum selections for clinical applications. Moreover, these ensemble algorithms work effectively through controlling overfitting and stabilizing the model. In spite of these advantages, ensembles, however, require higher computational time and complexity constraints. In totality, their definite functional aptness for disease prediction has effectively upended the entire paradigm standpoint of the actual disease prediction systems [18].

Integration of advanced algorithms centered on medical data is seen to lead to continued success, with the intelligent diagnosis of ailments, minimization of medical errors, provision of opportunities to infuse early warning signs of diseases, and maximize the quality of patient care. With the initial evolution of these technologies, forecasts suggest that the machine-learning-based disease prediction systems will subsequently play an important role within today's clinical scenes and precision medicine.

#### IV. FEATURE SELECTION AND FEATURE REDUCTION TECHNIQUES

Feature selection and feature reduction techniques, being crucial for creating a solid foundation in machine learning and data mining, help improve model effectiveness, reduce computational complexity, and eliminate irrelevant or redundant data. In scenarios like air quality prediction, healthcare analytics, fraud detection, and image processing, datasets have thus many characteristics that hardly contribute to predictive accuracy. The feature selection purely depends on these two; it is specified as the identification of important features from the original dataset. On the other hand, feature reduction transforms the data into a lower-dimensional space without losing significant information. Both of these processes will enhance prediction efficiency, reduce overfitting, and improve model interpretability. Filter-based feature selection methods evaluate the relevance of features using statistical measures decoupled from any machine learning algorithm. The methods turn out to be very computationally efficient and work best in high-dimensional settings. Some common methods involve correlation coefficient analysis [19], Chi-square test, Information Gain, Mutual Information, and ANOVA. The relevance score-based features are selected, and then the selected features are used to train the model. The training time is made smaller, making the method faster in execution and easier to scale. Any weakness is that groups of features may interact or have dependencies. Wrapper-based methods, however, act contrary to filter methods and make any decision by continuing to test varied subsets with respect to a specific machine-learning model. These methods are implemented through search methods such as forward selection, backward elimination, and recursive feature elimination, verifying which subset gives the maximum performance in terms of classification accuracy or error rate. Wrapper methods generally lead to better prediction accuracies as far as filter methods are concerned. The reason is they detect feature interactions; however, these other methods are highly computationally expensive and time-consuming, particularly where large datasets are to be processed, and too many attributes are found to handle. Principal Component Analysis (PCA) is one of the most widely utilized feature reduction techniques [20]. PCA transforms high-dimensional data into a small number of uncorrelated variables consisting of the principal components that explain maximum variance with minimal loss of information from the original dataset. PCA helps in reducing the dimensionality, interpreting complex data, and improving computation efficiency. PCA becomes somewhat problematic due to the formation of new transformed features, as interpretability tends to be compromised. Notwithstanding that, it may still prove to be very efficient given large and correlated features. Table 1 presents a comparative analysis of existing studies on machine learning-based disease prediction methods, highlighting the techniques used, application areas, and major findings of each research work.

**Table 1 Comparative Analysis of Existing Studies**

Ref. No.	Technique / Model	Application Area	Key Findings
[1]	Multi-perspective Machine Learning Ensemble	Heart Disease Prediction	Proposed an interpretable ensemble model with high prediction accuracy for cardiovascular diagnosis.
[2]	ML and Optimization Techniques	Healthcare Prediction	Discussed integration of machine learning and optimization methods for disease prediction systems.
[3]	Interpretable ML Framework	Alzheimer’s Disease Prediction	Developed imaging biomarker-based interpretable models for accurate diagnosis and prediction.
[4]	Hybrid ML and Optimization Models	Disease Prediction	Combined multiple ML techniques to improve disease prediction performance.
[5]	Machine Learning Algorithms	Medical Insurance Cost Prediction	Applied ML models for healthcare cost prediction and feature importance analysis.
[6]	Semi-Supervised Learning	Cardiovascular Disease Forecasting	Improved disease forecasting using self-learning semi-supervised approaches.
[7]	Binary Classification Evaluation	Cardiovascular Disease Prediction	Enhanced diagnostic capability through evaluation of binary classifiers.
[8]	Committee ML Methods with Grey Relational Analysis	Coal Calorific Value Prediction	Proposed a hybrid committee learning model with enhanced prediction accuracy.
[9]	Deep Learning and Graph Neural Networks	LncRNA-Disease Association Prediction	Achieved accurate disease association prediction using GNN-based deep learning.
[10]	AI-Based Medical Imaging Methods	Adipose Tissue Disease Analysis	Reviewed AI techniques for segmentation and disease prediction in medical imaging.
[11]	Hybrid ML with Feature Selection	Biomedical Disease Classification	Highlighted benefits of hyperparameter tuning and feature selection integration.
[12]	Boosting ML Algorithms	Composite Material Behavior Prediction	Improved prediction of material behavior using boosting-based ML models.

[13]	Explainable Machine Learning	Early-Onset Schizophrenia Prediction	Developed interpretable models considering primary and interaction effects.
[14]	ALADDIN ML Model	Liver Disease Prediction	Enhanced prediction of liver fibrosis using machine learning techniques.
[15]	Machine Learning for Real-World Data	Disease Prediction and Management	Systematic review of ML applications in healthcare data analysis.
[16]	Advanced Feature Engineering and Optimization	Cardiac Disease Prediction	Improved ML accuracy through optimized feature engineering methods.
[17]	ML and Deep Learning Techniques	Psychiatric Disorder Prediction	Reviewed EEG-based prediction methods using ML and DL models.
[18]	AI, ML, and Deep Learning	Disease Diagnosis and Prediction	Comprehensive review of AI-driven healthcare prediction systems.
[19]	Predictive Analytics with Big Data	Disease Diagnosis	Applied machine learning and big data analytics for healthcare prediction.
[20]	Ensemble ML Models with Class Balancing	Thyroid Disease Prediction	Improved thyroid disease prediction using ensemble techniques and balancing methods.
[21]	Optimized Machine Learning Algorithms	Heart Failure Prediction	Enhanced classification accuracy using optimized ML algorithms.

## V. CONCLUSION

In this review article, different feature selection and feature reduction techniques intended for machine learning and data analysis have been discussed comprehensively. The paper underscores the importance of selecting appropriate features and reducing the dimensionality of data, which improves forecasting accuracy, reduces computational complexities, and enhances overall model performance. Various methods such as filter, wrapper and hybrid techniques to select features, as well as Principal Component Analysis, Genetic Algorithm based feature selection, and novel deep learning methods as dimensionality reduction techniques have been discussed in detail. Every method is beneficial in one way or another depending on the dataset and application requirements for a given situation. It also illustrates that good feature engineering in high-dimensional datasets is of prime importance for intelligent systems to deliver results. These might be the topics of thorough research to deeply think about by designing frameworks for the processing of large- scale real-time data effectively in terms of automation and adaptability.

## REFERENCES

- [1] Miller, Sean T., et al. "Multi-perspective machine learning MPML: A high-performance and interpretable ensemble method for heart disease prediction." *Machine Learning with Applications* (2026): 100836.
- [2] Gauthaman, Shruthi A., et al. "Introduction to machine learning and optimization in healthcare." *Integrative machine learning and optimization algorithms for disease prediction*. IGI Global Scientific Publishing, 2026. 1-34.
- [3] Kang, Wenjie, et al. "An interpretable machine learning framework with data-informed imaging biomarkers for diagnosis and prediction of Alzheimer's disease." *Computerized Medical Imaging and Graphics* (2026): 102722.
- [4] Giriprasad, S., et al. "Integrating Multiple ML and Optimization Techniques for Superior Disease Prediction." *AI Model Design and Data Management for Disease Prediction*. IGI Global Scientific Publishing, 2026. 35-58.
- [5] Liu, Yanhang, and Shuqian Zhang. "Application of Machine Learning Algorithms for Medical Insurance Cost Prediction and Feature Analysis." *Emerging Markets Finance and Trade* 62.1 (2026): 94-110.
- [6] Tusher, Ekramul Haque, et al. "Semi-supervised learning: Assisted cardiovascular disease forecasting using self-learning approaches." *Journal of Advanced Research in Applied Sciences and Engineering Technology* (2026): 136-150.
- [7] Eltawil, Mohamed, et al. "Comment on Iacobescu et al. Evaluating Binary Classifiers for Cardiovascular Disease Prediction: Enhancing Early Diagnostic Capabilities. J. Cardiovasc. Dev. Dis. 2024, 11, 396." *Journal of Cardiovascular Development and Disease* 13.1 (2026): 46.
- [8] Lawal, Abiodun Ismail, et al. "A novel grey relational analysis-based committee of machine learning methods for enhanced prediction of coal calorific value." *Fuel* 406 (2026): 137070.
- [9] Farrokhi, Zahra, Jamshid Pirgazi, and Ali Ghanbari Sorkhi. "An effective deep learning and graph neural network approach for accurate prediction of LncRNA-disease associations." *Biomedical Signal Processing and Control* 112 (2026): 108431.
- [10] Alghamdi, Mohammed J., Muhammad Rashid, and Muhammad Arif. "From Segmentation to Disease Prediction: A Systematic Review of AI Methods for Adipose Tissue Analysis in Medical Imaging." *IEEE Access* 14 (2026): 9729-9757.

- [11] Dhanka, Sanjay, et al. "Advancements in hybrid machine learning models for biomedical disease classification using integration of hyperparameter-tuning and feature selection methodologies: A comprehensive review." *Archives of Computational Methods in Engineering* 33.1 (2026): 289-324.
- [12] Tariq, Aiman, Ayşe Polat, and Babür Deliktaş. "Boosting machine learning algorithms for predicting the macroscopic material behavior of continuous fiber reinforced composite." *Journal of Reinforced Plastics and Composites* 45.3-4 (2026): 781-799.
- [13] Lin, Chih-Wei, et al. "Exploring Primary and Interaction Effects of Minor Physical Anomalies: Development and Validation of Prediction Models Using Explainable Machine Learning Algorithms for Early-Onset Schizophrenia." *Schizophrenia Bulletin* 52.1 (2026): sbaf016.
- [14] Alkhouri, Naim, et al. "ALADDIN: a machine learning approach to enhance the prediction of significant fibrosis or higher in metabolic dysfunction-associated steatotic liver disease." *Official journal of the American College of Gastroenterology| ACG* 121.2 (2026): 362-374.
- [15] Alhumaidi, Norah Hamad, et al. "The use of machine learning for analyzing real-world data in disease prediction and management: systematic review." *JMIR Medical Informatics* 13.1 (2025): e68898.
- [16] Bilal, Omair, et al. "Boosting Machine Learning Accuracy for Cardiac Disease Prediction: The Role of Advanced Feature Engineering and Model Optimization." *The Review of Socionetwork Strategies* 19.2 (2025): 271-300.
- [17] Abir, Shake Ibna, et al. "Machine learning and deep learning techniques for EEG-based prediction of psychiatric disorders." *Journal of Computer Science and Technology Studies* 7.1 (2025): 46-63.
- [18] Sadr, Hossein, et al. "Unveiling the potential of artificial intelligence in revolutionizing disease diagnosis and prediction: a comprehensive review of machine learning and deep learning approaches." *European journal of medical research* 30.1 (2025): 418.
- [19] Chinta, P. C. R., et al. "Predictive Analytics for Disease Diagnosis: A Study on Healthcare Data with Machine Learning Algorithms and Big Data." *J Cancer Sci* 10.1 (2025): 1.
- [20] Zhong, Jiachen, and Yiting Wang. "Enhancing thyroid disease prediction using machine learning: A comparative study of ensemble models and class balancing techniques." (2025).
- [21] Ahmed, Marzia, et al. "Predicting the classification of heart failure patients using optimized machine learning algorithms." *IEEE Access* (2025).
- [22] Siddique, Md & Bin Seraj, Md. Masrafi & Adnan, Md Nasim & Galib, Syed. (2024). Artificial Intelligence for Infectious Disease Detection: Prospects and Challenges. 10.1007/978-3-031-59967-5\_1.