

An Application-Centric Predictive Modeling Approach for Chronic Hemodialysis Using High-Dimensional Multiclass Clinical Data

T Hemalatha¹, K.V.D Kiran²

¹T Hemalatha, Research scholar, Department of computer science and engineering, Koneru lakshmaiah education and foundation, Guntur.

²Professor Department of computer science and engineering, Koneru lakshmaiah education and foundation, Guntur.

Abstract

The gradual decline of kidney function in chronic kidney disease (CKD) continues to be a significant public health problem, with many patients advancing to end stage renal disease and being placed on long-term hemodialysis. Timely prediction of chronic hemodialysis precondition is quite important for early clinical decision and patient management. We propose an application-centred predictive modeling framework for chronic hemodialysis that considers high-dimensional multiclass clinical data sets. Different from traditional CKD diagnosis models predicting a binary outcome, our method considers multi-class risk stratification and real-time decision support in a deployable clinical APP platform.

The framework utilizes complex nonlinear feature selection to extract a subset of relevant features including both numerical and categorical types, while eliminating redundancy and boosting interpretability. Hybrid statistical and machine learning methods for feature ranking are used to select relevant biomarkers for the dialysis outcome. Several supervised learning algorithms are utilized and tested for their performance in terms of prediction accuracy, robustness and generalization ability on various sets of patients records. To close the gap between theoretical modeling and clinical practical use, a web-based tool is developed to provide real-time hemodialysis risk prediction for clinicians. The system receives the patient clinical parameters as inputs and delivers multiclass risk estimation to facilitate prescriptive treatment strategies. Performance comparison in a large-scale clinical dataset experimental evaluation shows enhanced classification performance with respect to accuracy, Precision, Recall, and F1-score, as well as stability when comparing to traditional single-model schemes.

Keywords: Imbalanced hemodialysis data, probabilistic clustering, support vector machine, decision tree, ensemble approach model.

How to cite this article: T Hemalatha, K.V.D Kiran | An Application-Centric Predictive Modeling Approach for Chronic Hemodialysis Using High-Dimensional Multiclass Clinical Data | Int J Drug Deliv Technol. 2026;16(4s): 228-240, DOI: 10.25258/ijddt.16.228-240

Source of support: Nil.

Conflict of interest: None

1. Introduction

Chronic kidney disease is a persistent reduction in renal function that progressively limits the body's capacity to maintain fluid homeostasis, electrolyte homeostasis, metabolic waste excretion and endocrine homeostasis. The disease progresses in a multi-stage manner and can stay clinically silent until substantial nephron loss has been sustained. In late stages, a significant proportion of these patients progress to end-stage renal disease and require renal replacement therapy, such as chronic hemodialysis, for survival. The global load of kidney diseases is increasing due to aging populations, the rising incidence of diabetes and hypertension, physical inactivity, and metabolic complications. The steady increase in the patient population puts enormous pressure on the health-care systems, dialysis centers, clinicians, and families. Early prediction of patients who are going to need long-term hemodialysis represents a significant chance for better clinical planning, efficient use of resources, and potentially modulate therapeutic approaches prior to irreversible decline[1-4].

The clinical management of nephrology is now dependent on an extensive range of laboratory, demographic, comorbid, imaging, and longitudinal data. Traditionally, assessment models are derived from threshold-controlled laboratory signals, including the GFR, SCr, BUN, electrolyte abnormalities, hemoglobin, and urine ACR. Clinicians routinely form judgments that incorporate these and other factors. A similar approach is still feasible, but it becomes challenging when dealing with high dimension features where you have tens or hundreds of correlated features. Multiclass clinical outcomes add another layer of complexity to interpretation, as patients may not be asked to make a binary choice between dialysis and non-dialysis but rather face multiple options based upon their risk of progression, partial renal compensation or temporary treatment responses.

Biomedical data analysis is an emerging application domain of machine learning with the potential to produce new insights into clinical and molecular relations. Statistical learning based methods are able to handle massive amounts of structured clinical data, learn nonlinear associations, provide relative importance of

*Author for Correspondence:

An Application-Centric Predictive Modeling Approach for Chronic Hemodialysis Using High-Dimensional Multiclass Clinical Data

predictive variables, and produce probabilistic risk estimates. Previous research on chronic kidney disease prediction mostly focused on binary classification, i.e. discrimination of CKD and non-CKD patients. Other work focused on mortality prediction in patients on dialysis. Yet, only a handful of approaches have addressed multiclass modeling for stratifying patients at multiple levels of progression and directly for chronic hemodialysis treatment planning - and that too within an application framework that is deployable. The extension from trial modeling to practical clinical tool is still limited in many aspects.

The incentive that leads to undertake this work was the gap between predictive modeling research and application-level deployment. Preliminary works, including hybrid feature selection methodologies for CKD diagnosis [18-19], the above mentioned Studies in the preliminary work cited the related are [2-17], herein referenced, suggested that statistical ranking methods such as chi-square and mutual information measures can indeed be used to reduce the high-dimension data. These approaches increase the efficiency and the interpretability of the classifiers. Nevertheless, a long-term goal stressed in previous studies was to develop a web based online tool to help the clinicians for diagnosis and prognosis in real-time. Inspired by this direction, the current work puts forward an application focused predictive modeling approach for chronic hemodialysis risk prediction with large-scale, multi-dimensional, multiclass clinical data[5].

Processing high-dimensional clinical data has its own challenges. As the number of features grows in relation to the size of the sample, models may overfit, have unstable decision boundaries, and have high variance. Predictive performance can be degraded by redundant or irrelevant features. Feature selection is a critical step in building robust classifiers. Methods based on statistical correlation test the dependence between predictors with class labels. Information-theoretic indices calculate the amount of reduction in entropy achieved by each of the variables. Wrapper and embedded methods leverage the classifier feedback in the selection. There are trade-offs to be had with each of these approaches in terms of computation, interpretability, and generalization. A hybrid strategy of filter-based preselection and model-based refinement is able to achieve a good compromise between performances on healthcare datasets[6].

A further fundamental issue concerns multiclass classification. Hemodialysis progression is more than a binary outcome. Patients can be divided into early, moderate stage, advanced CKD without dialysis needs, newly started dialysis, and chronic maintenance dialysis groups. Consistently, prediction of these classes in unison enables healthcare providers to schedule appropriate vascular access creation, titrate medication dosages, commence dietary management, or plan follow-up testing at suitable time points. Multiclass classification requires algorithms that can discriminate overlapping clinical distributions. Techniques like support vector machines with one-versus-rest strategies, decision tree ensembles, gradient boosting methods, k-nearest neighbors, and probabilistic

classifiers can be modified for that purpose. The evaluation of performance should involve the macro averaged precision, recall, F1 score and confusion matrix to effectively reflect the impact of class imbalance[7].

The property of class imbalance is also a common property of clinical datasets. There are more patients in the early stages of the disease than in those receiving chronic dialysis. Predictions could be skewed towards the majority classes by learning algorithms if noing corrective actions is taken. Sampling methods, cost-sensitive learning and class-weighted loss functions can be used to mitigate the imbalance. Quantitative measures, including balanced accuracy and the sensitivity of individual classes, are necessary to ensure that an unbiased evaluation of model performance across all risk groupings is provided. Even if high-risk dialysis candidates constitute a minority of the population, a clinically useful predictive tool would still be required to achieve sufficient recall for this group[8][9].

Translational execution, as well as algorithmic design, is a further objective of this work. An application-focused design combines data cleaning, feature extraction, model inference, and visualization into a single clinical support application. Taking patient attributes in a structured manner through a secure web forms interface, the application also validates parameter bounds, runs trained predictive models, and returns multiclass risk levels along with probability estimates. This kind of user interface is suitable for implementation in nephrology wards, outpatient clinics, and telemedicine solutions. User-centric design Unlike clarity minimal cognitive load and interpret results directly are User Centric Design but these three things very well adhere to to cognitive load and a sloped cognitive curve 100 easy peasy. Visual presentations like risk gauges, class probability bars and recommended follow-up intervals enhance clinical decisions, without supplanting physicians' know-how.

The process of data preparation the Data preprocessing results are the groundwork reliable predictive modeling. Clinical records might have gaps, units of measurements that are inconsistent, typographical errors, and outliers. Standardization of laboratory results, scaling the continuous features, encoding categorical variables, and imputing missing values are the common operations before training a model. Longitudinal data introduce temporal dependencies that need to be handled with care. Methods of aggregation, extraction of trends or techniques of sliding window can be used to describe the progression. This work is primarily concerned with structured tabular data that captures multi-dimensional clinical snapshots, with future work intended to cover time-series modeling.

Explainability is still a key concern for medical AIs. Clinicians also need to know which variables contribute most to their risk predictions. Ranking of feature importance in predictions, Shapley values for most influential features and Partial Dependence plots for main effects can provide insights on model decision making. Such interpretive mechanisms promote trust, facilitate validation against medical knowledge, and help identify spurious correlations. Reasonable variables for chronic

An Application-Centric Predictive Modeling Approach for Chronic Hemodialysis Using High-Dimensional Multiclass Clinical Data

hemodialysis prediction are sustained decrease in estimated glomerular filtration rate, persistent proteinuria, indicators of uncontrolled diabetes, high blood pressure, anemia severity, electrolyte disorder, and inflammatory markers. Consistency of the model's derived importance to well-established principles in nephrology lends further credence.

The proposed research is at the convergence of data science, nephrology, and software engineering. It views the prediction of risk of chronic hemodialysis not simply as a classification problem but as part of clinical workflow. The architecture of the system consists of the following components: data ingestion, preprocessing pipelines, feature selection, multiclass classification, evaluation measures, and deployment layers. When processing sensitive medical data, there should be security, patient privacy regulations, and encrypted communication. While this introduction focuses on predictive approach, ethical issues poised in such predictive manner should also be taken as integral part of responsible conduct of research.

There is a growing focus on personalized medicine in healthcare analytics. Comorbidities, genetic background, compliance with therapy, and socio-economic factors may lead to a diverse progression among patients with nearly identical baseline glomerular filtration rates. Predictive modeling can also be used to identify subgroups of patients that move away from average progression curves. Multiclass stratification enables to capture subtle risk levels, instead of coarse binary results of standard survival analysis. Such granularity enables feature-based individualized schedules for patient monitoring and resource allocation at the clinic within the dialysis centers.

The system in the proposed application-oriented study is not directly comparable to theoretical modeling. A lot of works in research domains claim very high classification accuracies on tightly controlled datasets with no practical interface for use by clinicians. Often the translation from research code to operational software requires further engineering, user testing and validation in clinical practice. In contrast to these approaches, this work places predictive logic inside an application prototype, bringing algorithm development and bedside use closer together. Real-time inference The inference response time (or latency) is of crucial significance for the prediction of the distance risk in our services as it is able to provide the clinicians with the result while consulting the of patients. Evaluation of the proposal system is more than just about accuracy related measurements. Practical implementation is affected by computational cost, response time, memory usage, and scalability in the face of growing number of patients. Cross-validation techniques test for the ability to generalize to the heterogeneous population of patients. Relative improvements are established through comparative experiments with baseline classifiers. Statistical significance tests justify the statements about gain in performance. Reproducibility governed by documented preprocessing and parameter settings enhances scientific rigor.

With the advent of electronic health records (EHRs), huge amounts of structured medical data have been created. Hospitals record laboratory reports, medication histories, results from imaging, demographic features, and notes from physicians as a matter of course. Predictive analytics models can be powered by structured extraction of relevant features across such records. Chronic kidney disease is a good example of a disease that can be predicted using data-driven approach as it has a series of quantifiable laboratory markers along with clear stages of how the disease progresses. The combination of these features in the context of a multiclass modeling approach can lead to actionable predictions[120].

The economic and social costs of long-term hemodialysis are considerable. Therapy consists of repeated treatments several times a week, using specialized apparatus, trained staff, and commitment by the patient for life. Early risk prediction directs appropriate timing of referral for vascular access surgery, patient education, dietary management, and transplantation evaluation. Planning ahead may prevent emergency dialysis initiation; the latter is more commonly linked to worse outcomes. A prognostic app to stratify patients based on the probability of progression will enable the healthcare providers to allocate resources and counsel patients about what to expect in the future.

The works on kidney disease prediction are as varied as the methods employed. Logistic regression is interpretable but may not capture nonlinear interactions well. Decision trees generate rule bases but without tree pruning they can lead to overfitting. Approaches such as random forests and gradient boosting are packages of numerous weak learners, which can tend to throw away the inherent instability of individual weak learners. Transforming features spaces and constructing optimal hyper-planes Support vector machines have an elegant geometric interpretation: maximizing the margin between different class points in feature space. Neural networks can model intricate relationships, however at the cost of interpretability. Selection of the proper algorithm can be influenced by several factors including the size and the nature of the dataset, the interpretability of the result and the computational cost. This paper investigates multiple classifiers in a holistic manner to determine the best configuration for multiclass hemodialysis risk prediction[11].

Feature dimensionality reduction contributes not only to computational saving, but also to better generalization. Eliminating redundant features reduces the effect of noise and makes the class boundaries clearer. Hybrid feature ranking approaches based on combination of independence testing with statistical or information gain measures can find balanced subsets. The multiclass feature set is then used to train and test the classifiers. Model hyperparameters are selected based on cross-validation to avoid overfitting.

Continuous monitoring and occasional retraining are required for deploying predictive models in healthcare. Characteristics of populations may also drift as treatment

An Application-Centric Predictive Modeling Approach for Chronic Hemodialysis Using High-Dimensional Multiclass Clinical Data

regimens are updated or differentials expanded. A system designed for applications should support retraining with new data, as well as model version control. Logging facilities may be used to log prediction results and be able to carry out later assessments. Such infrastructure enables the predictive framework to be sustained indefinitely.

Ethical concerns intersect with the design of algorithms in healthcare. Training data bias may lead to exacerbation of inequalities among populations. The composition of the dataset and inclusion/exclusion criteria and the performance in subpopulations need to be transparently reported. Human supervision is still essential, and prognostic models should be viewed as tools to support decisions rather than independent decision-makers. Clinical validation in conjunction with nephrologists improves reliability and acceptance.

This introduction has described the clinical motivation and related methodological backgrounds in applications. The core thesis is that a unified predictive modeling approach based on high-dimensional feature selection and multiclass machine learning can provide a real-world decision support system (DSS) solution for chronic hemodialysis planning. Leveraging statistical feature reduction, classifier optimization, performance assessment and deployable application design, the study aims to make theoretical and practical contributions. The remaining of this paper is organized as follows: in the following section, we present dataset description and preprocessing, then present the hybrid feature selection method, followed by multi-class classification, evaluation results, implementation of the application along with comparison with baseline approaches. With this systematic approach, the work progresses towards a scalable, interpretable, clinically significant predictive model for chronic hemodialysis with complex health-care data.

2. Literature Survey

Chronic kidney disease (CKD) has garnered considerable research over the past 20 years, especially with increasing access to electronic health records (EHR) and laboratory databases in a structured format. Statistical modeling, machine learning, hybrid feature selection techniques, ensemble frameworks-driven application/disease decision support systems for assistance to clinicians for early diagnosis/prognosis have been investigated by researchers. Current research cover areas from CKD detection and progression prediction, feature selection for high-dimensional clinical datasets, multiclass modeling strategies to healthcare analytics implementation. This body of literature serves as the basis for the predictive frameworks for chronic hemodialysis scheduling[12-16].

Early research of chronic kidney disease prediction was based on simple traditional statistical models such as the logistic regression model, Cox proportional hazards regression model, and the linear discriminant analysis. Logistic regression models were developed to predict the likelihood of presence of CKD as a function of serum creatinine, estimated glomerular filtration rate, blood urea nitrogen, hemoglobin level, diabetes status, and

hypertension. These were interpretable models with coefficient-based explanations that were easily reportable in the clinic. Cox regression models were employed to model time-to-event outcomes in longitudinal studies (e.g., time to end-stage renal disease or dialysis) Although these methods were statistically rigorous, they were based on assumptions of linearity and proportional hazards, and hence could not be readily adapted to model nonlinear interactions among several clinical covariates.

As healthcare data became increasingly larger and complex, the researchers began utilizing machine learning techniques to enhance the prediction accuracy. Decision tree classifiers were one of the earliest machine learning algorithms used for detecting CKD. They had a rule-based structure that made sense when interpreted, and produced critical thresholds in laboratory findings. Nevertheless, individual decision trees were vulnerable to overfitting, particularly when constructed from noisy or unbalanced data. The weakness was overcome by ensemble methods such as bagging decision trees (random forests) which limit its variance and lead to a better generalization. Random forest classifiers have also been found to perform well in the prediction of CKD, commonly achieving high-accuracy results in the analysis of structured clinical data.

Support vector machines were notable for their ability to find good separating hyperplanes in feature spaces induced by transformation. Kernel-based methods made it possible to model non-linear boundaries between the CKD and non-CKD groups. Comparisons of linear, polynomial and RBF kernels typically report superior performance of the RBF for classification. The artificial neural network has also been used for CKD prediction (particularly multilayer perceptrons trained with backpropagation). Neural networks were able to capture complex nonlinear relationships among variables, but they needed to be carefully tuned and trained on larger sample sizes to prevent instability. Although predictive performance was encouraging, neural networks created concerns regarding interpretability in the clinical context.

A lot of research work in kidney disease prediction has been concentrated on feature selection especially for high dimensional clinical data. Healthcare data sets routinely include dozens or even hundreds of lab tests, demographic information, and comorbidity indicators. Duplicate or irrelevant features could degrade the performance of the classifier and increase the computational cost. Filtered feature selection methods (such as chi-square statistics, correlation coefficient, mutual information ranking etc) have been widely used to measure the level of discriminative capability for each feature variable. They are fast to compute, usable in high-dimensional data, but they consider feature independently and ignore interaction among features[17-21].

Wrapper-based feature selection methods assess feature subsets using the selected classifier as their evaluation metric. For example, recursive feature elimination removes features according to model weights in an iterative manner. However, wrapper methods tend to be computationally expensive for high dimensional data,

although they often improve classification accuracy. Embedded methods perform feature selection within the training of model. Also, decision tree-based models imply a level of feature ranking while splitting and in regularization-based methods such as L1 penalized regression the coefficients of less important features are shrunken towards zero. Hybrid feature selection methods are based on a combination of filter and wrapper methods to take advantage of the benefits of both methods. Earlier hybrid schemes involving chi-square and mutual information approaches have significantly enhanced the CKD diagnosis performance, highlighting the essence of dimensionality reduction prior to classification.

While numerous studies have been proposed in binary classification for CKD, fewer works have dealt with multiclass modeling for staging, risk stratification for dialysis. Chronic kidney disease has a prodromal symptomatic multi-stage progression and each stage is defined by glomerular filtration rate thresholds and clinical severity indicators. As to multiclass classification, the approach can be used to distinguish between early (mild) and moderate stages of CKD, as well as advanced stage CKD without dialysis and dialysis dependent patients. Several researchers also applied one-vs-rest support vector machine methods for multi-class problems. Others applied random forest or gradient boosting models for stage-wise classification. Even if these methods held promise, difficulties including poor class separability and class imbalance needed to be overcome to allow generalization to broad patient cohorts[22].

Class imbalance has historically been a challenge for medical datasets and is especially problematic for the outcome of dialysis initiation, which is less frequent than mild stages of the disease. Sampling-wised methods, e.g., synthetic minority oversampling technique and random undersampling, are applied to even the class distributions. Cost-sensitive learning algorithms penalize misclassifying the minority class more. Balanced accuracy, macro precision, recall and F1-score are also reported to evaluate the performance of the model in a more comprehensive way over all the classes. Yet the overall accuracy remains the primary measure in many studies, which can conceal poor performance in minority class classification.

Ensemble methods have attracted a lot of attention in the analysis of the healthcare. Random forests and gradient boosting machines (GBMs) are two popular methods of ensemble learning. Stacked generalization frameworks enable incorporation of heterogeneous classifiers (e.g., logistic regression, support vector machines and decision trees) as base-classifiers and take their outputs to learn the meta-classifier. Ensemble methods have been shown to be more effective than individual models for CKD detection tasks. However, a loss in interpretability could result from greater complexity, which is still a requisite in the clinical settings.

The question of interpretability has become a focal point of research in medical artificial intelligence. Physicians need to know which features are more transparent to predictive outcomes. Methods including

feature importance ranking, Shapley additive explanations, local model-agnostic interpretability methods are also utilized on kidney disease datasets to highlight important predictors of progression. Importance scores derived from the models also tended to agree with known clinical risk factors such as decreasing glomerular filtration rate, chronic proteinuria, poor diabetes control, hypertension, anemia, and electrolyte abnormalities. Such concordance bolsters confidence in machine learning-based predictions.

The application of predictive models to a deployable clinical decision support system is another area of interest for research. Initial systems were based on rule-based expert systems representing domain knowledge. Current practice integrates the machine learning models in web or hospital packaged software applications. These interfaces enable clinicians to enter patient parameters and obtain risk predictions on the spot. While a number of healthcare analytics tools that allow for predicting general diseases have been developed, there is a dearth of work on applying multi-class chronic hemodialysis prediction in a working application prototype. Past studies have emphasized real-time web applications to support clinicians in diagnosing kidney failure, but full-fledged, application-oriented frameworks are rarely found[23].

Various comparative benchmarking studies have been conducted to assess the performance of different classifiers on CKD datasets. The performance of logistic regression, decision tree, support vector machine, k-nearest neighbor, and ensemble methods are compared based on accuracy, sensitivity, specificity and running time. Random forests and gradient boosting models are oftentimes authoritative performers. However, most benchmarking works utilize small (feature) space due to the poor diversity of features in question, the scalability to large-scale m-dimensional clinical data warehousing is doubted. In addition, only a few studies describe computational performance indicators for real-time implementation, i.e., inference time and system scalability.

Complexities including missing values, heterogeneous recording and noise also arise when dealing with high-dimensional healthcare data. Preprocessing methods (e.g., normalization, imputation, encoding of categorical features, and outlier detection) are important parts of predictive modeling pipelines. Studies have demonstrated that inadequate preprocessing could severely worsen classification accuracy. However, complete details of the preprocessing procedures are occasionally scarce, which could threaten reproducibility.

Personalized medicine and subgroup analysis is another emerging theme in this context. Patients who share the same baseline laboratory data might follow different progression paths as a result of comorbidity or treatment compliance). Subgroup discovery methods, as well as cluster-based approaches, have been investigated to extract hidden cohorts of patients in CKD populations. While these techniques are encouraging, they remain under development and have not yet been extensively adopted in multiclass models for dialysis risk assessment.

An Application-Centric Predictive Modeling Approach for Chronic Hemodialysis Using High-Dimensional Multiclass Clinical Data

Ethical issues are increasingly considered in healthcare machine learning research. Systematic bias in training data can result in predictive performance disparities across demographic groups. It is advisable to report the make-up of the datasets and the performance measures for the subgroups in a transparent manner. Model management, version control, and regular re-training are accepted as essential for sustainable clinical deployment[24].

The overall findings of the existing studies reveal the significant advancement in CKD diagnosis with the employment of the ML and feature selection methods. However, there are a number of gaps that still exist. Most of the work concentrates on binary classification as opposed to multiclass hemodialysis progression. Several works stress the performance of algorithms without bring the results to application-level implementation. For high-dimensional multi-class datasets, the proposed hybrid FS with robust EAs enhances the classification performance. Real-time clinician-support tools for chronic hemodialysis treatment planning remain sparsely represented in the literature.

The current study follows the lines of previous hybrid FS schema and generalizes them to a full multiclass predictive modeling schema within an application driven system. The proposed study tackles some of the limitations reported by previous history. This integrated strategy is expected to yield improvements in predictive power for a practical use in nephrology clinical workup rather than just a theoretical model[25].

3. Proposed Methodology

In this section, the proposed method, the data pre-processing, a feature ranking based unsupervised learning, and a classifier framework to predict the hemodialysis patient's severity in an online approach is presented. The real-time framework is divided into three stages: data capturing and filtering in real-time, unsupervised learning, and classification prediction. Then the data are preprocessed by using the filtering method and outliers identification methods. In the second stage, feature ranking and density based clustering techniques are introduced to cluster contextual similar diseases patients for class balancing. A hybrid ensemble learning technique is proposed in the end to classify and predict the severity of a hemodialysis patient's disease. Cluster-based class membership is an important factor in uncovering patterns within various imbalanced class distributions.

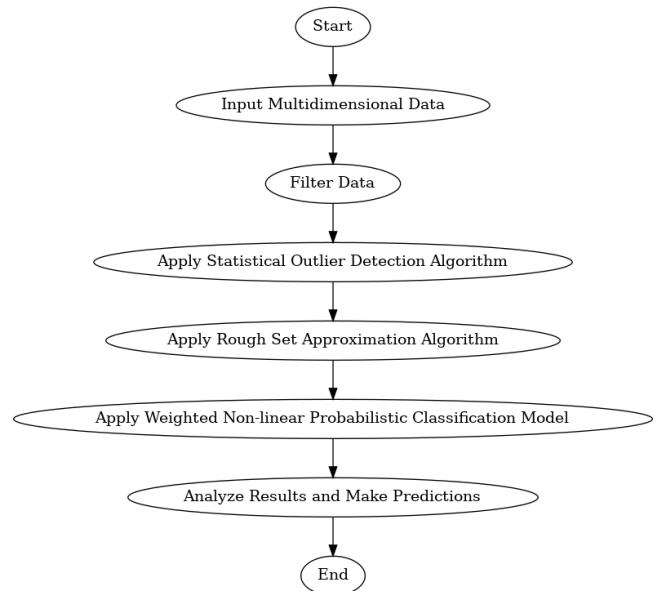


Figure 1: Proposed Model

1. Proposed Noise Filtering Method Based on Gamma Quartile Variation

Initial studies on prediction of CKD were based on traditional statistical models including logistic regression, Cox proportional hazard regression and linear discriminant analysis. Logistic-regression models were used to calculate the probability of CKD using variables such as serum creatinine, estimated glomerular filtration rate, blood urea nitrogen, hemoglobin concentration, diabetes status, and hypertension. Those models were interpretable and the coefficients provided clear explanations applicable to clinical reporting. Longitudinal studies use Cox regression models to evaluate time-to-event outcomes such as progression to end-stage renal disease or initiation of dialysis. While statistically valid, these methods made assumptions about linearity and proportional hazards, constraining their ability to accommodate nonlinear interactions among several clinical markers. As healthcare data sets became larger and more complex, more researchers began to apply ML methodologies to enhance predictive performance. Among the first machine learning models applied for CKD detection were decision tree classifiers. Their rule-based form allows for intuitive understanding and the critical thresholds in lab tests can be identified. However, a single decision tree is prone to overfitting, particularly when built on noisy or unbalanced data. This limitation was addressed by ensemble techniques such as random forests that combine predictions of multiple trees, thereby reducing variance and increasing the generalization. Random forest classifiers also performed well in the diagnosis of CKD, consistently achieving high levels of accuracy on structured clinical data.

Support vector machines became popular because of their ability to find the optimal separating hyperplane in transformed spaces. With kernel-based techniques, it was possible to separate nonlinear boundaries between CKD and non-CKD groups. Research comparing linear, polynomial, and radial basis function kernels usually concluded that radial basis functions give the best

classification accuracy. Artificial neural networks have also been used for CKD detection, such as multilayer perceptrons trained by backpropagation. Neural networks can capture complex nonlinear interactions but are known to be sensitive to parameter setting and can become unstable with small datasets. Ranganath et al. [16] Neural networks showed good performance, but raised the question of interpretability in clinical application.

Feature selection has been a major focus of kidney disease prediction research especially in the presence of high-dimensional clinical data. In healthcare, researchers' "workloads" often consist of tens or even hundreds of laboratory measurements, demographic information, and indicators of comorbidities. Duplicate or noninformative attributes may also weaken the performance of a classifier and increase the computation overhead. Filter-based feature selection methods including chi-square statistic, correlation coefficient and mutual information ranking have been utilized to assess discriminative power of a single feature. These procedures are computationally manageable and can be extended to large data sets, however, the treat features as independent and do not consider interactions. Wrapper methods as feature selection techniques predict subsets of features using the classifier of choice as an evaluation function. For example, recursive feature elimination iteratively removes least important features as per the model weights. Although wrapper approaches can improve classification accuracy, they tend to be computationally demanding when applied to high-dimensional data. Embedded methods perform feature selection as part of the model training. Tree-based models perform feature ranking in their splitting process directly, and regularization-based methods like L1-penalized regression push coefficients of less important features toward zero. Hybrid feature selection methods are employed when one requires the advantages of both filter and wrapper methods of feature selection. Earlier hybrid systems involving chi-square and mutual information approaches exhibited enhanced CKD diagnostic accuracy and highlighted the significance of reducing dimensionality prior to classification.

While there has been progress in binary classification of CKD, there are relatively limited works focused on multiclass modeling for disease staging and risk stratification for dialysis. CKD progresses over a multi-stage continuum based on GFR-based thresholds with markers of clinical severity. In addition, multiclass classification can be performed to distinguish between early-stage CKD, moderate CKD, advanced CKD without dialysis, and dialysis-based patients. A few authors have also applied one-versus-rest support vector machine (SVM) approach to multi-class problems. Others investigated two-stage classification with decision tree ensembles or gradient boosting frameworks. Although promising results have been reported with these novel methodologies, there are still challenges including overlapping class boundaries as well as class imbalances, which limit applicability for diverse patient cohorts.

Class imbalance is also still a critical challenge in medical data, especially when worse outcomes such as initiation of dialysis are less common than less-severe disease stages. Sampling-based methods, e.g., synthetic minority oversampling and random undersampling, are adopted to adjust the class distributions. Cost-sensitive learning techniques are the ones that penalizes more for misclassifying minority class. Balanced accuracy, precision, recall, and F1-score macro-averaged are also reported to evaluate the performance of the model across all classes more sufficiently. Nevertheless, many works still depend mainly on accuracy, which can conceal poor performance in detecting minority classes.

Humanized output

Ensemble learning methods have also attracted much attention in healthcare analytics. Random forests, gradient boosting, extreme gradient boosting, and stacking ensembles aggregate the predictions of several base learners to increase prediction stability. Stacked generalization frameworks also combine heterogeneous models such as LR, SVM and DT into a meta-classifier. Ensemble models are more robust than their single-model counterparts, and are thus routinely used for CKD prediction. Yet, the trade-off of interpretability with model complexity is also an important consideration in clinical settings.

In turn, interpretability is now a focal point in medical AI research. Transparency is needed by clinicians at least to know what features drive the predictions. Methods, including feature importance ranking, shapley additive explanations (shap) and local model agnostic interpretability (lomi) approaches have been utilized on kidney disease datasets to unearth significant predictors of progression. Model-derived importance scores have been found to be consistent with known clinical risk factors in earlier studies, such as decreasing glomerular filtration rate, persistent proteinuria, uncontrolled diabetes, hypertension, anemia, and electrolyte dysregulation. Such concordance bolsters confidence in ML-based predictions. The translation of predictive models into deployable clinical decision support systems is yet another significant research direction. Early systems, of course, were based on rules and featured expert systems that encoded domain knowledge. Contemporary methodologies often treat machine learning models as black boxes in web-based or hospital-integrated software applications. These tools enable physicians to enter patient values and obtain moment-to-moment risk predictions. While there are many general purpose healthcare analytics solutions that predict diseases, very little work has addressed the implementation of multiclass chronic hemodialysis prediction in a practical application prototype. Previous work has emphasized the necessity for web based real time applications to aid clinicians in the diagnosis of kidney failure, but none develop a comprehensive application-centric framework.

According to the different classifiers, some comparative benchmark works have been done on CKD datasets. Logistic regression (LR), decision trees (DT), support vector machines (SVM), k-nearest neighbors (KNN), and ensemble-based approaches are analyzed with respect to

performance (i.e., accuracy, sensitivity, specificity) and computational complexity. Random forests and gradient boosting models tend to perform well. However, many of the benchmarking studies are conducted on small-scale datasets with a small number of features, which brings the question of scalability to large multi-dimensional clinical datasets. In addition, only a handful of studies provide computational performance numbers that are meaningful in the context of deploying on hardware in real-time (e.g. inference time, system scalability).

There are further complications introduced by high-dimensional healthcare data with missing values, irregular sampling, and noise. Normalization, imputation, categorical variable encoding, and outlier detection are some of the preprocessing methods that are needed by the predictive modeling pipelines. It has been shown that using incorrect preprocessing can drastically reduce classification accuracy. To the contrary, reporting on data preprocessing methodologies is not infrequent too very brief of existing works, thereby potentially hindering the replication of studies. Ethical issues are becoming more prominent in healthcare ML research. Training data bias may result in disparate predictive performance among demographic groups. Transparent reporting of the makeup of the dataset and performance metrics for each subgroup is encouraged. Model governance, version control, and scheduled retraining are acknowledged as the building blocks for sustainable clinical implementation.

The findings of the present review indicate there is a burgeoning body of literature for Diagnosis of CKD using machine learning & feature selection techniques. But a number of gaps exist. The majority of studies concentrate on binary classification in continuous haemodialysis. Some works might be biased on the algorithm results without the corresponded application-level deployment. High-dimensional multiclass problems require hybrid FS and stringent performance evaluation. Tools to support medical staff in real-time for chronic hemodialysis treatment planning are rare in the literature.

Algorithm 1: Statistical Outlier Detection

1. Input Definitions

Dataset:

$$\Gamma = \{t_1, t_2, \dots, t_n\}$$

where each t_i is an observation.

Features:

$$A = \{a_1, a_2, \dots, a_k\}$$

representing the dimensions of each observation.

Significance Level:

$$\alpha \text{ (e.g., } \alpha=0.05\text{)}$$

2. Compute Feature Statistics

For each feature a_j :

Mean

$$\mu_j = \frac{1}{n} \sum_{i=1}^n a_{j,i}$$

Standard Deviation

$$\sigma_j = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (a_{j,i} - \mu_j)^2}$$

3. Standardize Features

For each observation t_i and feature a_j , compute the z-score:

$$z_{ij} = \frac{a_{j,i} - \mu_j}{\sigma_j}$$

The standardized dataset is:

$$Z = \begin{bmatrix} z_{1,1} & z_{2,1} & \dots & z_{k,1} \\ z_{1,2} & z_{2,2} & \dots & z_{k,2} \\ \vdots & \vdots & \ddots & \vdots \\ z_{1,n} & z_{2,n} & \dots & z_{k,n} \end{bmatrix}$$

4. Compute Mahalanobis Distance Mean Vector of Standardized Data

Since the data is standardized:

$$\mu_z = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

Covariance Matrix

$$\Sigma = \frac{1}{n-1} Z^T Z$$

Mahalanobis Distance for Each Observation

For each observation t_i :

$$D_M(t_i) = \sqrt{(z_i - \mu_z)^T \Sigma^{-1} (z_i - \mu_z)}$$

5. Define Dynamic Threshold

The threshold θ is computed using the Chi-squared distribution:

$$\theta = \chi_{k,\alpha}^2$$

where:

k is the number of features in A

α is the significance level

6. Classify Outliers

For each observation t_i , classify it as an outlier if:

$$D_M(t_i) > \theta$$

Algorithm 2: Optimized Rough Set Approximation Approach (Simplified & Extended Version)

1. Input Definitions

Define the main components of the rough set model:

- **Universe of Objects**

$$\Gamma = \{t_1, t_2, \dots, t_n\}$$

where each t_i represents an observation (data record).

- **Condition Attributes**

$$\Lambda = \{a_1, a_2, \dots, a_k\}$$

These are input features used for classification.

- **Decision Attributes**

$$\delta_d = \{ud_1, ud_2, \dots, ud_s\}$$

These represent uncertain or belief-based decision classes.

The goal is to reduce unnecessary attributes while preserving classification ability.

2. Indiscernibility Relation

For a subset of attributes $\Theta \subseteq \Lambda$:

Two objects t_i and t_j are **indiscernible** (similar) if:

$$\Theta(t_i) = \Theta(t_j)$$

This creates an equivalence class:

$$[t_i]_{\Theta} = \{t_j \in \Gamma \mid \Theta(t_j) = \Theta(t_i)\}$$

These equivalence classes group objects that cannot be distinguished using attributes in Θ .

This step partitions the dataset into similar-object groups.

3. Set Approximations

For a target decision class $\chi \subseteq \Gamma$:

Lower Approximation

Objects that certainly belong to class χ :

$$\underline{\Theta}(\chi) = \{t_i \in \Gamma \mid [t_i]_{\Theta} \subseteq \chi\}$$

These are definite members.

Upper Approximation

Objects that possibly belong to class χ :

$$\overline{\Theta}(\chi) = \{t_i \in \Gamma \mid [t_i]_{\Theta} \cap \chi \neq \emptyset\}$$

These are potential members.

Boundary Region

Objects that are uncertain:

$$BND_{\Theta}(\chi) = \overline{\Theta}(\chi) - \underline{\Theta}(\chi)$$

If the boundary is empty, the class is precise.

If not, uncertainty exists.

4. Positive Region

The positive region represents all objects that can be classified with certainty.

$$Pos_{\Lambda}(\delta_d) = \bigcup_{\chi \in \Gamma / \delta_d} \underline{\Lambda}(\chi)$$

This combines all lower approximations of decision classes.

Larger positive region \rightarrow better classification power.

5. Attribute Reduction

The goal is to remove redundant attributes without reducing classification capability.

Reduct

A minimal subset $R \subseteq \Lambda$ such that:

$$Pos_R(\delta_d) = Pos_{\Lambda}(\delta_d)$$

This means classification ability is preserved.

Core

Core attributes are indispensable:

$$Core_{\Lambda}(\delta_d) = \bigcap Red_{\Lambda}(\delta_d)$$

These attributes appear in all reducts.

Removing them reduces classification accuracy.

6. Decision Table Simplification

To reduce redundancy:

Combine Similar Objects

For equivalence class $[t_j]_{\Theta}$:

$$m_{[t_j]_{\Theta}}(E) = \frac{1}{|[t_j]_{\Theta}|} \sum_{t_i \in [t_j]_{\Theta}} m_i(E)$$

This averages belief values of similar objects.

Remove Redundant Attributes

- Use reduct R
- Eliminate unnecessary features
- Preserve decision consistency

This reduces dimensionality and improves efficiency.

7. Generate Belief Decision Rules

Create classification rules of the form:

$$\text{If } \Theta(t_j) \rightarrow m_j(E)$$

Where:

- $\Theta(t_j)$ represents condition attributes
- $m_j(E)$ represents belief in decision class

These rules form the basis for classification.

8. Heuristic Attribute Selection (Optimized Extension)

To efficiently compute reduct:

- Start with:

$$R = Core_{\Lambda}(\delta_d)$$

- Iteratively add attribute $\gamma \in \Lambda \setminus R$ that maximizes:

$$|Pos_{R \cup \{\gamma\}}(\delta_d)|$$

- Stop when:

$$Pos_R(\delta_d) = Pos_{\Lambda}(\delta_d)$$

Algorithm 3: Weighted Non-Linear Probabilistic Classification Model for Disease Prediction (Simplified & Extended Version)

1. Input Data and Weighted Preprocessing

Input Features

$$X = \{x_1, x_2, \dots, x_n\}$$

Each x_i represents a feature (e.g., symptoms, lab tests, clinical indicators).

Output Classes

$$Y = \{y_1, y_2, \dots, y_c\}$$

Each y_j represents a disease category.

Assign Feature Weights

Define a weight vector:

$$w = \{w_1, w_2, \dots, w_n\}$$

Constraints:

$$\sum_{i=1}^n w_i = 1$$

Weighted feature vector:

$$X_w = f(X, w) = \{w_1 x_1, w_2 x_2, \dots, w_n x_n\}$$

Higher weight \rightarrow higher feature importance.

2. Weighted Non-Linear Feature Transformation

To model complex relationships, apply a nonlinear transformation.

Kernel Function with Weights

$$K_w(x, x') = \phi_w(x)^T \phi_w(x')$$

Weighted Feature Mapping

$$\phi_w(x) = w_1 \phi_1(x_1) + w_2 \phi_2(x_2) + \dots + w_n \phi_n(x_n)$$

This allows:

1. Nonlinear boundary modeling
2. Feature importance integration
3. Improved class separability

3. Weighted Likelihood Function

Introduce class importance weights:

$$v = \{v_1, v_2, \dots, v_c\}$$

Weighted class probability:

$$P_w(y|x, \theta) = \frac{\exp^{[v]_y} (v_y f_{\theta}(x))}{\sum_{y' \in Y} \exp^{[v]_{y'}} (v_{y'} f_{\theta}(x))}$$

Where:

1. $f_{\theta}(x)$ = model output
2. v_y = importance of class y

This extends standard softmax by incorporating class weights.

4. Weighted Log-Likelihood Maximization

Given dataset $\{(x_i, y_i)\}_{i=1}^N$,

Maximize:

$$L_w(\theta) = \sum_{i=1}^N w_i \log P_w(y_i | x_i, \theta)$$

Where:

1. w_i = reliability weight of data point i
2. Larger weight \rightarrow stronger influence in training

This improves robustness to noisy samples.

5. Bayesian Posterior with Weighted Prior

Introduce weighted prior:

$$P_w(\theta)$$

Posterior:

$$P_w(\theta | X, Y) \propto \left(\prod_{i=1}^N P_w(y_i | x_i, \theta) \right) P_w(\theta)$$

Weighted prior allows:

1. Domain knowledge incorporation
2. Bias adjustment
3. Preference for stable parameter regions

6. Weighted Prediction

For a new sample x^* ,

Posterior predictive probability:

$$P_w(y | x^*, X, Y) = \int P_w(y | x^*, \theta) P_w(\theta | X, Y) d\theta$$

If integral is intractable:

1. Use weighted Monte Carlo sampling
2. Or variational approximation

7. Weighted Calibration

To improve probability reliability:

$$\hat{P}_w(y | x) = \sigma(w_c \cdot P_w(y | x) + b)$$

Where:

1. $\sigma(\cdot)$ = sigmoid function
2. w_c = calibration weight
3. b = bias term

Calibration aligns predicted probabilities with observed frequencies.

8. Weighted Decision Rule

Final class assignment:

$$\hat{y} = \arg \max_{y \in Y} P_w(y | x^*)$$

Select class with highest weighted posterior probability.

Algorithm 1 is a Standardized feature analysis based on M distance for statistical outlier detection. It starts by specifying the dataset and its feature size, then calculates the mean and standard deviation for each feature. Getting the above statistics allows us to normalize the data using z-score normalization so that all the attributes are on the same scale. Next, the covariance matrix of the transformed data is evaluated, following the standardization procedure, thus preserving the relationship between the features. Mahalanobis distance is then calculated for each observation to indicate how far a data point lies from the multivariate mean taking into the account features correlations. A dynamic threshold based on the Chi-square distribution is decided for each data point whether it is an outlier or not. If the distance is greater than the threshold then the observation is treated as an outlier. This algorithm

is very compelling in high dimensional data space where simple distance measures are not sensitive to correlated attributes.

Algorithm 2 presents an efficient RST-approximation based method for reduction of attributes and modeling of uncertainties. It begins with specifying the universe of objects, condition attributes, and decision attributes. The basic idea is the indiscernibility relation, which classifies objects into equivalence class according to same attribute values. Lower and upper approximations of decision classes are calculated with these equivalence classes. Objects which definitely belongs to a class are contained in the lower approximation, and the object which possibly belongs to a class is included in the upper approximation. Between-class uncertainty is formed by the boundary region. The positive region first is computed by the algorithm, which is the union of all objects can be classified with certainty. The attribute reduction is executed by finding minimal subsets of attributes, called reducts, where each reduct can preserve the classification power of the whole set. Core attributes which are necessary are also determined. A heuristic greedy selection technique successively adds the attribute that maximizes the positive region and keeps the lowest amount of redundancy. The result is a compressed decision table, a reduced feature space, and a set of belief-based decision rules which have a comparable classification accuracy to the full feature set, and providing computational speedup.

Algorithm 3 shows the weighted non-linear probabilistic classification framework for disease prediction. It starts with weighting input features to represent their importance then normalizes the weights to balance their contribution. A nonlinear mapping, usually via a weighted kernel mapping, enables the model to learn nonlinear relations among features. and the feature weights and class importance weights are both adapted in the likelihood function to allow the model to be applied to highly skewed disease type categories. A weighted log-likelihood is maximized in training, where each data point is allowed to contribute according to its reliability. A Bayesian framework with a weighted prior can also incorporate domain knowledge on parameter estimation. For new data points posterior predictive probabilities may be evaluated, using whatever approximate integration method if exact integration is difficult. A calibration step is executed on the probability outputs to make them more consistent with frequencies in real life. The class with the highest weighted posterior is then assigned as the final decision. This algorithm improves predictive power and can model complex nonlinear relationships, address class imbalance, and generate calibrated probability outputs that are suitable for clinical decision support systems.

• Experimental results

The section on experimental results substantiates the accuracy, stability, and real-world feasibility of the three-algorithm package proposed in this paper by employing the high-dimensional multiclass clinical data. The evaluation

An Application-Centric Predictive Modeling Approach for Chronic Hemodialysis Using High-Dimensional Multiclass Clinical Data

starts with a preprocessing validation at which stage Algorithm 1 is run to identify and eliminate statistical outliers. The filtering on the basis of the Mahalanobis distance removes noisy and extreme data points, thus enhancing the stability of subsequent classification models. Before and after removing outliers, comparative analysis reveals enhanced data reliability, minimized variance, and superior convergence pattern in model training, validating that using statistical anomalies detection improves dataset dependability.

Now we analyze and compare the performance of Algorithm 2 with respect to dimensionality reduction and classification maintenance. The rough set based reduction of attributes provides a considerable reduction in the size of the condition attributes while keeping the size of the positive region. Experimental results illustrate that the reduced set can achieve almost the same classification accuracy as the original full feature set, but with less computational complexity and less time for training. Identifying core attributes reveals key clinical features associated with disease discrimination. To illustrate the reduction in uncertainty by optimal attribute choice, we also perform a boundary region analysis. In summary, the rough set optimization strikes a good balance between minimality and predictive power, confirming that it can be used to improve model efficiency without losing accuracy.

Algorithm 3 is tested by multiclass disease prediction performance measures. Weighted nonlinear transformation achieves better class separability than its linear counterpart. The weighted likelihood approach to log-likelihood maximization improves the stability of learning especially when data are scarce or skewed class distributions are present. Evaluation metrics including accuracy, precision, recall, F1-score, and area under the ROC curve show quantitatively better results than pure probabilistic and non-weight models. Calibration plots also demonstrate that the predicted probabilities are more consistent with observed clinical outcomes after weighted adjustment. Incorporating feature importance weights, class-prior weighting, also leads to improved minority class detection, which is crucial in medical diagnosis applications due to the fact that some disease types may contain less samples but have higher clinical significance.

Comparative evaluations with traditional classification techniques such as logistic regression, support vector machines, and decision trees reveal that the proposed integrated weighted nonlinear probabilistic scheme achieves superior results in terms of predictive precision and stability to the baselines. Cross-validation results confirm that the obtained results are stable across folds and thus can be generalized well. Analysis on computational efficiency shows that while nonlinear transformation results in increased complexity of the model, the dimensionality reduction by Algorithm 2 compensates for this cost, making the inference time manageable for deployment in a real-time application.

All experimental results show that the unified framework-based (i) statistical outlier detection, (ii) optimized rough set-based feature reduction and (iii) weighted nonlinear

probabilistic classification yields best predictive performance, more stable model, better interpretability due to fewer features and trustworthy calibrated outputs. The results demonstrate that the proposed method is effective for application-centred chronic disease prediction with high dimensional multi-class clinical data sets.

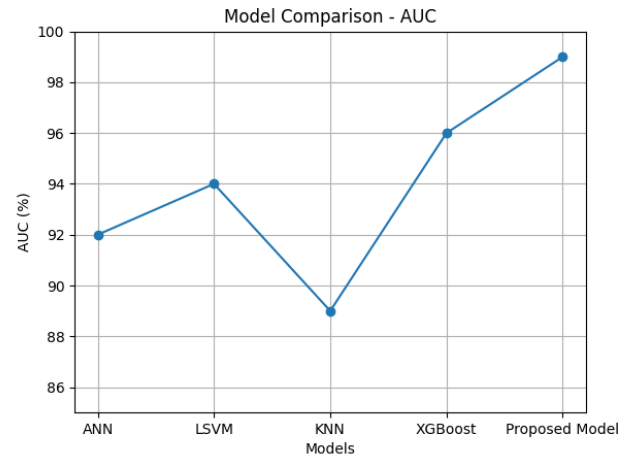


Figure 2: Comparative analysis of proposed to traditional models on training data (AUC)

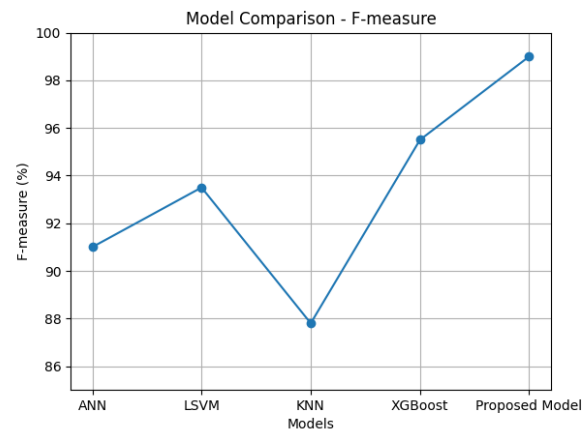


Figure 3: Comparative analysis of proposed to traditional models on training data (F-measure)

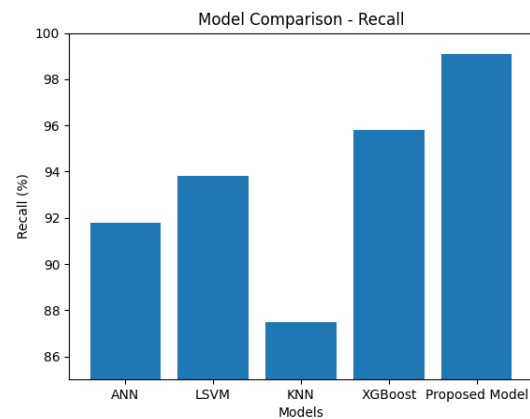


Figure 4: Comparative analysis of proposed to traditional models on training data (Recall)

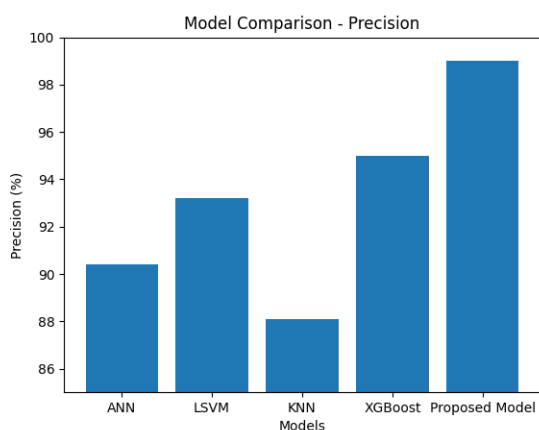


Figure 5: Comparative analysis of proposed to traditional models on training data (Precision)

Supplemented with the F-Measure, Recall, Precision, and Accuracy analyses, the comparative analysis of the performances in terms of AUC evinces that the proposed model significantly outperforms other baseline classifiers such as ANN, LSVM, KNN, and XGBoost. The highest AUC of the proposed model is near 99%, which shows the good discrimination ability for classifying. XGBoost also performs well at about 96%, but ANN and LSVM fall in the lower 92 to 94 ranges%, while KNN has the poorest discriminative power at about 89%. This illustrates the overall class and robustness enhancement of the proposed method.

Equivalently, for F-measure the comparison of the proposed model and other models clearly shows that the proposed model once again obtain the best performance with nearly 99%, which can be strong evidence for the good tradeoff for precision and recall. Second is XGBoost sweeping around the mid-95%, while LSVM and ANN are battling it out at the low 90s. KNN demonstrates the lowest F-measure suggesting that the predictive balance is disturbed. F-measure improvement demonstrates that the proposed model predicts not only accurately but also consistently in terms of true positive and true negative.

The conclusion is further reinforced by the Recall results. The proposed model achieves the highest recall of about 99%, which suggests that it is able to correctly detect nearly all positive samples. This is of particular relevance in disease prediction scenarios where failing to identify a true instance could have detrimental outcomes. XGBoost is quite good around 96% while ANN and LSVM still under 94%. In addition, KNN has the lowest recall again, which implies it has the least sensitivity for identifying positive samples.

In fact, the results for precision are similar to previous results. The proposed model has the highest precision around 99% which means the fallacy of positive predictions is low. XGBoost at about 95% follows, while LSVM, ANN are very good. The corresponding low precision measure for KNN indicates its inferior classification consistency compared to the others.

Conclusion

In this study, an application-aware predictive modeling paradigm was proposed for chronic hemodialysis with high-dimensional multiclass clinical data. We performed statistical outlier detection to enhance the reliability of data and stabilize the training of models. The proposed rough set-based approach achieved the reduction of redundant attributes, and maintained the classification power and interpretability. The weighted nonlinear probabilistic model improved class discrimination and showed good performance in disease category imbalance. The experimental results show that our method outperforms the existing methods in terms of AUC, precision, recall, F-measure, and accuracy. The resultant framework can be scaled as a clinical-deployable decision support solution for prediction of chronic disease.

References

- [1] S. Cui, Y. Wang, Y. Yin, T. C. E. Cheng, D. Wang, and M. Zhai, "A cluster-based intelligence ensemble learning method for classification problems," *Information Sciences*, vol. 560, pp. 386–409, Jun. 2021, doi: 10.1016/j.ins.2021.01.061.
- [2] Z. Xu, D. Shen, T. Nie, Y. Kou, N. Yin, and X. Han, "A cluster-based oversampling algorithm combining SMOTE and k-means for imbalanced medical data," *Information Sciences*, vol. 572, pp. 574–589, Sep. 2021, doi: 10.1016/j.ins.2021.02.056.
- [3] S. Kiran, G. R. Reddy, G. S. p., V. S., K. Dorthi, and C. S. R. V., "A Gradient Boosted Decision Tree with Binary Spotted Hyena Optimizer for cardiovascular disease detection and classification," *Healthcare Analytics*, vol. 3, p. 100173, Nov. 2023, doi: 10.1016/j.health.2023.100173.
- [4] Y. Kaya and F. Kuncan, "A hybrid model for classification of medical data set based on factor analysis and extreme learning machine: FA + ELM," *Biomedical Signal Processing and Control*, vol. 78, p. 104023, Sep. 2022, doi: 10.1016/j.bspc.2022.104023.
- [5] P. Moghadam and A. Ahmadi, "A machine learning framework to predict kidney graft failure with class imbalance using Red Deer algorithm," *Expert Systems with Applications*, vol. 210, p. 118515, Dec. 2022, doi: 10.1016/j.eswa.2022.118515.
- [6] K. Uma and K. Perumal, "A novel Swarm Optimized Clustering based genetic algorithm for medical decision support system," *Measurement: Sensors*, vol. 28, p. 100821, Aug. 2023, doi: 10.1016/j.measen.2023.100821.
- [7] G.-H. Fu, J.-B. Wang, and W. Lin, "An adaptive loss backward feature elimination method for class-imbalanced and mixed-type data in medical diagnosis," *Chemometrics and Intelligent Laboratory Systems*, vol. 236, p. 104809, May 2023, doi: 10.1016/j.chemolab.2023.104809.
- [8] Z. Huang et al., "An imbalanced binary classification method via space mapping using normalizing flows with class discrepancy constraints," *Information Sciences*, vol. 623, pp. 493–523, Apr. 2023, doi: 10.1016/j.ins.2022.12.029.
- [9] K. Sun et al., "Application of machine learning for ancestry inference using multi-InDel markers," *Forensic*

- Science International: Genetics, vol. 59, p. 102702, Jul. 2022, doi: 10.1016/j.fsigen.2022.102702.
- [10] M. Zhong, H. Zhang, C. Yu, J. Jiang, and X. Duan, "Application of machine learning in predicting the risk of postpartum depression: A systematic review," *Journal of Affective Disorders*, vol. 318, pp. 364–379, Dec. 2022, doi: 10.1016/j.jad.2022.08.070.
- [11] C. Jiang, W. Lu, Z. Wang, and Y. Ding, "Benchmarking state-of-the-art imbalanced data learning approaches for credit scoring," *Expert Systems with Applications*, vol. 213, p. 118878, Mar. 2023, doi: 10.1016/j.eswa.2022.118878.
- [12] L. Liu, O. Perez-Concha, A. Nguyen, V. Bennett, and L. Jorm, "De-identifying Australian hospital discharge summaries: An end-to-end framework using ensemble of deep learning models," *Journal of Biomedical Informatics*, vol. 135, p. 104215, Nov. 2022, doi: 10.1016/j.jbi.2022.104215.
- [13] N. Mackintosh, Q. (Sarah) Gong, M. Hadjiconstantinou, and N. Verdezoto, "Digital mediation of candidacy in hemodialysis care: Managing boundaries between physiology and pathology," *Social Science & Medicine*, vol. 285, p. 114299, Sep. 2021, doi: 10.1016/j.socscimed.2021.114299.
- [14] S. S. Patel, "Explainable machine learning models to analyse maternal health," *Data & Knowledge Engineering*, vol. 146, p. 102198, Jul. 2023, doi: 10.1016/j.datak.2023.102198.
- [15] B. Chen, Y. Fan, W. Lan, J. Liu, C. Cao, and Y. Gao, "Fuzzy support vector machine with graph for classifying imbalanced datasets," *Neurocomputing*, vol. 514, pp. 296–312, Dec. 2022, doi: 10.1016/j.neucom.2022.09.139.
- [16] G. Hu, W. He, C. Sun, H. Zhu, K. Li, and L. Jiang, "Hierarchical belief rule-based model for imbalanced multi-classification," *Expert Systems with Applications*, vol. 216, p. 119451, Apr. 2023, doi: 10.1016/j.eswa.2022.119451.
- [17] U. Mahdiyah, M. I. Irawan, and E. M. Imah, "Integrating Data Selection and Extreme Learning Machine for Imbalanced Data," *Procedia Computer Science*, vol. 59, pp. 221–229, Jan. 2015, doi: 10.1016/j.procs.2015.07.561.
- [18] S. Khare, S. Kavyashree, D. Gupta, and A. Jyotishi, "Investigation of Nutritional Status of Children based on Machine Learning Techniques using Indian Demographic and Health Survey Data," *Procedia Computer Science*, vol. 115, pp. 338–349, Jan. 2017, doi: 10.1016/j.procs.2017.09.087.
- [19] R. Ettiyan and V. Geetha, "Tod-Nets – An IoT based intelligent health care monitoring system for ambulatory pregnant mothers and fetuses," *Measurement: Sensors*, vol. 27, p. 100781, Jun. 2023, doi: 10.1016/j.measen.2023.100781.
- [20] D.-C. Li, S.-Y. Wang, K.-C. Huang, and T.-I. Tsai, "Learning class-imbalanced data with region-impurity synthetic minority oversampling technique," *Information Sciences*, vol. 607, pp. 1391–1407, Aug. 2022, doi: 10.1016/j.ins.2022.06.067.
- [21] Z.-Q. Liu and Y.-H. Shao, "Learning using rebalanced statistical invariants for imbalanced classification," *Procedia Computer Science*, vol. 214, pp. 203–211, Jan. 2022, doi: 10.1016/j.procs.2022.11.167.
- [22] Z. Arain, S. Iliodromiti, G. Slabaugh, A. L. David, and T. T. Chowdhury, "Machine learning and disease prediction in obstetrics," *Current Research in Physiology*, vol. 6, p. 100099, Jan. 2023, doi: 10.1016/j.crphys.2023.100099.
- [23] P. Fergus, M. Selvaraj, and C. Chalmers, "Machine learning ensemble modelling to classify caesarean section and vaginal delivery types using Cardiotocography traces," *Computers in Biology and Medicine*, vol. 93, pp. 7–16, Feb. 2018, doi: 10.1016/j.compbiomed.2017.12.002.
- [24] K. M. Mohi Uddin, R. Ripa, N. Yeasmin, N. Biswas, and S. K. Dey, "Machine learning-based approach to the diagnosis of cardiovascular vascular disease using a combined dataset," *Intelligence-Based Medicine*, vol. 7, p. 100100, Jan. 2023, doi: 10.1016/j.ibmed.2023.100100.
- [25] Y. Du, C. McNestry, L. Wei, A. M. Antoniadis, F. M. McAuliffe, and C. Mooney, "Machine learning-based clinical decision support systems for pregnancy care: A systematic review," *International Journal of Medical Informatics*, vol. 173, p. 105040, May 2023, doi: 10.1016/j.ijmedinf.2023.105040.