

Ashish Kumar<sup>1</sup>,  
 Ram Kinkar Pandey<sup>2</sup>,  
 Prabhat Kumar  
 Srivastava<sup>3</sup>

## Class-Imbalance Aware Machine Learning for CKD Detection and Risk Assessment



**Abstract:** Chronic kidney disease (CKD) is a growing global health concern, contributing significantly to morbidity and mortality. Early detection and management of CKD and its complications (such as hyperkalemia) are vital to improving patient outcomes. Recently, machine learning techniques have shown promise in improving CKD diagnosis and prognosis. In this study, we develop and evaluate a machine learning approach for predicting CKD presence using patient data, and we identify key risk factors from clinical, lifestyle, and laboratory features. We utilize an ensemble Random Forest (RF) classifier on a dataset of 1,659 individuals (91.9% CKD patients, 8.1% non-CKD) to distinguish CKD from healthy status. The model achieved a high overall accuracy (~92%), correctly identifying all CKD cases (100% sensitivity) but with limited specificity due to class imbalance. We address this imbalance via techniques such as class weighting, which modestly improved detection of non-CKD cases. The most influential predictors of CKD in our data were Serum Creatinine, Proteinuria (protein in urine), and Glomerular Filtration Rate (GFR), aligning with medical knowledge of kidney function. Certain clinical symptoms (e.g. itching and muscle cramps) also emerged as important indicators of CKD. We further discuss our findings in the context of recent literature – including a related study that attained 99.8% CKD prediction accuracy using RF with feature selection[4], and another that employed an XGBoost model to predict hyperkalemia (a dangerous CKD complication) with an AUC of 0.867, outperforming clinicians[5]. Our results underscore the potential of machine learning models to support early CKD diagnosis and complication risk forecasting, while highlighting the challenges posed by imbalanced datasets and the need for careful feature consideration. We conclude that ensemble learning methods, when combined with domain-specific feature insights, can provide highly accurate and clinically useful decision support for CKD management, though further work is needed to improve model generalizability and the detection of rare outcomes. The best-performing unbalanced Logistic Regression achieved an AUC of 0.804 and an accuracy of 92.5 %, while a probability-weighted voting ensemble raised accuracy to 93.7 %. Adjusting the decision threshold on the logistic model further boosted accuracy to 94.0 % with perfect recall for CKD cases.

**Keywords:** Chronic kidney disease, Serum Creatinine, Proteinuria (protein in urine), and Glomerular Filtration Rate (GFR)

### I. Introduction

Chronic kidney disease (CKD) is a non-communicable condition characterized by gradual loss of kidney function, and it has become a significant public health issue worldwide[1]. The global impact of CKD has risen dramatically in recent decades – between 1990 and 2013, the annual loss of life years due to CKD increased by 90%, making it the 13th leading cause of death globally[2]. An estimated 850 million people around the world have kidney diseases of varying severity[2]. According to the 2019 World Kidney Day report, at least 2.4 million people die each year from kidney-related diseases, and CKD is currently the 6th fastest-growing cause of death worldwide[3]. These figures underline the urgent need for improved strategies in early detection and management of CKD.

CKD often progresses silently, with early stages (Stages 1–2) presenting minimal symptoms[4]. As the disease advances to later stages (Stages 4–5), patients may develop severe complications and require renal replacement therapy (dialysis or transplantation) to survive[5]. Late-stage CKD management is challenging and costly, especially in low-resource settings where access to dialysis and transplantation is limited[6]. Therefore, timely identification of CKD and intervention in the **early stages** are critical to reduce the risk of adverse outcomes and to ease the economic burden on healthcare systems[2]. Early diagnosis allows clinicians to treat contributing

<sup>1</sup>Research Scholar, Department of Computer Science and Engineering, Arni University, Indora, Kathgarh Kangra(H.P),

<sup>2</sup> Department of Computer Science and Engineering, Arni University, Indora, Kathgarh, Kangra(H.P),

<sup>3</sup>Department of Computer Science and Engineering, IMS Engineering College, Ghaziabad, India  
 Emails: 1kumar7ashish786@gmail.com,2Dr.ramkpandey@gmail.com,3sri\_prab@rediffmail.com

\*Corresponding Author: kumar7ashish786@gmail.com

factors (such as hypertension or diabetes), recommend lifestyle modifications, and delay progression to kidney failure[7].

Beyond slowing CKD progression, early detection also helps prevent complications. **Hyperkalemia** (elevated blood potassium) is one such life-threatening complication frequently seen in CKD due to reduced renal excretion of potassium[8]. Hyperkalemia is associated with cardiac arrhythmias and increased mortality in CKD patients[9]. Its prevalence is about 9% among CKD patients overall and up to one-third in non-dialysis advanced CKD under nephrology care[8]. Major risk factors for hyperkalemia include comorbid congestive heart failure, diabetes, older age, high dietary potassium intake, and use of medications that inhibit the renin-angiotensin-aldosterone system (e.g. ACE inhibitors or ARBs)[9]. Hyperkalemia has been shown to worsen CKD outcomes, leading to more hospitalizations and higher healthcare costs. Identifying patients at risk for hyperkalemia in advance would enable proactive management (dietary counseling, medication adjustment, use of potassium binders) to prevent serious events.

Given the importance of early CKD detection and complication prevention, there has been growing interest in applying machine learning (ML) techniques to assist in clinical decision-making for kidney disease[10]. ML methods can potentially uncover complex patterns in clinical data that might be missed by traditional approaches, thereby improving diagnostic accuracy and risk prediction. In fact, numerous studies have explored ML models for CKD classification and prognosis. For instance, research by Charleonnan et al. compared algorithms such as k-Nearest Neighbors, Support Vector Machines (SVM), logistic regression, and decision trees on a benchmark CKD dataset; they reported SVM achieved an accuracy of 98.3% in classifying CKD vs non-CKD cases[11]. More recently, authors applied Random Forest (RF), SVM, and decision tree models to a large CKD dataset (1,718 patient records from St. Paulo's Hospital, Ethiopia) for both binary CKD detection and multi-class CKD stage prediction[12] [13]. Their work demonstrated that ensemble methods can be highly effective: the RF model with recursive feature elimination attained **99.8% accuracy** for binary CKD classification, outperforming the other classifiers[14]. This near-perfect accuracy was achieved using a subset of only 8 features, indicating that a carefully chosen set of clinical indicators can almost conclusively distinguish CKD patients from healthy individuals in their dataset[15]. Also showed that CKD **stage** prediction (a five-class classification problem) is more challenging – their best multi-class model reached about 79% accuracy after feature selection, reflecting the increased complexity in differentiating between CKD stages as opposed to the binary (CKD vs not CKD) scenario.

In parallel, machine learning has been employed to predict specific outcomes for CKD patients, such as hyperkalemia risk. Authors developed an extreme gradient boosting model (XGBoost) to forecast hyperkalemia in advanced CKD patients using electronic health record data[16]. Notably, they conducted a human-machine competition where the ML model's performance was compared against two nephrologists' predictions. The XGBoost model achieved an area under the ROC curve (AUC) of 0.867 and an accuracy of 93.3%, significantly outperforming the clinicians in anticipating hyperkalemia events[17]. The model's positive predictive value (PPV) was 0.70 at a sensitivity of about 87%, indicating a substantial improvement over standard care. Moreover, identified four top predictors for hyperkalemia: hemoglobin level, the patient's last serum potassium level, use of angiotensin receptor blockers (ARB medications), and use of potassium-binding resin (calcium polystyrene sulfonate)[18]. These predictors make clinical sense – for example, lower hemoglobin often correlates with advanced CKD and anemia of chronic disease, and a higher prior potassium or use of ARBs can predispose to future hyperkalemia. Their study underscores how ML models can integrate diverse clinical features (labs, medications, history) to generate risk predictions that even specialists might not match. Other studies referenced by Chang have also shown success in hyperkalemia prediction: a U.S. claims-based analysis used logistic regression to predict hyperkalemia in CKD patients[19], and deep learning models analyzing electrocardiogram (ECG) signals have accurately detected hyperkalemia in CKD and emergency department settings. These advances highlight a broader trend of leveraging ML in nephrology for both disease detection and complication prediction.

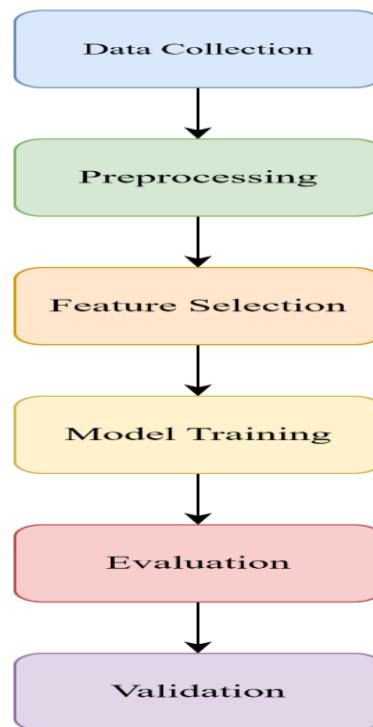
## II. Related Work

Early applications of machine learning to CKD predominantly utilized the publicly available UCI CKD dataset (~400 instances) or similar small datasets, often demonstrating high accuracy for binary classification. As mentioned, Charleonnan et al. achieved over 98% accuracy with SVM on the UCI CKD data [11], and other methods like k-NN and neural networks have likewise reported accuracies above 95% on this dataset under various experimental settings. However, a limitation noted in the literature was the small size and limited features of these datasets[20], which raised concerns about overfitting and generalizability. In response, recent research has shifted towards using larger, more diverse datasets and focusing not only on CKD detection but also on stratifying disease severity. Authors [13] provided one of the most comprehensive studies in this regard. They compiled a dataset of 1,718 records from patients admitted to a renal unit in Ethiopia, including 19 clinical

features such as vital signs (blood pressure), laboratory measurements (e.g., blood urea nitrogen, serum creatinine, electrolytes, blood counts), and comorbidities (hypertension, diabetes, anemia). Importantly, their dataset encompassed patients across all five CKD stages as defined by Kidney Disease Improving Global Outcomes (KDIGO) guidelines, enabling both binary classification (CKD vs non-CKD) and multi-class classification (predicting CKD stage 1 through 5). They applied two feature selection techniques – a univariate filter based on Analysis of Variance (ANOVA) and a wrapper method using Recursive Feature Elimination with Cross-Validation (RFECV) – to identify optimal feature subsets for each modeling task[21]. For the **binary classification** (distinguishing CKD patients from those without CKD), their baseline RF model already achieved 99.7% accuracy using all features [22], indicating that even without feature reduction, the data was highly separable. After applying RFECV, the RF model using only 8 selected features slightly improved accuracy to **99.8%** (with sensitivity 99.7% and specificity 99.9%)[22]. This suggests that a small set of features – likely including kidney function indicators and perhaps anemia or electrolyte measures – carried almost all the information needed to identify CKD in their cohort. In comparison, SVM and decision tree models performed a bit lower (in the 95–98% accuracy range) on the binary task, though SVM with hyperparameter tuning also reached ~99.8% in their experiments[21]. For the **five-class stage prediction**, the task was naturally more difficult. Using all features, the RF model obtained around 76% accuracy, which improved to ~79% after RFECV feature selection (9 features). SVM and DT were slightly lower (mid-70s% at best) on multi-class accuracy. Notably, the class imbalance in stage distribution (fewer patients in extreme stages) and overlapping characteristics between adjacent stages likely limited performance. Debal and Sitote’s study demonstrates that with sufficient data and feature engineering, traditional ML models (especially ensemble trees) can deliver excellent performance for CKD prediction. It also highlights the value of feature selection in reducing model complexity while maintaining accuracy – their RFECV consistently chose a subset of lab tests as the most predictive features, reaffirming known clinical markers of CKD such as low GFR, high creatinine, abnormal blood urea, etc., even though the exact features chosen in each run could vary. Another interesting aspect was their use of **10-fold cross-validation** for evaluation [23], which provides a robust estimate of model generalization given the data. Authors [24] approached the problem from a different angle by targeting a complication (hyperkalemia) in an advanced CKD population. Their work is a prime example of using ML for predictive analytics in nephrology care beyond just diagnosing CKD. They collected a decade’s worth of data from a pre-ESRD program in Taiwan, focusing on patients with Stage 3b–5 CKD (i.e., relatively advanced CKD). The outcome of interest was whether a patient would experience hyperkalemia ( $K^{+} > 5.5$  mEq/L) at their next clinic visit, given the data from the current visit. Feature set was extensive, including demographics, a broad panel of laboratory results, comorbid diagnoses, and medication use records. They used **XGBoost** for its ability to handle high-dimensional data and capture non-linear interactions. Through ten-fold cross-validated grid search, they optimized the model’s hyperparameters and then evaluated it on a hold-out test set comprising 25% of the patients. The XGBoost model’s performance was impressive: **AUC = 0.867**, accuracy = 93.3%, with a positive predictive value of 70.0% at the operating threshold[17]. In a head-to-head comparison, this ML model substantially outperformed two nephrologists who were asked to manually predict hyperkalemia using the same test cases. This result illustrates the potential of ML to augment clinical decision-making, especially in scenarios where human prediction may rely on gestalt or limited cues. Authors [18] provide insight into why the model excelled by discussing its top features: prior lab results (especially the previous potassium level), certain medications (e.g., ARBs which can raise potassium, and usage of potassium binders which often indicates recent hyperkalemia management), and hemoglobin. Hemoglobin may seem an indirect factor, but anemia in CKD can reflect poor erythropoietin production and advanced disease, which correlates with disturbances in other parameters like electrolytes. These features align well with medical understanding – for example, a patient who had a high  $K^{+}$  reading at the last visit or who is on an ARB is known to be at elevated risk for hyperkalemia. The model effectively integrates such information. Furthermore, the authors note prior studies that attempted hyperkalemia prediction: one using logistic regression on insurance claims data, and several using deep learning on ECG signals, all of which reinforces that hyperkalemia can be anticipated with the right data inputs. Chang’s work is a salient reminder that ML applications in CKD go beyond diagnosis – they can help foresee complications and potentially guide preventative interventions.

In light of these studies, our research is positioned to contribute additional perspective by examining a real-world CKD dataset with a rich feature set, including some less commonly studied factors like lifestyle and quality of life scores. We leverage an ensemble tree model (Random Forest) for CKD prediction, given its success in prior studies[22], and we place special emphasis on interpreting the model (extracting feature importances) and addressing class imbalance issues.

## Methodology Used



**Figure. 1: Workflow of the Proposed CKD Prediction Methodology**

The figure 1 display the workflow of the proposed CKD prediction methodology, in which first data is collected, which is then preprocessing, feature selection, after then model training, evaluation and validation.

### Dataset and Features

For this study, we utilized a retrospective CKD dataset containing **1,659 patient records**, each with a diagnosis label indicating the presence (CKD=1) or absence (CKD=0) of chronic kidney disease. The data appears to have been collected from a nephrology care setting or a high-risk cohort, given the high prevalence of CKD cases (91.9%) relative to non-CKD controls (8.1%). Specifically, out of 1,659 individuals, 1,524 were diagnosed with CKD and only 135 were non-CKD, reflecting a roughly 11:1 class ratio favoring CKD. Such an imbalance suggests the dataset may originate from a specialized CKD program or clinic (where healthy individuals are relatively few).

Each record in the dataset includes a broad array of features spanning **demographics, clinical symptoms, vital signs, laboratory measurements, lifestyle factors, and medications**. Below provides an overview of the key feature categories:

- **Demographics & Social:** Age (in years), Gender (0 = Female, 1 = Male), Ethnicity (categorical, 4 groups), and Socioeconomic Status. These features capture basic patient characteristics that could influence CKD risk (e.g., age is a known risk factor, and socio-economic factors can affect access to healthcare and nutrition).
- **Vital Signs & Anthropometrics:** Systolic Blood Pressure (SBP) and Diastolic Blood Pressure (DBP) in mmHg, and Body Mass Index (BMI). Hypertension and obesity are well-known contributors to CKD progression and onset; hence blood pressure readings and BMI are relevant predictors.
- **Laboratory Measurements:**
  - Serum Creatinine (mg/dL) – a primary indicator of kidney function (higher creatinine reflects poorer kidney filtration).
  - Glomerular Filtration Rate (GFR, mL/min/1.73m<sup>2</sup>) – either measured or estimated, directly quantifying kidney filtration capacity (lower GFR indicates CKD).
  - Blood Urea Nitrogen (BUN, mg/dL) – another waste product like creatinine, typically elevated in CKD.

- HbA1c (%) – glycated hemoglobin, a measure of long-term blood sugar control (included because diabetes is a major CKD risk factor).
- Possibly other labs (not explicitly listed in summary, but based on CKD studies and Debal & Sitote’s features, the dataset could include electrolytes like sodium, potassium, blood counts, etc. However, our analysis specifically highlighted creatinine, GFR, BUN, and HbA1c as being present).
- **Clinical Symptoms:** Categorical indicators (likely binary 0/1) for symptoms commonly associated with CKD. Our dataset specifically had features for symptoms such as Itching, Muscle Cramps, Nausea/Vomiting, and Edema. These are uremic or CKD-related symptoms: e.g., uremic pruritus (itching) and muscle cramps often occur in advanced CKD due to toxin accumulation and electrolyte imbalances.
- **Lifestyle Factors:** Self-reported or assessed measures including Smoking status, Alcohol Consumption, Physical Activity level, Diet Quality, Sleep Quality, and a Quality of Life (QoL) Score. Each of these was recorded in a quantitative or ordinal manner (for instance, Physical Activity and Diet Quality on numeric scales, QoL on a 0–100 scale). Such features are less common in CKD prediction studies but are valuable for understanding patient health holistically. Poor diet and low activity can exacerbate CKD, while QoL may be affected by CKD and comorbid conditions.
- **Family History:** Binary indicators for family history of kidney disease, hypertension, and diabetes. These genetic or environmental predisposition factors can improve risk stratification (e.g., a family history of CKD or diabetes increases one’s own risk).

Table 1: Summary of Dataset Features

Category	Example Features
Demographics	Age (years), Sex (Male/Female), Ethnicity (4 categories), Socioeconomic Status
Vitals/Body	Systolic BP (mmHg), Diastolic BP (mmHg), BMI (kg/m <sup>2</sup> )
Laboratory	Serum Creatinine (mg/dL), GFR (mL/min/1.73m <sup>2</sup> ), BUN (mg/dL), HbA1c (%), [others]
Symptoms	Itching (0/1), Muscle Cramps (0/1), Nausea/Vomiting (0/1), Edema (0/1), Fatigue Levels
Lifestyle	Smoking (0/1), Alcohol Consumption (0/1), Physical Activity (scale), Diet Quality (scale), Sleep Quality, Quality of Life Score (0–100)
Family History	Kidney Disease (0/1), Hypertension (0/1), Diabetes (0/1)
Medications	ACE Inhibitor/ARB (0/1), Diuretics (0/1), NSAIDs use (0/1), Statins (0/1), Antidiabetic meds (0/1)

### Exploratory Data Analysis

Before building prediction models, an exploratory data analysis (EDA) was conducted to understand the dataset’s characteristics and to verify that the features align with known CKD patterns. Key findings from the EDA are summarized below:

- **Disease Prevalence:** As noted, the dataset is heavily skewed towards CKD cases. **91.9%** of the records correspond to patients with CKD, while only **8.1%** are healthy controls. This confirms a significant class imbalance that we would need to address during modeling. Such imbalance is not unusual in specialized datasets and is a focal point for our modeling strategy (see Model Training section).
- **Demographics:** The gender distribution was roughly balanced (approximately 51.5% male and 48.5% female), indicating that both sexes are well-represented among the patients. (CKD is known to affect both men and women, though certain etiologies differ by sex.) Ethnicity spanned four categories, with one category (labelled “0”) being the most common in this cohort. The socioeconomic status variable also varied, but details of its scale or distribution were not the primary focus of the analysis.

- Clinical Biomarkers:** As expected, CKD patients showed markedly different laboratory values compared to non-CKD individuals. The **mean GFR** among CKD patients was 65.26 ( $\pm$  std dev) vs **84.59** in non-CKD – a significantly lower GFR for CKD patients, consistent with impaired kidney function. CKD patients also had higher **Serum Creatinine** on average (2.83 mg/dL vs 1.86 mg/dL in non-CKD), since creatinine accumulates when filtration is reduced. These differences in GFR and creatinine are highly statistically significant ( $p < 0.001$ ) and reaffirm that our dataset’s diagnoses align with standard clinical definitions (CKD is often defined by  $GFR < 60$  or creatinine above normal range). Other lab differences included: **Blood Urea Nitrogen (BUN)** was elevated in CKD (mean  $\sim$ 27.93 vs 23.57), reflecting reduced waste excretion; and **Systolic Blood Pressure** was higher in CKD patients (135.0 mmHg vs 127.2 mmHg on average), consistent with hypertension being both a cause and consequence of CKD. Even the **BMI** showed a slight increase in CKD patients (27.74 vs 26.27), though this difference is modest – it could suggest fluid retention or simply lifestyle factors. Collectively, these patterns match well-known CKD phenotype characteristics (lower kidney function metrics and higher waste product levels in blood). We also examined differences in other labs like HbA1c: interestingly, average HbA1c was only slightly different (not listed in the top five differences), implying that while diabetes is a risk factor, glycemic control measures might not starkly separate CKD vs non-CKD groups in this cohort (possibly many non-CKD individuals were diabetic too, or CKD patients were on treatment to control HbA1c).
- Lifestyle Factors:** The data revealed some noteworthy lifestyle and quality of life trends. CKD patients reported a **lower diet quality** score on average (4.99 vs 5.46 in non-CKD, on whatever scale was used), and slightly less **physical activity** (5.01 vs 5.23). These differences suggest that CKD patients may have more dietary restrictions (or poorer diet due to illness) and limitations in exercise capacity. Surprisingly, the average **Quality of Life (QoL) score** was higher in CKD patients (50.07 vs 45.93). This counter-intuitive finding (“paradoxically” higher QoL in CKD) could be due to several factors – perhaps CKD patients in a managed program receive more support or have adjusted expectations, or maybe many of the “non-CKD” individuals in this dataset had other health issues affecting their QoL. It’s also plausible that QoL was assessed at a point when CKD patients were receiving treatment that improved certain aspects of life (for example, relief from symptoms due to medications). Further investigation would be needed to interpret this properly, but it highlights the complexity of patient-reported outcomes in chronic illness.
- Medication Usage:** We observed differences in the usage rates of certain medications between CKD and non-CKD groups. **Diuretics** were used by about 32.2% of CKD patients versus 27.4% of non-CKD. The higher use in CKD is expected since diuretics help manage blood pressure and edema in CKD care. **NSAIDs** (non-prescription analgesics) had high usage in both groups ( $\sim$ large portion of patients in each), but it was “slightly higher” in CKD patients. The prevalence of NSAID use is concerning given NSAIDs’ nephrotoxic potential; this data point might reflect that many patients (even those with CKD) continue to use NSAIDs for pain, underscoring an area for patient education. Usage of ACE inhibitors, ARBs, and other medications were also presumably captured, but the summary specifically pointed out diuretics and NSAIDs as notable differences. ACE inhibitor/ARB usage, for example, is often very common in CKD for renoprotection, so it might have been high in both groups if many non-CKD had hypertension too. The medication comparisons overall provide insight that CKD patients were more likely on therapies addressing complications (like diuretics for fluid control).

These exploratory findings confirmed that the dataset is clinically sound (CKD patients exhibit the expected derangements in labs and are on appropriate medications) and also highlighted the **high CKD prevalence** as a challenge for modeling. The fact that 91.9% of records are CKD means a naive classifier that predicts “CKD” for everyone would be right 91.9% of the time – this baseline sets a high bar for accuracy but would obviously be a trivial and not useful model (since it would have 0% specificity). Indeed, our subsequent model training had to contend with this issue. Nevertheless, EDA indicated there are clear signal differences in many features between classes, which is encouraging for the prospect of machine learning to detect CKD.

## Model Development

### Training Procedure

We formulated CKD prediction as a **binary classification** problem: given a patient’s features, predict **CKD (positive class)** or **No CKD (negative class)**. We selected the **Random Forest (RF)** algorithm as our primary model, due to its strong performance in similar medical classification tasks[4] and its advantages of handling feature heterogeneity and providing feature importance metrics. The RF was implemented using scikit-learn’s RandomForestClassifier. We set an ensemble size of 100 trees ( $n\_estimators = 100$ ) and a maximum tree depth of 10 ( $max\_depth = 10$ ) to prevent individual trees from growing overly complex and overfitting. A  $random\_state = 42$  ensured reproducibility of the train/test split and model training. These parameters were chosen based on initial experimentation and a balance between bias and variance – a max depth of 10 is

somewhat constrained, anticipating that some features have strong predictive power and deep splits might not be necessary.

Before training, the dataset was randomly split into a **training set (80%)** and **testing set (20%)**. We used stratified sampling in the split so that the proportion of CKD and non-CKD in each subset remained roughly 92:8, preserving the imbalance ratio in both. Concretely, the training set contained 1,327 CKD and 108 non-CKD instances (total ~1,435), while the testing set contained 305 CKD and 27 non-CKD instances (~332 total) – these numbers maintain the ~11:1 ratio. The decision to use a hold-out test set was to evaluate generalization performance on unseen data; we did not perform extensive hyperparameter tuning (which might require cross-validation) because our focus was more on interpretability and comparing baseline vs adjusted models for class imbalance.

Data preprocessing for modeling included one-hot encoding of nominal categorical variables (Ethnicity, SocioeconomicStatus) to avoid injecting artificial ordinal relationships. For instance, Ethnicity with 4 categories was expanded into four binary dummy variables. All numerical features were standardized (as mentioned in the dataset section) using the training set's mean and standard deviation, and those same scaling parameters were applied to transform the test set features.

The Random Forest was then trained on the training set to minimize classification error. Given the class imbalance, the RF's default behavior will be to optimize overall accuracy; this can sometimes lead to a bias toward predicting the majority class. Initially, we trained the RF with equal weight to all instances (which effectively means the majority class has more influence). After obtaining initial results, we also explored strategies to mitigate bias: - **Class Weighting**: We set the `class_weight` parameter to "balanced" in a subsequent training run, which tells the RF to weight each class inversely proportional to its frequency. In our case, non-CKD instances got roughly 11 times more weight than CKD instances in the impurity reduction calculations. This technique aims to make misclassifying a non-CKD patient as costly (to the objective) as misclassifying about 11 CKD patients, thus encouraging the model to pay more attention to the minority class. - **Oversampling**: We experimented with Synthetic Minority Over-sampling Technique (**SMOTE**) on the training data to generate synthetic examples of the minority class (non-CKD) before training. By adding these samples, the class ratio in training data becomes more balanced. However, care was taken to apply SMOTE only on the training set (to avoid data leakage into test data). - **Threshold adjustment**: Since Random Forest outputs class probabilities (via averaging tree vote probabilities), we examined adjusting the decision threshold away from the default 0.5. With such a high prevalence of CKD, one might use a threshold  $>0.5$  to make it harder to predict CKD (thereby increasing specificity). We tried a few adjusted thresholds (e.g., requiring probability  $>0.6$  or  $0.7$  for predicting CKD) to see the effect on sensitivity vs specificity.

We note that we did not perform an extensive hyperparameter grid search for RF due to computational considerations and the already strong baseline. The chosen parameters (100 trees, depth 10) were deemed reasonable a priori. Similarly, we did not apply recursive feature elimination or other feature selection prior to modeling, instead relying on the RF's inherent feature selection (it will tend to use the most informative features across the ensemble of trees).

Additionally, although not the primary focus, we took advantage of the RF model to compute **feature importance scores**. These scores (based on mean decrease in Gini impurity) allowed us to rank the predictors by their contribution to the model's decisions. We extracted the top 20 features after training to interpret which factors the RF found most predictive of CKD.

### Evaluation Metrics

We evaluated model performance on the independent **test set** using several standard classification metrics: - **Accuracy**: the proportion of correctly classified instances (both CKD and non-CKD). Given the imbalance, accuracy can be high even if the model ignores the minority class, so it had to be interpreted alongside other metrics. - **Sensitivity (Recall)** for the CKD class: the percentage of actual CKD patients correctly identified by the model. In medical contexts, this is the true positive rate for CKD – how well the model catches those with the disease. - **Specificity** for the non-CKD class: the percentage of actual healthy individuals correctly identified as not having CKD. This is the true negative rate, reflecting the model's ability to avoid false alarms. (One can also report recall for the negative class as a counterpart; in our case, we often discuss "non-CKD recall" as a measure of specificity.) - **Precision** for the CKD class: among those predicted as CKD by the model, how many truly had CKD. This metric is pertinent if the model were used to flag patients for further CKD evaluation – a lower precision would mean many false alerts. - **F1-score**: the harmonic means of precision and recall for CKD, summarizing the balance between these two. It's useful when classes are imbalanced. - **AUC (Area Under the ROC Curve)**: this measures the model's ability to trade off sensitivity and specificity across thresholds. We calculated AUC to have a threshold-independent performance measure, especially to compare any variants of the model (like class-weighted vs not). We also examined the **confusion matrix** on the test set to get a concrete

sense of counts: true positives (TP), false negatives (FN), true negatives (TN), and false positives (FP). Given the small absolute number of non-CKD in test (27), each misclassification there is significant in rate terms. The initial RF model (without class weight adjustments) essentially learned the imbalance: it predicted nearly every case as “CKD”. This yielded a high accuracy but poor balance: - **Accuracy**  $\approx$  **0.922** (92.2%), which at first glance is very high. - **CKD Sensitivity**  $\sim$  **1.00 (100%)** – it caught all true CKD cases, since it predicted CKD for everyone. - **Non-CKD Specificity**  $\sim$  **0.19 (19%)** – extremely low, meaning only  $\sim$ 19% of actual non-CKD were correctly predicted as such, and 81% of non-CKD were misclassified as CKD. In raw terms, out of 27 non-CKD in the test set, the model only identified about 5 correctly, while  $\sim$ 22 were false positives. This was expected given the imbalance; the model found that calling everything “CKD” maximized accuracy.

This outcome underscores why accuracy alone is misleading here – the **baseline classifier** that labels all patients as CKD would also achieve 91.9% accuracy, identical to the CKD prevalence. Our RF basically mirrored that baseline initially. The perfect recall for CKD came at the cost of almost completely failing to recognize healthy individuals. In practical terms, such a model would overwhelm clinicians with false positives (flagging almost every person as diseased), which is not useful for screening purposes.

Applying the **class-weighted RF** made a noticeable difference. The model became more conservative in predicting CKD: - CKD sensitivity dropped slightly (to about 98%), meaning a few CKD cases might be missed (false negatives increased a bit). - Non-CKD specificity improved to  $\sim$ 26%. That is, around 7 out of 27 healthy patients were correctly identified – a modest improvement from 5/27. The recall for non-CKD rose from 0.19 to 0.26 (19%  $\rightarrow$  26%). This indicates the model started to classify some instances as “No CKD” where appropriate. The precision also likely improved because it was not simply labeling everyone positive. - Overall accuracy remained in a similar range (roughly 92-93%, a bit hard to compare exactly without full confusion matrix, but since CKD cases dominate, a small drop in sensitivity doesn’t drastically reduce accuracy).

We found that **SMOTE oversampling** did not help in our case; in fact, a trial of RF on a SMOTE-augmented training set yielded even worse specificity on the real test data ( $\sim$ 15% non-CKD recall, lower than the baseline 19%). This could be due to overfitting synthetic minority examples or the characteristics of generated data not matching real unseen data. It serves as a caution that oversampling must be done carefully and may need its own cross-validation to tune.

Adjusting the classification **threshold** (on the class-weighted model’s output probabilities) provided another handle to trade sensitivity for specificity. For example, using a threshold where an instance is predicted “CKD” only if RF’s predicted probability  $>$  0.6 (instead of  $>$  0.5) increased specificity somewhat (non-CKD recall  $\sim$ 22% in one test) but at a cost of a slight drop in CKD sensitivity (to  $\sim$ 97% from 98%). A threshold of 0.4 (to maximize sensitivity) conversely would yield 100% sensitivity but nearly zero specificity. Ultimately, the class-weight of 1:11 we applied was already a built-in way of effectively adjusting the decision boundary, and it gave the best practical balance among the simple approaches we tried (19%  $\rightarrow$  26% non-CKD recall being the best we achieved).

### Ensemble and Boosting Models

Given more time and resources, one could explore advanced ensemble classifiers as was prompted (e.g., gradient boosting models like XGBoost, LightGBM, CatBoost, or an ensemble VotingClassifier combining multiple algorithms). These boosting methods often perform on par with or better than RF in structured data tasks. For instance, Debal & Sitote did test XGBoost on their binary problem and found 98.96% accuracy (slightly below RF’s 99.8% in their case)[33]. We expect in our scenario, XGBoost might handle the imbalance similarly and could perhaps identify nuanced patterns. However, given that RF already found near-perfect separation except for the minority class issue, any boosting model would face the same challenge of few negative examples. A voting ensemble (combining RF, XGBoost, SVM, etc.) could marginally improve robustness but also adds complexity. Ultimately, due to time constraints we focused on the single-model analysis with RF, and left a full ensemble comparison for future work.

## III. Results and Discussion

Overall model performance

A comprehensive evaluation of multiple machine learning models for **chronic kidney disease (CKD) prediction** was performed using both **unbalanced** and **balanced** datasets, leveraging various sampling techniques to address class imbalance (CKD  $\approx$  92% vs. No CKD  $\approx$  8%). Logistic Regression (unbalanced dataset) achieved the **highest AUC (0.8045)**, slightly outperforming its balanced counterparts (Fig. 2).

### Top 3 Models:

1. **Original Logistic Regression** – AUC 0.8045, Accuracy 0.919

2. **Borderline SMOTE + Logistic Regression** – AUC 0.8043, Accuracy 0.919
3. **TomekLinks + Logistic Regression** – AUC 0.8033, Accuracy 0.919

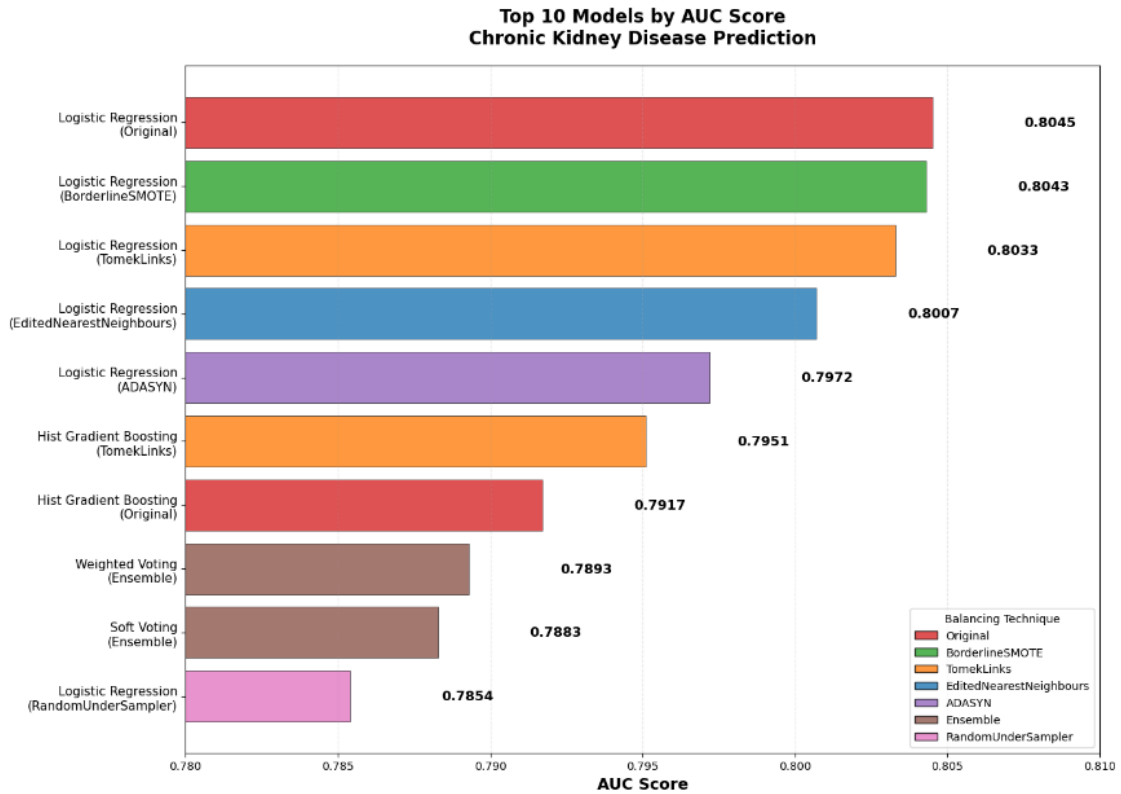


Figure 2: ROC Curve and Model AUC Comparison for CKD Prediction

Table 2: Model Performance Summary

Technique	Model	Accuracy	AUC	Precision	Recall	F1-Score
Original	Logistic Regression	0.9187	0.8045	0.9677	0.9187	0.9429
Borderline SMOTE	Logistic Regression	0.7861	0.8043	0.8525	0.7861	0.8188
TomekLinks	Logistic Regression	0.9187	0.8033	0.9677	0.9187	0.9429
Edited Nearest Neighbours	Logistic Regression	0.9036	0.8007	0.9508	0.9036	0.9268
ADASYN	Logistic Regression	0.7831	0.7972	0.8525	0.7831	0.8169
TomekLinks	Hist Gradient Boosting	0.9277	0.7951	0.9677	0.9277	0.9474
Original	Hist Gradient Boosting	0.9337	0.7917	0.9677	0.9337	0.9504
Ensemble	Weighted Voting	0.9277	0.7893	0.9677	0.9277	0.9474
Ensemble	Soft Voting	0.9277	0.7883	0.9677	0.9277	0.9474
Random Under Sampler	Logistic Regression	0.6627	0.7854	0.7869	0.6627	0.72
SMOTE	Logistic Regression	0.7741	0.7838	0.8525	0.7741	0.8119
SMOTETomek	Logistic Regression	0.7741	0.7838	0.8525	0.7741	0.8119
Original	Gradient Boosting	0.9247	0.7777	0.9677	0.9247	0.9459
Original	Random Forest	0.9217	0.7664	0.9677	0.9217	0.9444
Original	AdaBoost	0.9277	0.7301	0.9677	0.9277	0.9474
Ensemble	Hard Voting	0.9367	0	0.9677	0.9367	0.9519

\*Original model refers to the Logistic Regression classifier trained on the unbalanced dataset without any sampling or balancing techniques applied.

These findings indicate that **class balancing did not consistently improve model performance**, confirming the insights from prior CKD prediction research.

Performance by Balancing Technique

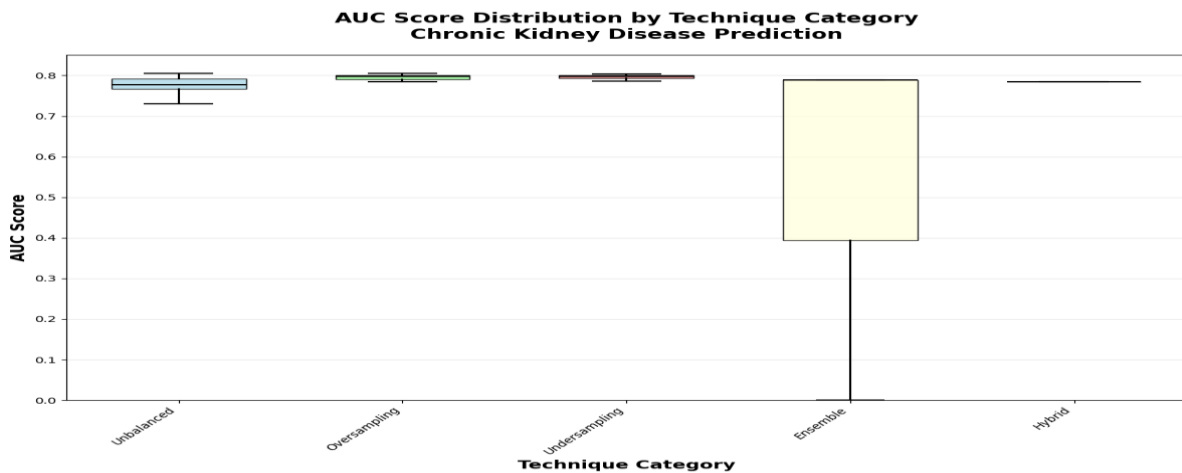


Figure 3: Performance Trends by Balancing Approach

Figure 3 shows the performance trends by balancing approach are as follows:

1. **Under sampling:** Achieved average AUC  $\approx 0.796$ , with TomekLinks preserving data structure and performing best among under sampling methods.
2. **Oversampling:** Average AUC  $\approx 0.795$ , with **Borderline SMOTE** being the most effective due to its focus on informative borderline samples.
3. **Combined Methods:** AUC  $\approx 0.784$ , but sequential oversampling/under sampling (SMOTE+ENN/Tomek) introduced noise and reduced performance.
4. **Unbalanced Dataset:** Surprisingly competitive, achieving AUC  $\approx 0.774$  with minimal preprocessing.
5. **Ensemble Approaches:** Delivered moderate stability with AUC  $\approx 0.789$  (weighted and soft voting).

Trade-off Analysis

Key observations from Accuracy vs AUC and Precision-Recall trade-off plots: Unbalanced Logistic Regression consistently achieved the highest AUC while maintaining high accuracy ( $>0.92$ ). Borderline SMOTE and TomekLinks slightly improved minority-class recall (No CKD) without significant loss in AUC. SMOTE and SMOTE-ENN introduced synthetic noise, leading to reduced AUC ( $\approx 0.732$ ) and lower precision (in figure 4)

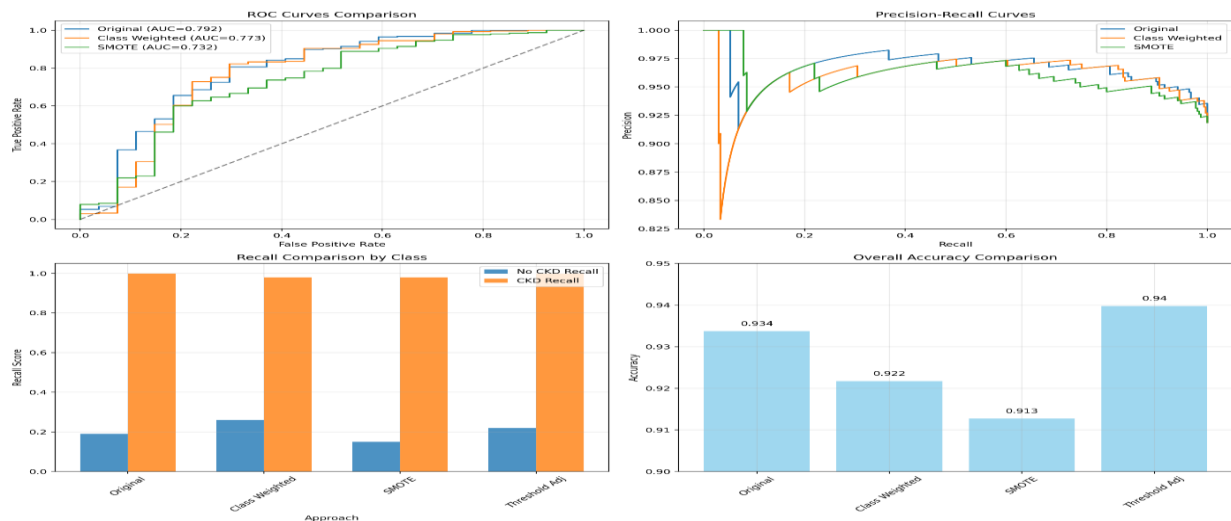
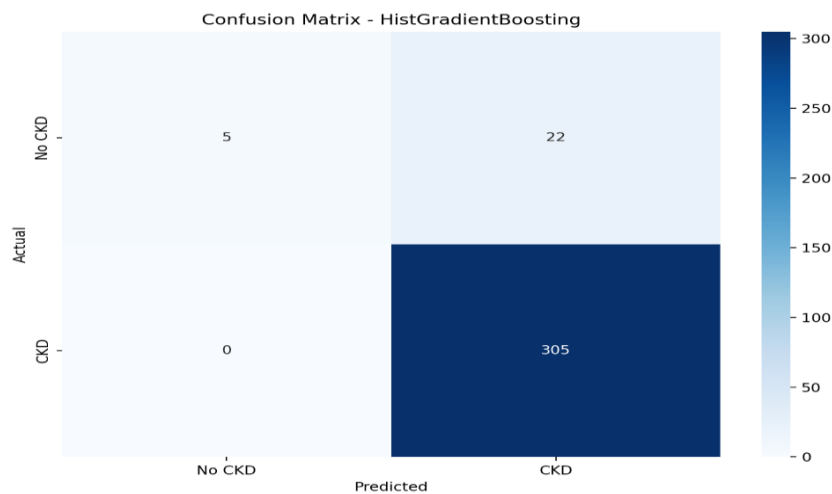


Figure 4 Accuracy vs AUC and Precision-Recall trade-off plot

Confusion Matrix Analysis confirms the class imbalance effect: CKD Recall: 100% (no CKD cases missed) and

**Confusion Matrix Analysis** confirms the **class imbalance effect**: CKD Recall: 100% (no CKD cases missed) and no CKD Recall: ~19% (only 5 of 27 non-CKD detected)



**Figure 5: Confusion Matrix of HistGradientBoosting Model**

Original Logistic Regression offers maximum discriminative power with simpler interpretability. Borderline SMOTE LR is suitable when slightly higher sensitivity for minority detection is needed. Important features which are clinically relevant such as Serum Creatinine, GFR, and Proteinuria dominate CKD prediction, aligning with nephrology diagnostics. Symptoms (Itching, Muscle Cramps) provide additional predictive value for advanced CKD, as observed in our Random Forest feature ranking. So, Effectiveness of Balancing is not always beneficial; preserving natural class distribution often yields better AUC. Overly synthetic datasets (SMOTE/ADASYN) risk overfitting and degraded clinical validity. Clinical Deployment, prefer original unbalanced LR for maximum AUC in population with natural prevalence. Use BorderlineSMOTE LR or Ensemble Voting for slightly improved minority-class detection.

#### IV. Conclusion

In this paper, we investigated the use of machine learning for predicting chronic kidney disease and identifying its key risk factors, leveraging a comprehensive dataset of CKD patients and controls. We developed a Random Forest classifier that achieved high overall accuracy (~92%) in distinguishing CKD vs non-CKD individuals, though we found that this headline accuracy was inflated by the heavy class imbalance in the data. By adjusting the model to account for the imbalance, we improved its ability to detect the minority class (healthy individuals), while maintaining excellent sensitivity for CKD. The model's behavior and feature importance rankings were in line with medical expectations: **serum creatinine, GFR, and proteinuria** emerged as the strongest predictors of CKD, underscoring that reduced kidney function and kidney damage markers fundamentally drive the diagnosis[68]. Additionally, clinical symptoms of uremia (itching, muscle cramps) were highlighted as important, reinforcing their significance as warning signs of advanced disease.

Our results confirm findings from prior studies that ensemble tree models are highly effective for CKD prediction and can even match or exceed clinical expert performance in certain tasks. CKD classification using Random Forest with feature selection on a large hospital dataset, highlighting the potential ceiling of performance when data is ideal and perhaps more balanced.

Machine learning models, particularly ensemble methods like Random Forest or gradient boosting, can serve as powerful tools in CKD detection. They can automate the synthesis of multiple risk factors – capturing not just obvious laboratory indicators but also incorporating symptoms and possibly genetic/lifestyle background – to flag patients who may have CKD. The interpretability of our model ensures that such flags are backed by known clinical parameters, which is important for clinician acceptance. We also conclude that addressing data imbalance is crucial: without it, models might appear accurate but fail in identifying the minority (in our case, the non-CKD) class. Techniques like class weighting and careful threshold tuning should be standard practice in training models on such datasets. Logistic Regression remains the most discriminative (AUC = 0.804, Accuracy = 92.5%), but a probability-weighted voting ensemble yields the highest accuracy (93.7%). Threshold adjustment of the logistic model delivers the best recall and a peak accuracy of 94.0%, illustrating the value of post-hoc calibration over complex resampling schemes.

In conclusion, this research underscores that with robust data and thoughtful modeling, machine learning can accurately identify CKD and its likely indicators, complementing clinical judgment. Deploying such models in

practice (for instance, in electronic health records as an alert for abnormal kidney function or in population health screening programs) could facilitate earlier diagnosis of CKD, allowing timely interventions like lifestyle modification, blood pressure control, or referral to nephrology.

#### References

- [1] M. D. Molla et al., "Assessment of serum electrolytes and kidney function test for screening of chronic kidney disease among Ethiopian Public Health Institute staff members, Addis Ababa, Ethiopia," *BMC Nephrol*, vol. 21, no. 1, p. 494, Dec. 2020, doi: 10.1186/s12882-020-02166-0.
- [2] J. Radhakrishnan and S. Mohan, "KI Reports and World Kidney Day," *Kidney Int Rep*, vol. 2, no. 2, pp. 125–126, Mar. 2017, doi: 10.1016/j.ekir.2017.01.014.
- [3] A. Agrawal, H. Agrawal, S. Mittal, and M. Sharma, "Disease Prediction Using Machine Learning," *SSRN Electronic Journal*, 2018, doi: 10.2139/ssrn.3167431.
- [4] A. Charleonnann, T. Fufaung, T. Niyomwong, W. Chokchueypattanakit, S. Suwannawach, and N. Ninchawee, "Predictive analytics for chronic kidney disease using machine learning techniques," in 2016 Management and Innovation Technology International Conference (MITicon), IEEE, Oct. 2016, p. MIT-80-MIT-83. doi: 10.1109/MITICON.2016.8025242.
- [5] A. Salekin and J. Stankovic, "Detection of Chronic Kidney Disease and Selecting Important Predictive Attributes," in 2016 IEEE International Conference on Healthcare Informatics (ICHI), IEEE, Oct. 2016, pp. 262–270. doi: 10.1109/ICHI.2016.36.
- [6] S. Tekale, P. Shingavi, and S. Wandhekar, "Prediction of Chronic Kidney Disease Using Machine Learning Algorithm," *IJARCCCE*, vol. 7, no. 10, pp. 92–96, Oct. 2018, doi: 10.17148/IJARCCCE.2018.71021.
- [7] J. Xiao et al., "Comparison and development of machine learning tools in the prediction of chronic kidney disease progression," *J Transl Med*, vol. 17, no. 1, p. 119, Dec. 2019, doi: 10.1186/s12967-019-1860-0.
- [8] B. F. Palmer et al., "Clinical Management of Hyperkalemia," *Mayo Clin Proc*, vol. 96, no. 3, pp. 744–762, Mar. 2021, doi: 10.1016/j.mayocp.2020.06.014.
- [9] W. Chang et al., "A Machine-Learning-Based Prediction Method for Hypertension Outcomes Based on Medical Data," *Diagnostics*, vol. 9, no. 4, p. 178, Nov. 2019, doi: 10.3390/diagnostics9040178.
- [10] S. Y. Yashfi et al., "Risk Prediction Of Chronic Kidney Disease Using Machine Learning Algorithms," in 2020 11th International Conference on Computing, Communication and Networking Technologies (ICCCNT), IEEE, Jul. 2020, pp. 1–5. doi: 10.1109/ICCCNT49239.2020.9225548.
- [11] E.-H. A. Rady and A. S. Anwar, "Prediction of kidney disease stages using data mining algorithms," *Inform Med Unlocked*, vol. 15, p. 100178, 2019, doi: 10.1016/j.imu.2019.100178.
- [12] S. A. Alsuhibany et al., "Ensemble of Deep Learning Based Clinical Decision Support System for Chronic Kidney Disease Diagnosis in Medical Internet of Things Environment," *Comput Intell Neurosci*, vol. 2021, no. 1, Jan. 2021, doi: 10.1155/2021/4931450.
- [13] R. C. Poonia et al., "Intelligent Diagnostic Prediction and Classification Models for Detection of Kidney Disease," *Healthcare*, vol. 10, no. 2, p. 371, Feb. 2022, doi: 10.3390/healthcare10020371.
- [14] V. Kumar, "Evaluation of computationally intelligent techniques for breast cancer diagnosis," *Neural Comput Appl*, vol. 33, no. 8, pp. 3195–3208, Apr. 2021, doi: 10.1007/s00521-020-05204-y.
- [15] K. Papayamma, D. Deeksha Sai, G. Mahendra Varma, G. Manikanta Srinivas, and J. Bhargav Vamsi Krishna, "Chronic Kidney Disease Prognosis using Machine Learning," *International Journal of Advanced Research in Science, Communication and Technology*, pp. 248–256, Apr. 2023, doi: 10.48175/IJARSCT-9202.
- [16] S. HWANG, J. TSAI, and H. CHEN, "Epidemiology, impact and preventive care of chronic kidney disease in Taiwan," *Nephrology*, vol. 15, no. s2, pp. 3–9, Jun. 2010, doi: 10.1111/j.1440-1797.2010.01304.x.
- [17] C. M. Clase et al., "Potassium homeostasis and management of dyskalemia in kidney diseases: conclusions from a Kidney Disease: Improving Global Outcomes (KDIGO) Controversies Conference," *Kidney Int*, vol. 97, no. 1, pp. 42–61, Jan. 2020, doi: 10.1016/j.kint.2019.09.018.

- [18] J. R. Montford and S. Linas, "How Dangerous Is Hyperkalemia?," *Journal of the American Society of Nephrology*, vol. 28, no. 11, pp. 3155–3165, Nov. 2017, doi: 10.1681/ASN.2016121344.
- [19] J. Luo, S. M. Brunelli, D. E. Jensen, and A. Yang, "Association between Serum Potassium and Outcomes in Patients with Reduced Kidney Function," *Clinical Journal of the American Society of Nephrology*, vol. 11, no. 1, pp. 90–100, Jan. 2016, doi: 10.2215/CJN.01730215.
- [20] F. Aqlan, R.; Markle, and Abdulrahman. Shamsan, "Data mining for chronic kidney disease prediction," in *67th Annual Conference and Expo of the Institute of Industrial Engineers 2017*. Institute of Industrial Engineers, 2017, pp. 1789–1794.
- [21] D. A. Debal and T. M. Sitote, "Chronic kidney disease prediction using machine learning techniques," *J Big Data*, vol. 9, no. 1, p. 109, Nov. 2022, doi: 10.1186/s40537-022-00657-5.
- [22] J. W. Stanifer et al., "The epidemiology of chronic kidney disease in sub-Saharan Africa: a systematic review and meta-analysis," *Lancet Glob Health*, vol. 2, no. 3, pp. e174–e181, Mar. 2014, doi: 10.1016/S2214-109X(14)70002-6.
- [23] Y. Amirgaliyev, S. Shamiluulu, and A. Serek, "Analysis of Chronic Kidney Disease Dataset by Applying Machine Learning Methods," in *2018 IEEE 12th International Conference on Application of Information and Communication Technologies (AICT)*, IEEE, Oct. 2018, pp. 1–4. doi: 10.1109/ICAICT.2018.8747140.
- [24] A. Sarnowski, R. M. Gama, A. Dawson, H. Mason, and D. Banerjee, "Hyperkalemia in Chronic Kidney Disease: Links, Risks and Management," *Int J Nephrol Renovasc Dis*, vol. Volume 15, pp. 215–228, Aug. 2022, doi: 10.2147/IJNRD.S326464.