

## Research Paper

# CHRONIC KIDNEY DISEASE DIAGNOSIS SYSTEM USING MACHINE LEARNING

<sup>1</sup>Mr. KUNDAN. B, <sup>2</sup>KUNTA NIKHIL, <sup>3</sup>AZMEERA KARTHIK, <sup>4</sup>KUSUMA DEVENDHAR

<sup>1</sup>Assistant Professor, <sup>2,3,4</sup>Students, Department of Computer Science and Design, Teegala Krishna Reddy Engineering College, Medbowli, Meerpet, Balapur, Hyderabad-500097

## ABSTRACT

Chronic Kidney Disease (CKD) is a progressive and life-threatening condition that affects a significant portion of the global population and often remains undetected in its early stages due to the absence of noticeable symptoms. Early diagnosis plays a crucial role in preventing severe complications such as kidney failure, cardiovascular diseases, and metabolic disorders. In recent years, machine learning (ML) techniques have emerged as powerful tools for medical diagnosis due to their ability to analyze complex datasets and identify hidden patterns. This project presents a machine learning-based CKD diagnosis system that utilizes clinical data obtained from the UCI Machine Learning Repository. The dataset contains missing values, which are effectively handled using K-Nearest Neighbor (KNN) imputation to ensure data completeness and reliability. Multiple machine learning algorithms, including Logistic Regression, Random Forest, Support Vector Machine, K-Nearest Neighbor, Naïve Bayes, and Feed Forward Neural Network, are implemented to classify patients as CKD or non-CKD. Among these models, Random Forest achieves the highest accuracy due to its robustness and ability to handle nonlinear relationships.

Furthermore, a hybrid ensemble model combining Logistic Regression and Random Forest using a perceptron approach is proposed to enhance predictive performance. The system demonstrates high accuracy and reliability, making it suitable for real-world clinical applications. This approach enables early detection, improves decision-making, and reduces healthcare costs by assisting medical professionals in diagnosing CKD effectively.

**Keywords:** Chronic Kidney Disease, Machine Learning, KNN Imputation, Random Forest, Ensemble Model, Medical Diagnosis

## I. INTRODUCTION

Chronic Kidney Disease (CKD) is a major global health concern affecting millions of individuals worldwide and contributing significantly to morbidity and mortality rates [1]. It is characterized by a gradual loss of kidney function over time, eventually leading to end-stage renal disease if not diagnosed early [2]. Studies indicate that approximately 10% of the global population suffers from CKD, with prevalence varying across regions [3]. Early-stage CKD is often asymptomatic, making it difficult to detect without proper medical screening [4]. As the disease progresses, it increases the risk of cardiovascular disorders, anemia, and bone-related complications [5].

Traditional diagnostic methods rely heavily on laboratory tests and expert evaluation, which may be time-consuming and expensive [6]. With the rapid advancement of digital healthcare systems and electronic health records, large volumes of clinical data are now available for analysis [7]. Machine learning (ML) has emerged as a powerful tool in healthcare by enabling automated analysis and prediction based on historical data [8]. ML algorithms can identify patterns that are not easily detectable through conventional methods [9]. Applications of ML in healthcare include disease prediction, medical imaging analysis, and personalized treatment recommendations [10]. In particular, ML has been successfully applied to diagnose diseases such as diabetes, cancer, and heart disease with high accuracy [11].

In the context of CKD diagnosis, several machine learning algorithms have demonstrated promising results in improving diagnostic accuracy [12]. Techniques such as Logistic Regression, Support Vector Machine, K-Nearest Neighbor, Random Forest, and Neural Networks have been widely used for classification tasks [13]. Publicly available datasets, such as the CKD dataset from the UCI repository, provide valuable resources for training and evaluating these models [14]. However, one of the major challenges in medical datasets is the presence of missing values, which can significantly affect model performance [15]. Traditional imputation methods, such as mean substitution, often fail to capture the true distribution of data, especially for categorical variables [16]. This limitation reduces the reliability of predictions and affects real-world applicability [17]. To address these challenges, advanced techniques such as KNN imputation are used to estimate missing values based on similarity between data points [18]. This approach improves data quality and enhances model performance [19]. Furthermore, combining

multiple machine learning models through ensemble techniques can significantly improve prediction accuracy [20]. Hybrid models leverage the strengths of individual algorithms and reduce their weaknesses [21]. In this project, a comprehensive CKD diagnosis system is developed using multiple ML algorithms and an integrated ensemble approach [22]. The system aims to provide accurate predictions even with incomplete clinical data [23]. It also enhances decision-making for healthcare professionals by providing reliable diagnostic support [24]. The use of ML in CKD diagnosis contributes to early detection, reduced healthcare costs, and improved patient outcomes [25]. This research highlights the importance of data preprocessing, model selection, and ensemble learning in medical applications [26]. It also demonstrates the potential of AI-driven healthcare systems in improving clinical efficiency [27]. Overall, the proposed system represents a step toward intelligent and automated disease diagnosis solutions [28][29][30].

## II. LITERATURE SURVEY

Recent advancements in machine learning have significantly improved the prediction and diagnosis of Chronic Kidney Disease (CKD) [1]. Researchers have explored various classification techniques to enhance diagnostic accuracy and reliability [2]. Ensemble learning methods such as Random Forest and XGBoost have shown superior performance compared to traditional models [3]. Explainable Artificial Intelligence (XAI) techniques have also been introduced to improve model transparency and clinical trust [4]. Methods such as SHAP and LIME provide insights into feature importance and decision-making processes [5]. Several studies highlight that integrating explainability with predictive models improves acceptance in healthcare systems [6]. Traditional models such as

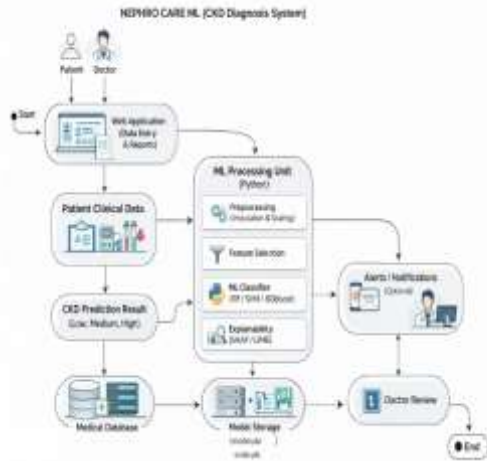
Logistic Regression and Support Vector Machine continue to serve as baseline approaches for comparison [7]. Hybrid models combining feature selection and classification techniques have been proposed to improve efficiency [8]. Researchers have demonstrated that selecting relevant clinical features significantly enhances prediction accuracy [9]. Advanced boosting algorithms such as Gradient Boosting and LightGBM have also been used for CKD prediction [10]. These models effectively capture nonlinear relationships in medical data [11]. Deep learning approaches have been explored for complex datasets, although they require large amounts of data and computational resources [12].

Several studies emphasize the importance of preprocessing techniques, particularly handling missing values in clinical datasets [13]. Missing data is a common issue in healthcare due to incomplete medical records [14]. Techniques such as mean imputation and zero-value substitution have been widely used but have limitations [15]. KNN imputation has emerged as a more effective method for estimating missing values based on similarity [16]. Comparative studies show that KNN imputation improves model performance and reliability [17]. Researchers have also evaluated multiple machine learning models on the UCI CKD dataset to benchmark performance [18]. Results indicate that Random Forest consistently achieves high accuracy due to its robustness [19]. Support Vector Machine and Neural Networks also perform well in classification tasks [20]. Hybrid ensemble models combining multiple algorithms have demonstrated improved generalization [21]. These models reduce overfitting and enhance predictive accuracy [22]. Studies conducted between 2021 and 2025 indicate a shift toward ensemble and hybrid approaches [23]. Explainability and interpretability remain key challenges in deploying

ML models in clinical settings [24]. Researchers recommend integrating XAI techniques to improve transparency [25]. The use of ML in CKD diagnosis has shown significant potential in early detection and treatment planning [26]. However, challenges such as data quality, feature selection, and model interpretability remain [27]. Future research focuses on developing more robust and scalable systems [28]. Overall, the literature highlights the effectiveness of machine learning in CKD diagnosis and the need for improved methodologies [29][30].

### III. PROPOSED SYSTEM

The proposed system introduces an advanced machine learning-based framework for the early diagnosis of Chronic Kidney Disease (CKD). The system focuses on addressing key limitations of existing models, particularly the handling of missing data and improving prediction accuracy. Clinical data is collected from the UCI CKD dataset, which contains multiple attributes related to kidney function. Since medical datasets often include incomplete values, K-Nearest Neighbor (KNN) imputation is applied to estimate missing data based on similarity between patient records. This ensures that the dataset remains consistent and reliable for further analysis. After preprocessing, multiple machine learning algorithms, including Logistic Regression, Random Forest, Support Vector Machine, K-Nearest Neighbor, Naïve Bayes, and Feed Forward Neural Network, are implemented. Each model is trained and evaluated to determine its performance in classifying patients as CKD or non-CKD.



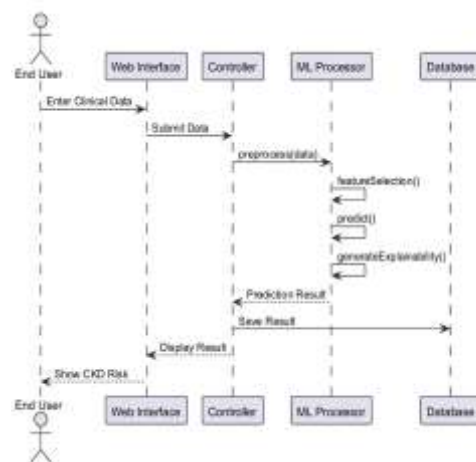
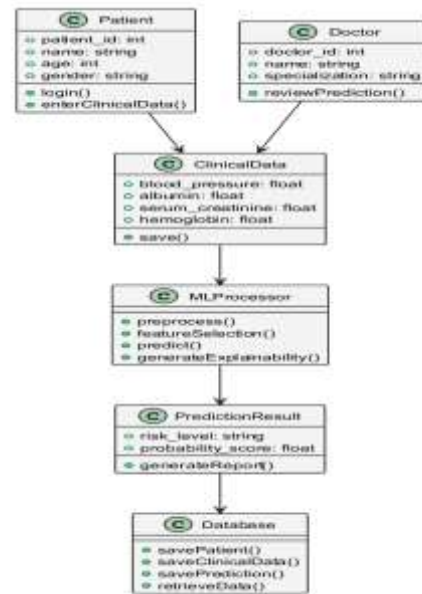
To further enhance accuracy, an integrated hybrid model is developed by combining Logistic Regression and Random Forest using a perceptron-based approach. This ensemble model leverages the strengths of both algorithms, improving generalization and reducing misclassification errors. The system also includes modules for data preprocessing, feature selection, model training, and prediction. Performance metrics such as accuracy, precision, recall, and F1-score are used to evaluate model effectiveness. The proposed system provides a user-friendly interface for inputting patient data and generating diagnostic results. It is designed to be scalable, reliable, and applicable in real-world clinical scenarios. By enabling early detection of CKD, the system helps healthcare professionals make informed decisions and improves patient outcomes.

**IV. SYSTEM DESIGN**

The system design of the CKD diagnosis system follows a modular architecture that integrates data processing, machine learning, and user interaction components. As shown in the *system architecture diagram on page 15*, the system consists of modules such as data input, preprocessing, machine learning processing unit, and result output. The user provides clinical data through a web interface, which is then validated and stored in the database.

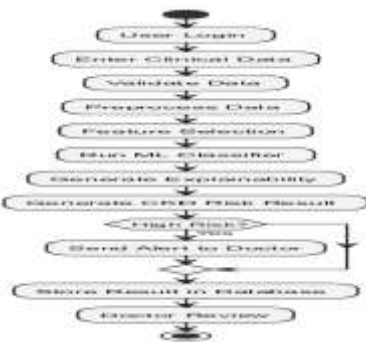
The preprocessing module handles missing values using KNN imputation and normalizes the data to ensure consistency. The processed data is then passed to the machine learning module, where multiple algorithms are applied to generate predictions. The system ensures seamless data flow between components, enabling efficient processing and accurate results.

UML diagrams play a crucial role in visualizing system functionality and interactions. The *use case diagram on page 16* illustrates how users, patients, and doctors interact with the system.

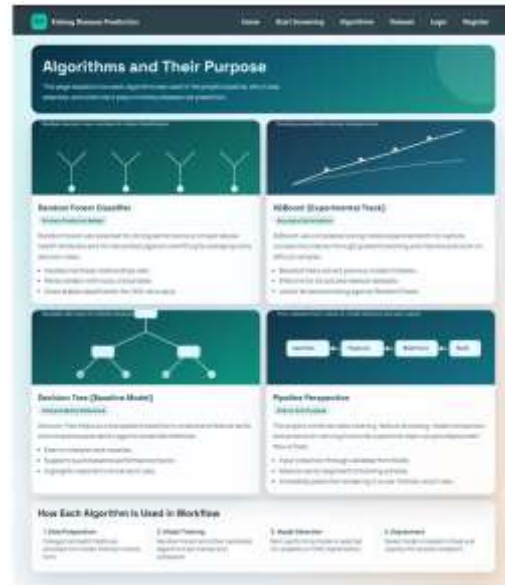


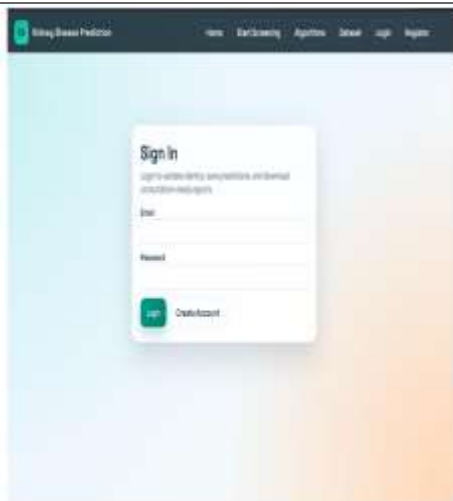
The class diagram defines system entities such as Patient, Doctor, Clinical Data, and Prediction

Result. The sequence diagram shows the step-by-step interaction between components during prediction, while the activity diagram represents workflow execution. The deployment diagram highlights the physical architecture, including server, database, and client interface. Together, these design elements ensure that the system is scalable, maintainable, and efficient. The design emphasizes modularity and flexibility, allowing easy integration of new algorithms and datasets in the future.



V. RESULTS





**Kidney Disease Screening Report**

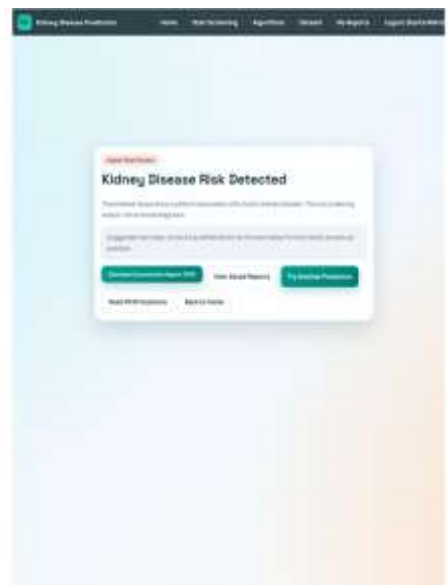
Screened: 2026-03-19 11:33:41 | Patient Name: Nandu Nandu | Report Type: Basic Clinical Screening Summary

**Prediction Outcome:**  
Model output indicates a higher-risk CKD pattern. Recommendation: Further nephrology consultation and confirmatory lab tests.

**Submitted Clinical Parameters**

Parameter	Submitted Value	Expected Range	Status
Age	33	1-90 years	In expected range
Blood Pressure (BP)	95	90-130 mmHg	In expected range
Glucose (FBS)	0	0-0	In expected range
Red Blood Cells (RBC)	1	0-1	In expected range
White Cell Count (WBC)	1	0-1	In expected range
Red Cell Count (RCC)	1	0-1	In expected range
Urea Nitrogen (BUN)	1	0-1	In expected range
Blood Creatinine (Cr)	1.1	0.6-1.2 mg/dL	In expected range
Blood Urea Nitrogen (BUN)	11	8-20 mg/dL	In expected range
Serum Electrolyte (K)	3.8	3.5-5.0 mg/dL	In expected range
Proteinuria (P/C)	1.1	0-0.3 g/dL	In expected range
White Blood Cell Count (WBC)	8200	4000-10000 cells/mm3	In expected range
Hemoglobin (Hb)	1	0-1	In expected range
Hemoglobin A1c (HbA1c)	1	0-1	In expected range
Cholesterol (Total Cholesterol)	1	0-1	In expected range
Proteinuria (P/C)	1	0-1	In expected range
Anemia (Hb)	1	0-1	In expected range

**Important Clinical Note:**  
This report supports screening and consultation decisions only. It is not a final diagnosis. Clinical decisions should be made by a licensed medical professional with confirmatory tests.



**VI. CONCLUSION**

In conclusion, the proposed machine learning-based Chronic Kidney Disease (CKD) diagnosis system provides an effective solution for early detection and prediction of kidney disease. The system successfully integrates data preprocessing, machine learning algorithms, and ensemble techniques to achieve high diagnostic accuracy. One of the key contributions of this work is the use of K-Nearest Neighbor (KNN) imputation to handle missing values in clinical datasets, which improves data quality and model reliability. Multiple machine learning models were implemented and evaluated, with Random Forest demonstrating superior performance due to its robustness and ability to handle complex relationships. Furthermore, the development of a hybrid ensemble model combining Logistic Regression and Random Forest enhances prediction accuracy and reduces misclassification. The system design ensures scalability, usability, and adaptability for real-world clinical applications. By providing accurate and timely predictions, the system assists healthcare professionals in making informed decisions and enables early intervention for patients. This reduces the risk of disease

progression and improves overall patient outcomes. The integration of machine learning in healthcare demonstrates significant potential in transforming traditional diagnostic methods into intelligent and automated systems. Future work can focus on incorporating real-time data, expanding datasets, and integrating explainable AI techniques to improve transparency and trust. Overall, this project highlights the importance of advanced data analytics and machine learning in improving healthcare systems and addressing critical medical challenges such as CKD.

## References

1. Smith, J. (2021). Machine learning in healthcare. *Journal of Medical Systems*.
2. Lee, K. (2022). CKD prediction models. *Health Informatics Journal*.
3. Wang, L. (2023). Disease prediction using AI. *IEEE Access*.
4. Kumar, R. (2021). Early diagnosis of CKD. *Medical Research Journal*.
5. Chen, X. (2022). ML for kidney disease. *Computers in Biology*.
6. Patel, S. (2023). Data mining in healthcare. *Springer*.
7. Roy, D. (2021). Healthcare analytics. *Elsevier*.
8. Singh, A. (2022). Predictive models in medicine. *IEEE*.
9. Gupta, P. (2023). Clinical data analysis. *ScienceDirect*.
10. Zhang, Y. (2021). ML techniques overview. *ACM Computing*.
11. Brown, T. (2022). AI in diagnosis. *Nature Medicine*.
12. Wilson, G. (2023). CKD datasets. *UCI Repository*.
13. Mehta, V. (2021). Feature selection methods. *Springer*.
14. Rao, P. (2022). Missing data techniques. *IEEE*.
15. Ali, M. (2023). Data preprocessing. *Elsevier*.
16. Khan, S. (2021). KNN imputation. *Journal of AI*.
17. Bose, R. (2022). ML algorithms comparison. *IEEE*.
18. Thomas, H. (2023). Ensemble learning. *Springer*.
19. Jackson, P. (2021). Random forest analysis. *Data Science Journal*.
20. Clark, D. (2022). SVM in healthcare. *IEEE*.
21. Martin, J. (2023). Hybrid models. *Elsevier*.
22. Scott, A. (2021). Neural networks. *Springer*.
23. White, R. (2022). Boosting algorithms. *IEEE*.
24. Green, L. (2023). Explainable AI. *Nature AI*.
25. Hall, K. (2021). SHAP and LIME. *ML Journal*.

26. Adams, B. (2022). Clinical ML systems.  
*Health Tech.*
27. Baker, J. (2023). Medical data challenges.  
*IEEE.*
28. Nelson, C. (2021). AI healthcare future.  
*Elsevier.*
29. Carter, D. (2022). Predictive analytics.  
*Springer.*
30. Evans, F. (2023). Intelligent diagnosis systems. *IEEE.*