Chronic Kidney Diseases Predications Using Machine Learning

M.Periyakaruppan¹, A.Bala Ayyappan², Dr.T.Gobinath³ ¹Dept of Computer Science Engineering ² Artificial Intelligence and Data Science (AIDS) ³Associate Professor, Dept of Computer Science Engineering ^{1, 2, 3} Chettinad College of Engineering and Technology, NH67, Karur-Trichy Highway, Puliyur CF PO, Karur, 639114, Tamilnadu, India.

Abstract- Kidney diseases are an increasingly important global health problem for millions of people a year, for which early detection and accurate prediction play a major role in improving patient outcomes and reducing the burden on healthcare systems. This project seeks to develop a predictive model of kidney disease diagnosis based on machine learning that assumes control over clinical data through highly advanced analytical techniques. The dataset was preprocessed so that missing values were handled, features standardized, and the most relevant predictors chosen. Exploratory data analysis (EDA) was also conducted to look for any type of patterns and relationship among clinical attributes. Two machine learning algorithms were used for the construction of the predictive models - Random Forest and Decision Tree. Training and validation of models through a structured approach ensure robust evaluation through metrics such as accuracy, precision, recall, and confusion matrices. Data visualization techniques were also employed to enhance interpretability and generate actionable insights in the dataset. Results The results have shown that the implemented models work well, with the Random Forest algorithm performing best in terms of accuracy and reliability. Therefore, this project serves as a possible application of ML in the clinical field to alert healthcare professionals when kidney diseases would commence. This early detection could ultimately serve to improve a patient's situation while providing an interface for predictive analytics in clinical decision-support systems. The Accuracy of LightGBM is 95%. XG Boost is 94%.

Keywords- Machine Learning, Artificial Intelligence in Healthcare, Patient Monitoring, Disease Classification, Packed Cell Volume (PCV)

I. INTRODUCTION

Kidneys are vital organs responsible for filtering waste, balancing body fluids, and regulating essential functions like blood pressure and red blood cell production. When the kidneys become damaged or lose their function over time, it leads to kidney diseases. These conditions can be acute (short-term) or chronic (long-term), with Chronic Kidney Disease (CKD) being the most common and serious form. If left untreated, CKD can progress to kidney failure, requiring dialysis or a transplant. Chronic Kidney Disease (CKD) is a global health concern affecting millions of individuals worldwide. It is characterized by a gradual loss of kidney function over time, leading to severe complications such as cardiovascular diseases and kidney failure if left untreated [1]. Early detection and diagnosis of CKD are crucial for effective intervention and treatment, thereby reducing morbidity and mortality rates associated with the disease [2]. With the advent of machine learning, predictive models have been increasingly utilized in the medical field to assist in disease detection and prognosis. These models leverage vast amounts of patient data to identify patterns and provide accurate predictions, aiding healthcare professionals in making informed decisions [3]. This research employs a Random Forest classifier, a widely used ensemble learning method known for its robustness and accuracy in handling medical datasets [4]. The model is trained on a dataset comprising key clinical attributes such as blood pressure, blood glucose levels hemoglobin, and other laboratory test results. By preprocessing and encoding categorical variables, the model enhances its ability to generalize and provide reliable predictions.

II. LITERATURE REVIEW

Chronic Kidney Disease (CKD) is a progressive medical condition where kidney function deteriorates over time. Traditional diagnostic methods rely on laboratory tests such as serum creatinine and estimated glomerular filtration rate (eGFR). However, the advent of machine learning (ML) techniques has significantly improved early detection, classification, and prognosis of CKD. Various studies have explored ML algorithms for automated CKD prediction models, focusing on feature selection, model accuracy, and clinical interpretability. Early diagnosis of kidney diseases is essential to prevent irreversible damage. However, many patients remain undiagnosed until kidney function is severely impaired. Therefore, innovative diagnostic approaches, including machine learning and predictive modeling, are

increasingly being explored to improve early detection and management of kidney diseases. Classification techniques, Feature selection, and Ensemble model are the most significant and vital tasks in machine learning and data mining. A lot of research has been conducted to apply data mining and machine learning classification technique, feature selection method and ensemble model on different medical datasets to classify disease datasets. Many of them show good classification accuracy. Classification algorithms are supervised learning method their class is known to predict the objective class level. The classification used to categorize datasets into training and test set. Data mining, classification is commonly used in healthcare application to classify patient dataset [5]Chronic Kidney Disease (CKD) is a progressive medical condition characterized by a gradual decline in kidney function. Traditional diagnostic methods primarily rely on laboratory tests such as serum creatinine levels and estimated glomerular filtration rate (eGFR). However, these approaches often fail to detect CKD in its early stages, leading to delayed diagnoses when kidney function is significantly impaired [6] Classification techniques, feature selection methods, and ensemble models are critical tasks in machine learning and data mining for medical diagnosis. Extensive research has been conducted to apply these techniques to various medical datasets, achieving promising classification accuracy[7].Classification algorithms, as supervised learning methods, operate by learning from labeled data to predict the target class. In healthcare applications, classification is widely used to categorize patient datasets into training and test sets, facilitating the development of robust diagnostic models[8]. The advent of machine learning (ML) techniques has significantly enhanced early detection, classification, and prognosis of CKD. Several studies have explored ML algorithms for automated CKD prediction models, emphasizing feature selection, model accuracy, and clinical interpretability[9]

III. METHODOLOGY

3.1 Dataset Collection

The Dataset here we use is the available CKD Dataset from UCI repository. It contains 400 samples of two different classes. Out of 25 attributes, 11 are numeric and 13 are nominal and one is class attribute. The data set contains number of missing values. Here the information of dataset uses the patient's data like age, blood pressure, specific gravity, albumin, sugar, red blood cells etc.[10]

3.2 Data Preprocessing

3.3 Machine Learning ModelsWe implemented the following classification models:

- 1. Random Forest Classifier An ensemble learning method that combines multiple decision trees.[12]
- 2. Support Vector Machines (SVM) Finds the best boundary between classes.[13]
- 3. K-Nearest Neighbors (KNN) Classifies based on the closest data points.
- Gradient Boosting Models (e.g., XGBoost, LightGBM,AdaBoost) – Improve performance using multiple weak learners.[14]

3.4 Flow diagram

Fig1: Flow Diagram represents a machine learningbased Chronic Kidney Disease (CKD) prediction pipeline. It begins with Data Selection and Loading, where relevant patient records containing medical attributes such as blood pressure, blood sugar levels, and creatinine levels are gathered and loaded into a processing environment. Next, in the Data Preprocessing step, the raw data is cleaned by handling missing values, normalizing numerical features, and encoding categorical variables to ensure consistency. Following this, Feature Selection is performed to identify the most important predictors using techniques like correlation analysis or Principal Component Analysis (PCA), which helps improve model accuracy and efficiency. In the Classification step, a machine learning model-such as Random Forest is applied to classify whether a patient has CKD or not based on the selected features.[15]



Figure1: Flow Diagram



Chronic Kidney Disease (CKD) presents a significant health burden, often going unrecognized, underdiagnosed, and inadequately treated. Many individuals with CKD experience distressing symptoms, including skin irritation (pruritus) and psychological challenges, which frequently occur in clusters. Addressing one symptom may help alleviate others, highlighting the importance of a comprehensive symptom management approach. To improve patient well-being, healthcare experts and kidney advocacy organizations emphasize effective symptom management strategies. Their goal is to delay the progression of CKD and reduce the need for dialysis. For patients requiring dialysis, a gradual transition and the option for home-based dialysis can offer more personalized care. Recent discussions in the medical community suggest replacing the term "kidney failure" with Kidney Dysfunction Requiring Dialysis (KDRD) to remove negative connotations and focus on patient-centered treatment approaches. A patient-centric approach ensures that individuals with CKD actively participate in their treatment decisions, with healthcare providers prioritizing their values, preferences, and unique needs. Key factors include equitable access to care, clear communication, and education tailored to the patient's health literacy. Through shared decision-making among patients, caregivers, and medical professionals, CKD management can be optimized to align with individual health goals.[16]

3.6 Heatmap

Fig2: The heatmap highlights key predictors for CKD, including hemoglobin (hemo), packed cell volume (pcv), specific gravity (sg), and serum creatinine (sc). Serum creatinine (0.81) and blood urea (0.70) show strong correlations with CKD, while lower specific gravity (-0.79) is associated with the disease. In contrast, features like appetite and potassium (pot) have minimal impact.[17]



Figure 2:HeatMap

| \mathbf{T}_{2} | ւե | 1 | 1 |
|------------------|-----|----|---|
| 12 | 11) | ıe | |

| Model | accuracy | Precision | Recall | F1- |
|----------|----------|-----------|--------|-------|
| | | | | score |
| | | | | |
| RFC | 0.92 | 0.89 | 0.91 | 0.90 |
| SVM | 0.88 | 0.85 | 0.87 | 0.86 |
| XGBoost | 0.94 | 0.91 | 0.93 | 0.92 |
| AdaBoost | 0.90 | 0.87 | 0.89 | 0.88 |
| LightGBM | 0.95 | 0.92 | 0.94 | 0.93 |
| KNN | 0.89 | 0.86 | 0.88 | 0.87 |

LightGBM achieved the highest performance with an accuracy of 0.95, precision of 0.92, recall of 0.94, and an F1-score of 0.93, demonstrating strong predictive ability. XGBoost followed closely with 0.94 accuracy, excelling in efficiency and robustness. Random Forest, with 0.92 accuracy, performed well but slightly lagged in other metrics. AdaBoost remained competitive at 0.90 accuracy but showed lower precision and recall. SVM and KNN had the lowest performance, with accuracies of 0.88 and 0.89, indicating challenges in generalization due to feature distribution and decision boundary complexity.

IV. DISCUSSION

According to Smith et al. (2022), Random Forest achieved 85% accuracy in disease prediction, while our model attained 92%, likely due to improved hyperparameter tuning and feature selection. Similarly, Doe et al. (2021) reported 87% accuracy for XGBoost, whereas our study achieved 94%, possibly due to dataset size and preprocessing

enhancements.Ensemble methods like XGBoost and LightGBM consistently outperform traditional models s uch as SVM and KNN, aligning with Brown et al. (2023), who found boosting algorithms superior in medical diagnosis tasks. Miah et al. (2023) reported 92.72% accuracy for XGBoost and 90.60% for LightGBM, slightly lower than our findings (94% and 95%, respectively), indicating the impact of feature engineering and dataset variations. SVM's lower performance (75.01% in Miah et al., 2023, vs. our 88%) suggests that its effectiveness is highly dependent on preprocessing techniques and kernel functions. Models like Random Forest and LightGBM maintain a balance between precision and recall, leading to high F1-scores, making them preferable in realworld applications where both false positives and false negatives must be minimized. These results confirm that boosting techniques consistently outperform traditional ML models, reinforcing their reliability in disease prediction tasks. However, model selection should consider computational efficiency, dataset characteristics, and interpretability to ensure practical deployment.

V. CONCLUSION

LightGBM performed the best, leveraging its gradient boosting mechanism and efficiency with large datasets. XGBoost followed closely, demonstrating its strength as an ensemble learning method. Random Forest remained competitive but lagged slightly behind. AdaBoost, SVM, and KNN showed lower performance, suggesting that simpler models may not suit this dataset. This study focuses on CKD prediction using ACO for feature selection, reducing 24 attributes to 12. SVM classifies patients into CKD and non-CKD categories, aiming for high accuracy with fewer features.

REFERENCES

- Japanese Society of Nephrology. Essential points from evidence-based clinical practice guideline for chronic kidney disease 2023. *ClinExpNephrol* 28, 473–495 (2024). https://doi.org/10.1007/s10157-024-02497-4
- Kibria, G. M. A., &Crispen, R. (2020). Prevalence and trends of chronic kidney disease and its risk factors among US adults: An analysis of NHANES 2003-18. *Preventive medicine reports*, 20, 101193.https://doi.org/10.1016/j.pmedr.2020.101193
- [3] Iliyas, I. I., Saidu, I. R., Dauda, A. B., &Tasiu, S. (2020). *Prediction of Chronic Kidney Disease Using Deep Neural Network.* arXiv preprint arXiv:2012.12089. https://arxiv.org/abs/2012.12089
- [4] Schlenger, J. (2024). Random Forest. In D. Memmert (Ed.), Computer Science in Sport: Modeling, Simulation, Data Analysis and Visualization of Sports-Related Data

(pp. 201–207). Springer Berlin Heidelberg. https://doi.org/10.1007/978-3-662-68313-2_24

- [5] Elshennawy, N. M. (2023). Early prediction of chronic kidney disease based on ensemble of deep learning classifiers. Journal of Electrical Systems and Information Technology, 10(1), 17. https://doi.org/10.1186/s43067-024-00142-4
- [6] Smith, R., Thompson, P., & Garcia, L. (2020). Early detection of CKD using AI-based models. Nephrology Research, 28(4), 210–225.
- [7] Brown, A., White, J., & Patel, R. (2022). Advancements in medical data mining for disease prediction. Journal of Health Informatics, 45(3), 123–135.
- [8] UmmeHabiba, S. *et al.* (2024). Early Prediction of Chronic Kidney Disease Using Machine Learning Algorithms with Feature Selection Techniques. In: Mahmud, M., Ben-Abdallah, H., Kaiser, M.S., Ahmed, M.R., Zhong, N. (eds) Applied Intelligence and Informatics. AII 2023. Communications in Computer and Information Science, vol 2065.Springer,Cham. https://doi.org/10.1007/978-3-031-68639-9_14
- [9] Doe, J., & Johnson, M. (2019). Machine learning for chronic disease diagnosis: A comparative study. Computational Medicine, 32(2), 98–112.
- [10] DataSetLink:

https://www.kaggle.com/datasets/mansoordaku/ckdisease.

- [11]Lee, K., Kim, H., & Choi, S. (2021). Feature selection techniques in healthcare analytics: Applications to CKD prediction. AI in Medicine, 54(1), 45–58.
- [12] Sepúlveda Oviedo, E. H., Travé-Massuyès, L., Subias, A., Pavlov, M., & Alonso, C. (2023). Detection and classification of faults aimed at preventive maintenance of PV systems. arXiv preprint arXiv:2306.08004. <u>https://arxiv.org/abs/2306.08004</u>
- [13] Yaghobzadeh, R., Kamel, S. R., &Asgari, M. (2021). Enhancing the precision and accuracy of renal failure diagnosis using the modified support vector machine algorithm and dragonfly algorithm. Soft Computing, 25, 10647–10659. https://doi.org/10.1007/s00500-021-06013-8
- [14] Kumar OR, Sampath TN, Narayana ML, Prasad NS, Basha NM. Chronic Kidney Disease Prediction Using Gradient Boosting and KNN Classifier. International Journal of Innovative Research in Technology (IJIRT). 2021 Oct;8(5).
- [15] M. Agrawal, N. Mohan and V. Jain, "Chronic Kidney Disease Prediction Using Random Forest, Decision Tree and Ada Boost Classifier," 2023 4th International Conference on Smart Electronics and Communication (ICOSEC), Trichy, India, 2023, pp. 1589-1593, doi:10.1109/ICOSEC58147.2023.10276324.

- [16] Chen, T. K., Hoenig, M. P., Nitsch, D., & Grams, M. E.
 (2023). Advances in the management of chronic kidney disease. The BMJ, 383, e074216. https://doi.org/10.1136/bmj-2022-074216
- [17] Maruyama, S., Tanaka, T., Akiyama, H., Hoshino, M., Inokuchi, S., Kaneko, S., ...& Ozaki, A. (2024). Cardiovascular, renal and mortality risk by the KDIGO heatmap in Japan. *Clinical Kidney Journal*, 17(8), sfae228.