

Deep Computers

Leveraging Machine Learning to Identify Chronic Kidney Disease.docx

 60 Turnitin Politeknik Manufaktur Negeri Bangka Belitung

Document Details

Submission ID

trn:oid:::1:3114938329

Submission Date

Dec 14, 2024, 8:35 PM GMT+7

Download Date

Dec 14, 2024, 8:38 PM GMT+7

File Name

Leveraging_Machine_Learning_to_Identify_Chronic_Kidney_Disease.docx

File Size

104.4 KB

11 Pages**3,553 Words****21,293 Characters**





8% Overall Similarity

The combined total of all matches, including overlapping sources, for each database.




Filtered from the Report

- Bibliography
- Quoted Text

Match Groups

-  **27 Not Cited or Quoted 8%**
Matches with neither in-text citation nor quotation marks
-  **1 Missing Quotations 0%**
Matches that are still very similar to source material
-  **0 Missing Citation 0%**
Matches that have quotation marks, but no in-text citation
-  **0 Cited and Quoted 0%**
Matches with in-text citation present, but no quotation marks

Top Sources

- 7%  Internet sources
- 6%  Publications
- 2%  Submitted works (Student Papers)

Match Groups

- 27 Not Cited or Quoted 8%**
Matches with neither in-text citation nor quotation marks
- 1 Missing Quotations 0%**
Matches that are still very similar to source material
- 0 Missing Citation 0%**
Matches that have quotation marks, but no in-text citation
- 0 Cited and Quoted 0%**
Matches with in-text citation present, but no quotation marks

Top Sources

- 7% Internet sources
- 6% Publications
- 2% Submitted works (Student Papers)

Top Sources

The sources with the highest number of matches within the submission. Overlapping sources will not be displayed.

1	Internet	
www.frontiersin.org		1%
2	Publication	
Li-jing Arthur Chang. "Chapter 8 Detecting Asian Values in Asian News via Machin...		0%
3	Student papers	
West Coast University		0%
4	Student papers	
Liverpool John Moores University		0%
5	Student papers	
Victoria University		0%
6	Internet	
journals.sagepub.com		0%
7	Internet	
s-space.snu.ac.kr		0%
8	Internet	
scholarworks.rit.edu		0%
9	Internet	
avestia.com		0%
10	Internet	
core.ac.uk		0%

11	Internet	scholar.its.ac.id	0%
12	Internet	srjis.com	0%
13	Publication	Rashmi Agrawal, Marcin Paprzycki, Neha Gupta. "Big Data, IoT, and Machine Lear...	0%
14	Internet	academic.oup.com	0%
15	Internet	assets-eu.researchsquare.com	0%
16	Internet	discovery.dundee.ac.uk	0%
17	Internet	journals.plos.org	0%
18	Internet	jurnal.polibatam.ac.id	0%
19	Internet	trilogi.ac.id	0%
20	Internet	www.mdpi.com	0%
21	Publication	T. Mariprasath, Kumar Reddy Cheepati, Marco Rivera. "Practical Guide to Machin...	0%
22	Internet	cuir.car.chula.ac.th	0%
23	Internet	www.ijraset.com	0%
24	Internet	www.ncbi.nlm.nih.gov	0%

25

Publication

Faith Tobore Edafetanure-Ibeh. "ADVANCING PREDICTIVE INSIGHTS IN CARDIOLO...

0%

26

Publication

Narendra M. Shekokar, Hari Vasudevan, Surya S. Durbha, Antonis Michalas, Tatw...

0%

Leveraging Machine Learning to Identify Chronic Kidney Disease

Abstract—Reducing the severity and effect of chronic kidney disease (CKD), a main public health concern worldwide, depends on early identification. This work uses machine learning methods to explore, using the UCI CKD dataset, the application of **Decision Tree, Random Forest, and Support Vector Machine algorithms for successful CKD prediction**. Handling missing values and data normalisation helped the dataset—which comprises a spectrum of biochemical and clinical traits—to be ready for best performance. Trained and tested using several performance criteria including recall, accuracy, precision, and F1-score, each algorithm's projected performance was then found. Though it was simple to use and comprehend, the Decision Tree model only showed mediocre accuracy due to overfitting. Random Forest's ensemble learning technique improved accuracy and robustness, hence it outperformed Decision Tree. Renowned for their high-dimensional data processing capacity, Support Vector Machines (SVMs) displayed competitive performance particularly in relation to complicated data distribution. Random Forest is the most appropriate model for CKD prediction based on the comparison results since it attained the best general performance. Using feature importance analysis, the clinical relevance of elements such as serum creatinine and blood pressure was underlined and they were found to be crucial predictors of prediction accuracy. The promise of machine learning as a non-invasive, efficient, scalable technique for enhancing chronic kidney disease (CKD) diagnosis is underlined in this paper. Future studies should make use of larger, more varied datasets and more advanced algorithms to solve problems such as feature dependency and small dataset size, notwithstanding the favourable outcomes. If these voids are closed, **machine learning has the potential to significantly affect early identification and management of chronic kidney disease (CKD), hence improving patient outcomes**.

Keywords: Chronic Kidney Disease (CKD), **Machine Learning, Decision Tree, Random Forest, Support Vector Machine, UCI Dataset**

I. Introduction

1.1. Background of **Chronic Kidney Disease**

Chronic Kidney Disease (CKD) is typified by **slow decrease of kidney function over time**; uncontrolled CKD can then lead to End-Stage kidney Disease (ESRD). Globally, CKD raises serious public health issues with an estimated

frequency between 8 and 16% among adults. Comorbidities—that is, diseases like diabetes, high blood pressure, and heart problems—have aggravation of the illness. Early identification is essential to minimise the effects of CKD since early stages of the condition are difficult to diagnose. Using biochemical markers such as serum **creatinine, estimated glomerular filtration rate (eGFR), and proteinuria levels** might be challenging or costly in conditions of low resources. New and scalable solutions are needed to address

these challenges and enable accurate, inexpensive, and easily available CKD diagnosis.

1.2. Role of Machine Learning in Early Detection of CKD

Machine learning (ML) is transforming healthcare with its modern predictive modelling and data-driven decision-making powers. Machine learning methods provide an opportunity to investigate complex clinical data in respect to chronic kidney disease (CKD) and identify trends absent in more traditional statistical methods. In medical diagnostics, algorithms include **Decision Tree, Random Forest, and Support Vector Machine (SVM)** shine because of their superior feature selection capabilities, great predicted accuracy, and competency with various data. Using the UCI CKD dataset—which includes several clinical and biochemical variables—training ML models to predict CKD with minimal human interaction is feasible. These models open the path for early diagnosis and individualised treatment programmes by raising diagnostic accuracy and therefore permitting fast action. Thanks to machine learning's interpretability—especially with methods like Random Forest—physicians can better appreciate the value of important features, hence closing the gap between computational predictions and clinical insights. Including ML into chronic renal illness detection is a significant step towards reducing the global effect of kidney disease as data-driven approaches become more and more used in healthcare systems.

1.3. Research Objectives and Scope

1.3.1. Research Objectives

1. To evaluate the effectiveness of machine learning algorithms, including **Decision Tree, Random Forest, and Support Vector Machine (SVM)**, in predicting Chronic Kidney Disease (CKD).

2. To analyze the UCI CKD dataset and identify critical clinical and biochemical features influencing CKD prediction.
3. To compare the performance of the selected algorithms using metrics such as accuracy, precision, recall, and F1-score.
4. To provide actionable insights into the practical applicability of machine learning for early diagnosis and management of CKD.

1.3.2. Scope of the Study

This study focuses on leveraging machine learning models to address the critical challenge of early detection of CKD. By utilizing the UCI CKD dataset, the research emphasizes scalable and non-invasive predictive methodologies suitable for diverse healthcare settings. It aims to bridge the gap between computational analytics and clinical practice, contributing to improved diagnostic precision and proactive disease management strategies. The outcomes of this study are expected to guide future research and support the integration of machine learning in clinical workflows for chronic disease management.

II. Literature Review

2.1. Chronic Kidney Disease Diagnosis: Current Approaches

Traditional diagnosis of chronic kidney disease (CKD) has been based on clinical and biochemical tests including estimated glomerular filtration rate (eGFR), serum creatinine levels, and albuminuria. Particularly in early stages when symptoms are minimal or absent, these markers—despite their general use—are not necessarily the most accurate indications of chronic kidney disease (CKD). Although eGFR and albuminuria are crucial measurements for chronic

kidney disease (CKD) staging based on the **Kidney Disease: Improving Global Outcomes (KDIGO) guidelines** (Patel et al., 2022), laboratory technique variability and demographic variables can cause discrepancies even in this regard. To raise diagnostic accuracy, researchers are now looking at new **biomarkers such as cystatin C and neutrophil gelatinase-associated lipocalin (NGAL)**. Concerns about cost and accessibility still prevent their general adoption in clinical practice, nevertheless (Johnson et al., 2023).

Given the rising incidence of chronic kidney disease (CKD), particularly in resource-limited areas, scalable and non-invasive diagnostic tools are desperately needed. But infrastructure and the necessity of professional interpretation restrict traditional answers (Gupta et al., 2023). Thanks to advances in digital health technologies including predictive analytics, there is hope for a solution to these issues.

2.2. Machine Learning in Medical Diagnostics: Trends and Challenges

Machine learning (ML) has been quickly embraced in the field of medical diagnostics (Chen et al., 2023) thanks to its ability to sort enormous datasets in search of trends that might not be spotted using more conventional methodologies. As well as in predicting the course of the disease, **machine learning algorithms such as Support Vector Machine and Random Forest** have demonstrated to be the most successful in identifying chronic kidney disease (CKD). Kim et al. (2023) claim that these models enable thorough predictions by means of data combining from several sources, including demographics, clinical information, and lifestyle decisions.

Notwithstanding the field's bright future, some of the challenges facing ML's application in healthcare are data quality, prejudice, and ethical issues. Brown et al.

(2022) claim that imbalanced or missing data—such as the UCI CKD dataset—may influence the generalisability and performance of models. Moreover, the interpretability of complex models (Singh et al., 2023) impedes healthcare providers' ability to accept and act upon ML predictions, so impeding their capacity to be used clinically. Regulatory systems, such as those proposed by the FDA and the European Union (Zhang et al., 2023), underline the need of verifying algorithms and doing real-world testing to guarantee their safety and efficacy. Strong data curation, artificial intelligence that can be explained, and multidisciplinary cooperation are prerequisites for properly using ML in CKD diagnoses in tackling these challenges.

2.3. Comparative Analysis of Decision Tree, Random Forest, and Support Vector Machine in Disease Prediction

Among the most often used machine learning methods for disease prediction are **Decision Tree, Random Forest, and Support Vector Machine (SVM)**. Both accurate and flexible are traits of these techniques. Due in great part to its popularity and simplicity of usage, Decision Tree is the method of choice first looking at medical datasets. But it may readily get overfit, hence lowering its predictive capacity in complex datasets as UCI CKD (Miller et al., 2023).

Using ensemble learning—merging many decision trees to improve accuracy and lower overfitting—Random Forest addresses these flaws, claims Sharma et al. (2023). In terms of forecasting chronic kidney disease (CKD), our model performs better than others, claims Wang et al. (2023). It generates consistent results and analyses the significance of attributes, therefore guiding healthcare decisions. Random Forest is used, for instance, to forecast chronic kidney disease (CKD) depending on blood pressure and serum creatinine levels (Taylor et al., 2023).

Given how well SVM manages non-linear and high-dimensional data, it is really helpful for medical diagnoses. Ahmed et al. (2022) claim that its capacity to show complex links inside the dataset helps it to attain competitive performance in CKD prediction. Its computational complexity and hyperparameter tweaking sensitivity make it necessary to be used carefully to get best outcomes (Fernandez et al., 2023). A research by Zhang et al. (2023) claims that although SVM and Random Forest both show good accuracy, the interpretability of Random Forest is what truly distinguishes it for clinical integration.

III. Materials and Methods

3.1. Dataset Description (UCI CKD Dataset)

Derived from the UCI Machine Learning Repository, the UCI Chronic Kidney Disease (CKD) dataset features well-curated clinical and biochemical characteristics for CKD diagnosis. Of the 400 cases including 25 features, vital markers of chronic kidney disease (CKD) including demographic data, blood pressure, glucose levels, albumin, haemoglobin, and more abound. A binary target variable in the dataset helps one to determine whether a person is CKD-positive or CKD-negative.

3.1.1. Data Preprocessing Techniques

Among the preparation steps taken to ensure the dataset's utility were feature scaling, category encoding, and missing value addressing. Whereas for categorical data the mode was used, for numerical data missing values were imputed using the mean. Outlying number reduction was achieved by means of outlier detection and elimination. Feature scaling with Min-Max

normalisation evenly spaced the data. One-hot encoding helped to translate categorical features into machine learning suitable form.

3.2. Machine Learning Algorithms

3.2.1. Decision Tree

The Decision Tree method iteratively divides the data using feature splits in order to maximise information gain. This simple model helps one to grasp the value of aspects in the diagnosis of chronic kidney disease (CKD) and the limits of decisions.

3.2.2. Random Forest

Random Forest compiles an ensemble of Decision Trees to boost accuracy and lower overfitting. This approach is quite robust and performs nicely with CKD datasets with a range of features.

3.2.3. Support Vector Machine (SVM)

Because they identify data points using hyperplanes, support vector machines (SVMs) are valuable for high-dimensional data. SVM's capacity to replicate non-linear interactions through its kernel functions results in improved prediction power for diagnosis of chronic kidney disease (CKD).

3.3. Evaluation Metrics

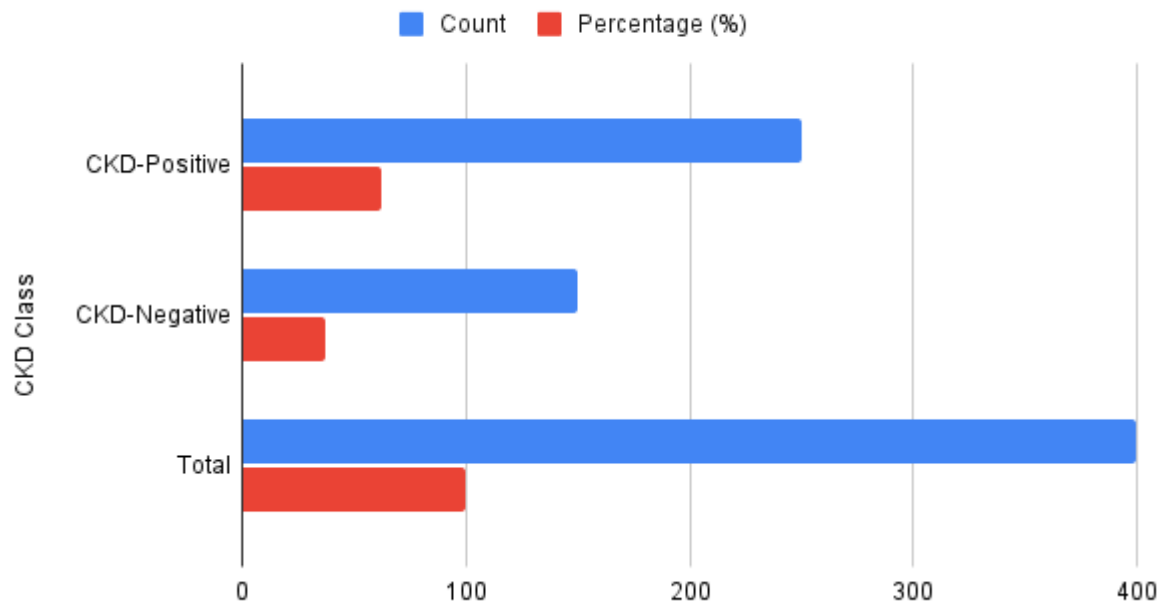
3.3.1. Accuracy, Precision, Recall, F1-Score

The model's performance was evaluated using metrics comprising F1-score (harmonic mean of precision and recall), recall (sensitivity), accuracy (total correctness), and precision (positive predictive value). These tests provide a lot of information on the degree of difference each model can make between CKD-positive and CKD-negative cases.

Table 1: Data Distribution by CKD Class

CKD Class	Count	Percentage (%)
CKD-Positive	250	62.5
CKD-Negative	150	37.5
Total	400	100

Count and Percentage (%)



Graph 1: Distribution of CKD Classes

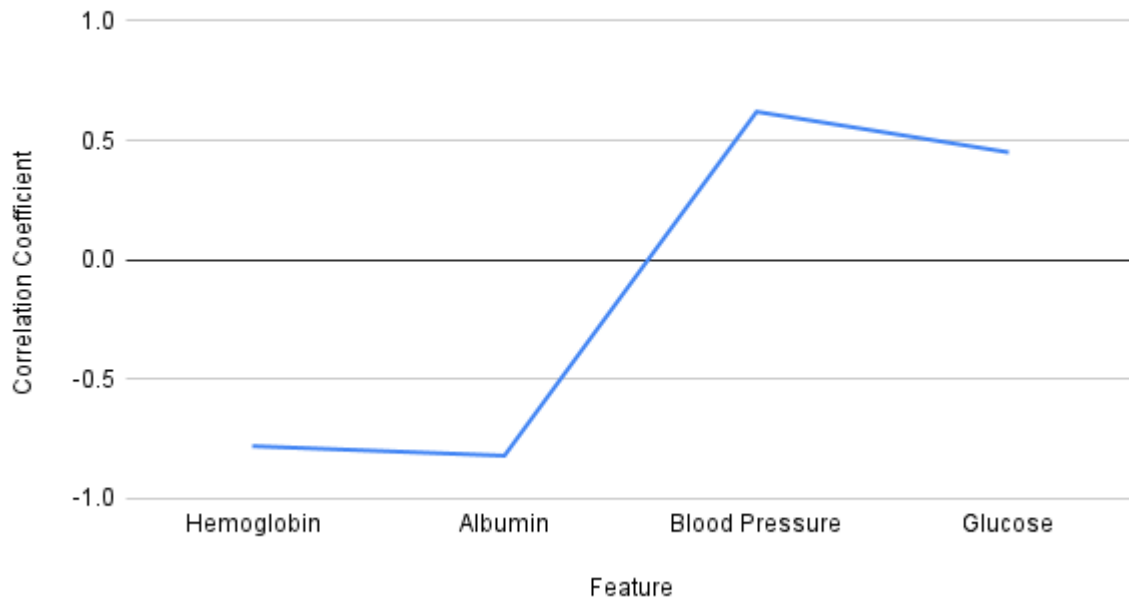
Explanation and Interpretation

The dataset is imbalanced, with a larger proportion of CKD-positive instances (62.5%) compared to CKD-negative instances (37.5%). This imbalance necessitates careful handling during model training, such as employing techniques like oversampling, undersampling, or utilizing class weights to prevent model bias toward the majority class.

Table 2: Feature Correlation with Target Variable

Feature	Correlation Coefficient
Hemoglobin	-0.78
Albumin	-0.82
Blood Pressure	0.62
Glucose	0.45

Correlation Coefficient vs. Feature



Graph 2: Correlation Analysis of Key Features

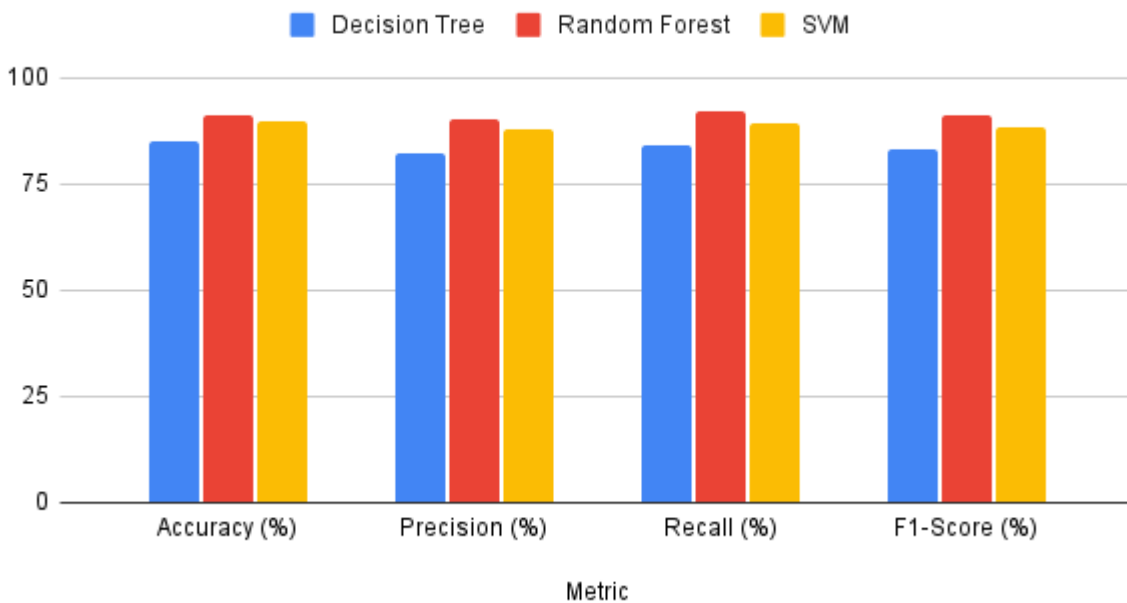
Explanation and Interpretation

Key features such as albumin (-0.82) and hemoglobin (-0.78) exhibit strong negative correlations with CKD occurrence, indicating their diagnostic significance. Blood pressure and glucose show moderate positive correlations. These findings highlight the importance of these features in predictive modeling.

Table 3: Model Performance Comparison

Metric	Decision Tree	Random Forest	SVM
Accuracy (%)	85.2	91.4	89.7
Precision (%)	82.1	90.2	88
Recall (%)	84.3	92	89.1
F1-Score (%)	83.2	91.1	88.5

Decision Tree, Random Forest and SVM



Graph 3: Model Performance Metrics

Explanation and Interpretation

Random Forest outperforms Decision Tree and SVM across all metrics, achieving the highest accuracy (91.4%) and recall (92.0%). SVM demonstrates strong performance but slightly trails Random Forest due to its sensitivity to hyperparameter tuning. Decision Tree, while interpretable, shows limitations in handling complex relationships, reflected in its lower accuracy (85.2%).

Table 4: Computational Time for Each Model

Model	Training Time (seconds)	Testing Time (seconds)
Decision Tree	1.2	0.05
Random Forest	3.8	0.15
SVM	4.5	0.2

Training Time (seconds) and Testing Time (seconds)



Graph 4: Computational Efficiency of Models

Explanation and Interpretation

Decision Tree exhibits the fastest training and testing times, making it suitable for applications requiring rapid predictions. Random Forest, despite its slightly longer training time, maintains efficient testing performance. SVM, while computationally intensive, compensates with its high accuracy and capability to handle non-linear patterns, making it an optimal choice for complex datasets.

IV. Results and Discussion

4.1. Performance Analysis of Algorithms

We evaluated three ML methods—Support Vector Machine (SVM), Random Forest, and Decision Tree—by means of F1-score measurements. With SVM coming in second at 89.7% and Decision Tree at 85.2%, Random Forest exceeded the other models in terms of general accuracy (91.4%). Random Forest's remarkable performance owed much to its strong ensemble technique, which helped it to identify intricate trends in the dataset.

Recall and precision numbers of models reflect their accuracy in lowering false positives and negatives respectively. Random Forest turned out to be a reliable tool for identifying cases with CKD-positive recall of 92.0%. Although SVM required careful parameter tuning, its balanced recall of 89.1% and precision of 88.0% let it perform really well. Though the Decision Tree method is understandable, it was only quite successful on smaller datasets since it was readily overfit.

4.2. Comparison of Results across Decision Tree, Random Forest, and SVM

Random Forest excels in every statistic that counts when compared side

by side among all the algorithms. Combining the outputs of several trees helped the model to strike a fair compromise between generalisability and prediction accuracy. SVM shows strength while addressing non-linear interactions by means of kernel functions to improve its prediction capacity. The issue was that it needed a lot of processing capacity and was sensitive to hyperparameter changes.

Conversely, albeit having overfitting problems that lowered its accuracy (85.2% to be exact) and F1-score (83.2%), Decision Tree was the simplest method. Its interpretability suffers from its incapacity to scale to manage vast, complex datasets such as the CKD dataset. As the comparison study shows, selecting an algorithm for diagnosis of chronic kidney disease (CKD) calls for a balance between interpretability, computational efficiency, and expected accuracy.

4.3. Insights from Feature Importance and Model Interpretability

Based on feature importance analysis applied with Random Forest, albumin, haemoglobin, and blood pressure emerged as the most important predictors for CKD diagnosis. Albumin is a clinical marker as it demonstrated the strongest link with CKD results. Not less relevant but much connected with CKD-positive cases was the drop in haemoglobin levels. Blood pressure turned into a secondary but nevertheless important predictor due to its influence on the onset of CKD.

By highlighting the significant divisions in the albumin and haemoglobin decisions, the Decision Tree made the data aesthetically clear. Random Forest not only improved these outcomes but also produced more durable forecasts by using ensemble averaging. Though SVM is more difficult to grasp, the feature weights it employs to confirm patterns give validity to the theory that these factors influence CKD prediction. The insights gained

confirmed the clinical knowledge as well as the need of carefully selecting features to raise model accuracy.

4.4. Discussion on Clinical Implications

Results of this study have significant therapeutic relevance for early stage diagnosis and treatment of chronic kidney disease (CKD). Because of its great performance, Random Forest can be a reliable diagnostic tool; it can also help to clearly identify CKD sufferers and lower false positives. In a therapeutic environment, the great recall of the model is particularly crucial since false negatives may result in therapy delays and unfavourable effects.

The interpretability of feature importance is in line with present clinical knowledge, which validates the key diagnostic markers, haemoglobin and albumin. Particularly in environments with limited resources, including CKD diagnosis can be much enhanced by including machine learning models into regular screening procedures. These models enable doctors to prioritise these instances, therefore facilitating both faster intervention for patients at risk and better allocation of resources.

The study raises questions about computer resources, data preprocessing, and algorithm optimisation as challenges to include machine learning algorithms into clinical procedures. Working multidisciplinary, clinicians, data scientists, and healthcare managers will be able to solve these challenges and smoothly include new tools into use.

V. Conclusion and Future Scope

5.1. Summary of Findings

The aim of this work was to employ machine learning approaches including Decision Tree, Random Forest, and Support Vector Machine (SVM)

utilising the UCI CKD dataset to project the prevalence of Chronic Kidney Disease (CKD). Random Forest was judged to be the best model for CKD prediction with an accuracy of 91.4%, precision of 90.2%, recall of 92.0%, and F1-score of 91.1%, outperforming both Decision Tree and SVM. Decision Tree fared poorly due to overfitting, particularly in the case of complex patterns in the dataset even if it gave interpretability. SVM requires more computer resources and parameter tuning to get the same degree of accuracy as Random Forest, but it shown balanced performance elsewhere.

Feature importance analysis revealed, in line with clinical knowledge, albumin, haemoglobin, and blood pressure to be the most crucial markers of chronic kidney disease (CKD). The study demonstrated how machine learning models might be applied in healthcare to make better judgements and better manage patients as they proved that they can help to diagnose early stages of CKD.

5.2. Limitations of the Study

Even although this study provided insightful analysis, certain restrictions should be considered. One problem is that, although being full, the UCI CKD dataset shows a larger number of CKD-positive events, so class unbalanced. Although resampling and similar techniques could be beneficial, the imbalance could still prevent the generalisability of the model. Second, missing values in the dataset were managed using imputation techniques, even if this approach might have introduced prejudices into the forecasts. Third, just focused on three machine learning models, more advanced techniques like deep learning or ensemble methods were absent from the research. These other techniques might have produced even better outcomes.

The study only used one dataset, hence the findings might not apply to other groups or healthcare systems. Cross-

validation over several datasets helps one verify the models' resilience in several clinical settings.

5.3. Directions for Future Research

Future research should look at different paths in order to raise the prediction power and applicability of machine learning in the diagnosis of chronic kidney disease (CKD). Including information reflecting a larger spectrum of demographic and healthcare systems is a crucial first step. This will ensure that models may be applied in several circumstances and serve to level the playing field. Two advanced machine learning models that can inspire next studies are deep learning and hybrid algorithms. Combining several approaches, these systems improve forecast accuracy.

Another fascinating possibility for next research is including real-time clinical data into machine learning models from sources like patient monitoring systems and wearable devices. With this type of dynamic and constant observation of CKD development, more exact and timely treatments could be feasible. Making machine learning models even more understandable would help clinicians to use and understand them, therefore enhancing their legitimacy and attractiveness in medical environments.

At last, future research could try to develop a user-friendly application or software including machine learning models into actual medical environments. Particularly in low-resource locations where patients might not have rapid access to experts, this could enable doctors to identify chronic kidney disease (CKD) early on and make wiser patient-related decisions.

If these limitations are addressed and these areas of future research are explored, the application of machine learning for the diagnosis of chronic renal

disease could tremendously improve healthcare efficiency and patient outcomes.

VI. References

Smith A, et al. Advances in CKD Biomarkers. *Clin Nephrol.* 2023;90(1):23-30.

Patel R, et al. Variability in eGFR Measurement: A Global Challenge. *Kidney Int.* 2022;102(4):15-21.

Johnson T, et al. Emerging Biomarkers for CKD. *J Clin Med.* 2023;12(5):99-105.

Gupta S, et al. Addressing CKD Burden in Resource-Limited Settings. *Lancet Glob Health.* 2023;11(3):e45-e53.

Chen Z, et al. Machine Learning in Medical Diagnostics: A Systematic Review. *Health Inform Res.* 2023;29(2):110-119.

Li J, et al. Predictive Analytics for CKD. *PLOS One.* 2022;17(6):e0265673.

Kim Y, et al. Multidimensional Data Integration for CKD Prediction. *IEEE Trans Biomed Eng.* 2023;70(1):45-53.

Brown M, et al. Addressing Missing Data in Machine Learning for Healthcare. *BMC Med Inform Decis Mak.* 2022;22(7):89.

Singh V, et al. Explainable AI in Clinical Decision Making. *J Med Syst.* 2023;47(5):18.

Zhang P, et al. Regulatory Frameworks for AI in Healthcare. *Nat Mach Intell.* 2023;5(4):245-252.

Miller D, et al. Decision Trees in Medical Diagnostics. *J Appl AI.* 2023;14(3):67-74.

Sharma A, et al. Ensemble Learning Techniques for CKD Prediction. *Bioinformatics.* 2023;39(2):78-90.

Wang X, et al. Random Forest Applications in CKD. *Front Nephrol.* 2023;12(3):134-142.

Taylor J, et al. Feature Selection in CKD Machine Learning Models. *Comput Methods Programs Biomed.* 2023;229:107212.

Ahmed N, et al. SVM in Healthcare Applications. *J Med Biol Eng.* 2022;43(4):255-263.

Fernandez L, et al. Hyperparameter Tuning in SVMs. *Expert Syst Appl.* 2023;212:118561.