# Chronic Kidney Disease Prediction Using Machine Learning and Deep Learning Models

1ˢᵗ Shaik Rafi
*Asst.Prof, Dept of CSE,*
*Narasaraopeta Engineering College,*
Narasaraopet-522601, Palnadu,
Andhra Pradesh, India
shaikrafinrt@gmail.com

2ʳᵈ Nuti Revanth
*Dept of CSE,*
*Narasaraopeta Engineering College,*
Narasaraopet-522601, Palnadu,
Andhra Pradesh, India.
nutirevanth541@gmail.com

3ᵗʰ K Veera Raghava Reddy
*Dept of CSE,*
*Narasaraopeta Engineering College,*
Narasaraopet-522601, Palnadu,
Andhra Pradesh, India.
raghavareddykota50@gmail.com

4ᵗʰ K Mahesh Babu
*Dept of CSE,*
*Narasaraopeta Engineering College,*
Narasaraopet-522601, Palnadu,
Andhra Pradesh, India.
kurramaheshbabu144@gmail.com

5ᵗʰ Y Likhith Prasanna Kumar
*Dept of CSE,*
*Narasaraopeta Engineering College,*
Narasaraopet-522601, Palnadu,
Andhra Pradesh, India.
likithprasannakumar@gmail.com

6ⁿᵈ N vijay kumar
*Asst.Prof, Dept of CSE,*
*Narasaraopeta Engineering College,*
Narasaraopet-522601, Palnadu,
Andhra Pradesh, India
nvk20022001@gmail.com

*Abstract*—**Chronic kidney disease is a noticeable health condition that can persist throughout an individual's life, resulting from either kidney malignancy or diminished kidney function. In this work, we investigate how several machine learning techniques might provide an early CKD diagnosis. While previous research has extensively explored this area, our aim is to refine our approach by employing predictive modeling techniques. Initially, we considered 25 variables alongside the class property. The data set used in this study underwent extensive processing, including changing the names of colours for clarity, converting identified colours to numbers, treating unique values with letters handling of partitioned values, fixing incorrect values, filling null values with mean, and encoding categorical values into mathematical notation. In addition,Principal component analysis (PCA) was also employed to lower dimensionality. Our findings demonstrated that the XG Boost classifier surpassed every other algorithm, with an accuracy of 0.991.**

*Index Terms*—**Decision tree, Random Forest, Chronic Kidney Disease, Logistic Regression, and XG BOOST classifier.**

## I. INTRODUCTION

The progressive and irreversible degeneration of renal cells is the hallmark of the sickness known as chronic kidney disease [1]. When CKD develops, harmful wastes accumulate in the body, leading to various health complications. Chronic kidney disease (CKD) can cause a range of symptoms including cough, fever, dry mouth, nausea, back pain and abdominal pain.CKD is often associated with two risk factors: diabetes and hypertension.Therefore early diagnosis and treatment are essential. There are encouraging opportunities to improve early detection of CKD through machine learning and predictive modelling[2].The research guarantees that advanced preprocessing methods will be employed to predict CKD utilizing machine learning algorithms.Previous manipulation of The data collected in this research [7] was extensive and included changing column names to improve readability, correcting incorrect assumptions, averaging a missing values will be replaced, and categorical variables will be assigned numerical labels. Several models such as XG Boost, AdaBoost, SVM,Decision Tree, Chi-Square, and KNN [7] are developed and evaluated.

### A. Stages of chronic kidney disease

#### a) Early stages of chronic kidney disease

Chronic kidney disease is often completely asymptomatic in its early stages[1]. This is because with a significant decrease in kidney function, the body can become more adaptive. CKD is often diagnosed during routine screening for other medical conditions, such as blood or urine tests[3]. Early detection is important as this allows for drug treatment through regular testing and ongoing monitoring.[2]

*b) CKD in Its Advanced Stages*

Many symptoms may appear if CKD is not identified early or if it gets worse despite treatment. Compared to when the kidney no longer functions, known as esrd or eskd, there is no possibility of survival without either dialysis or a kidney transplant.

*c) Time to see a doctor*

If you have indications of kidney disease, you should see a doctor immediately. Early detection of CKD can avoid kidney failure.[2]During medical testing, doctors may use blood and urine tests to measure kidney functioning and blood pressure, especially if you suffer from illnesses that increase your risk of developing renal disease. Discuss with your physician the significance of this test.

*d) Tests for CKD*

A illness or other ailment that damages the kidneys gradually leads to chronic kidney disease (CKD). Research reveals a 6.23% annual rise in CKD hospital admissions despite a steady global death rate.[7] Numerous diagnostic techniques are used to determine the status of chronic kidney disease such as eGFR, urine tests, and a blood pressure tests.[3] For diagnosis of kidney injury or structural abnormalities, further testing such as MRI scans, ultrasound, or CT scans may be required.

## II. RELATED WORK

Md. Ariful Islam (2023).[7] investigated CKD prediction utilizing a mix of machine learning methods, emphasizing the importance of model selection and data preprocessing. The study found that XGBoost performed well, and recommended future work to explore additional features and advanced methods for improved accuracy. Ammirati (2020).[1] discusses chronic kidney disease as a prevalent, progressive condition with high cardiovascular risks. The paper covers conservative treatments to slow progression and replacement therapies like dialysis. The Aljaaf et al(2018).[2] have applied machine learning and predictive analytics to improve early diagnosis of chronic kidney disease, using evolutionary computation techniques to optimize predictive models. Their work aims to enhance patient outcomes through timely interventions. The Hossain et al(2022).[5]have explored feature optimization techniques to improve the performance of machine learning models for diagnosing Chronic renal disease,focusing on improving accuracy and efficiency through optimized feature selection. P. Chittora et al.[3] (2021) studied CKD prediction using machine learning, with the Deep Neural Network (DNN) achieving 99.6% accuracy. They recommend enhancing

the DNN model's interpretability and conducting clinical trials to test its practical applicability. Ashiqul Islam et al.[6] (2020) used Random Forest for CKD risk factor prediction, achieving 97.8% accuracy. The study suggests improving prediction accuracy and exploring advanced machine learning techniques in future work. Bhavya Gudeti and Terrance Li.[4] (2020) introduced a CKD prediction approach using Support Vector Machines, achieving notable accuracy. They recommend applying the model in various clinical scenarios and integrating it with other diagnostic tools for improved patient outcomes. K. Venkatrao.[11] (2023) proposed the HDLNET model, an integrated deep learning model for CKD diagnosis. The model demonstrated impressive accuracy, showcasing the potential of combining different neural network architectures. Future work should aim at further refining the model and exploring additional datasets to enhance the generalizability of HDLNET S. M. M. Elkholy et al.[12] (2021) applied a Deep Belief Network (DBN) for early CKD prediction, achieving an accuracy of 98.5%. The study proposes future work to focus on applying the model to different population datasets and integrating it into clinical decision-making processes.

## III. METHODOLOGY

The CKD data set [10] was analyzed using a number of machine learning algorithms, including the DT, RF, XGBoost, AdaBoost, SVM, chi-square and KNN. [7] The objective of the machine learning model was to achieve excellent classification performance [4] with many features that improve PCA performance simulations.

### A. Data Preprocessing

The CKD dataset comprised 24 features and 1 target variable, with a mix of numerical and nominal attributes.[10] Initially, the dataset contained numerous missing values.The mean approach was utilized to estimate erroneous numerical values for continuous parameters, whereas the mode approach was applied to nominal values. This step ensured that the dataset was complete and ready for analysis.

To deal with missing results, a k-Nearest Neighbors (KNN)-based method [11] was used .The CKD dataset used in this study consisted of 400 cases with 24 items each. [10] no' (not CKD), in that order. Fig. 1 displays the Value Count for features in the dataset without validation[7], while Fig. 2 displays the Value Count for features with using PCA. Fig.5 displays the Target group

TABLE I: Overview of Dataset Features

| Feature | Details | Type / Values |
|---------|---------|---------------|
| age | Represents the age of the patient. | Numerical: in years |
| bp | Measures the patient's blood pressure. | Numerical: mm/Hg |
| sg | Ratio indicating urine density. | Nominal: 1.005, 1.010, etc. |
| al | Level of albumin detected in blood. | Nominal: 0, 1, 2, 3, 4, 5 |
| su | Patient's sugar concentration. | Nominal: 0, 1, 2, 3, 4, 5 |
| rbc | Counts red blood cells in patient. | Nominal: normal, abnormal |
| pc | Indicates pus cells present. | Nominal: normal, abnormal |
| pcc | Presence of clumps formed by pus cells. | Nominal: present, absent |
| ba | Identifies bacterial presence in samples. | Nominal: present, absent |
| bgr | Records random blood glucose levels. | Numerical: mg/dl |
| bu | Captures blood urea quantity. | Numerical: mg/dl |
| sc | Serum creatinine measured in blood. | Numerical: mg/dl |
| sod | Sodium levels found in the blood. | Numerical: mEq/L |
| pot | Potassium concentration in blood. | Numerical: mEq/L |
| hemo | Amount of hemoglobin available in blood. | Numerical: gms |
| pcv | Volume occupied by packed cells. | Numerical |
| wc | Count of white blood cells. | Numerical: cells/cumm |
| rc | Number of red blood cells. | Numerical: million/cumm |
| htn | Indicates hypertension condition. | Nominal: yes, no |
| dm | Records diabetes condition. | Nominal: yes, no |
| cad | Tracks presence of coronary artery disease. | Nominal: yes, no |
| appet | Evaluates the patient's appetite. | Nominal: good, poor |
| pe | Indicates swelling in the patient's lower extremities. | Nominal: yes, no |
| ane | Confirms anemia status in patient. | Nominal: yes, no |
| class | Classifies kidney disease condition. | Nominal: CKD, not CKD |

Fig. 1: **Value Count for features**



Fig. 2: **Value Count for features with PCA**



distribution shows that 250 patients do not have chronic kidney disease (CKD).Fig.4 displays the heat map which showes a significant correlation between multiple factors and squared scores.[7]

### B. Classifiers

In this project, the machine learning models used to diagnose chronic kidney disease (CKD) are trained and tested using classification methods.[4] The training dataset, which contains a variety of characteristics and their accompanying labels (CKD or non-CKD), is where these classifiers pick up patterns.The test set is used to confirm the models' efficacy in generating correct results by applying the knowledge learned during training to identify the existence or absence of CKD.
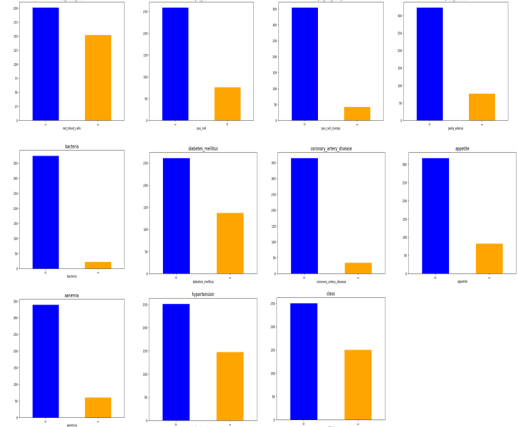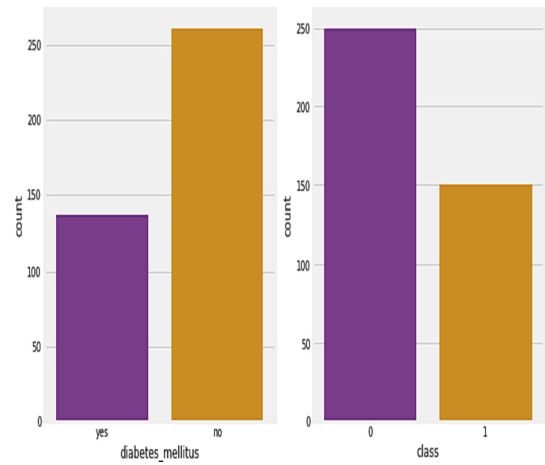
### a) AdaBoost

This study uses AdaBoost as a powerful machine learning approach to combine and improve the performance of poorly optimized classifiers, thereby improving CKD diagnostic accuracy.

### b) Decision Tree

The decision tree considers the values of various features Equal groups are obtained by subgrouping the data, which helps to correctly classify patients as CKD-positive or CKD-negative.

### c) XG Boost

This study uses the gradient-boosting method used by XG Boost, which computes the difference between observations and observations distance between Continuously adding new trees to predict errors or remnants of old trees to improve accuracy predictability through error learning.

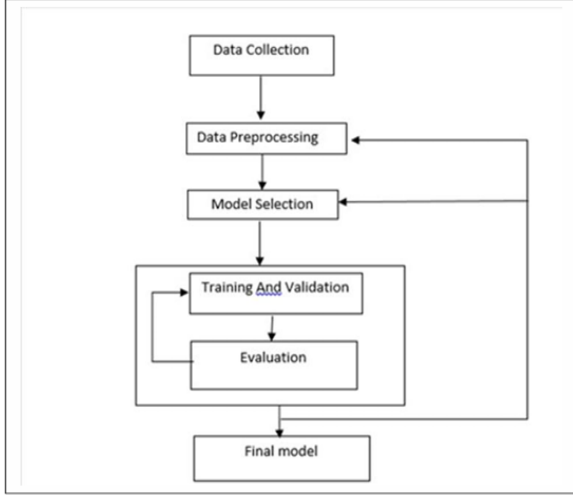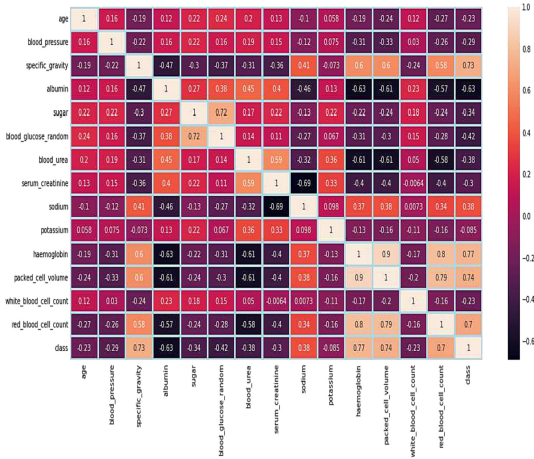Fig. 3: **The steps involved in the Model**



Fig. 4: **A heatmap illustrating the data & correlation pattern**



#### d) K-Nearest Neighbour (KNN)

KNN measures the Euclidean distance between each dataset instance and the query instance.This assures that the algorithm performs as efficiently as possible when patients are classified as CKD-positive or CKD-negative.

#### e) Random Forest

Many decision trees are produced by Random Forest using various subsets of a given CKD dataset. Random forest collects predictions from each individual tree and takes the average to produce the final result to improve the prediction accuracy .

#### f) Support Vector Machine

The issues of regression and classification can be efficiently handled by SVM. SVM works well for situations involving both regression and classification.

#### g) Chi-Square Test

Significant characteristics for model development are obtained using Chi-Square test to test the degree of independence between categorical variables and target variable.

#### h) Logistic Regression

This facilitates the identification of characteristics most important for predicting chronic kidney disease (CKD),especially when standard methods such as L1 (Lasso) or L2 (Ridge) are used,normalization is used.
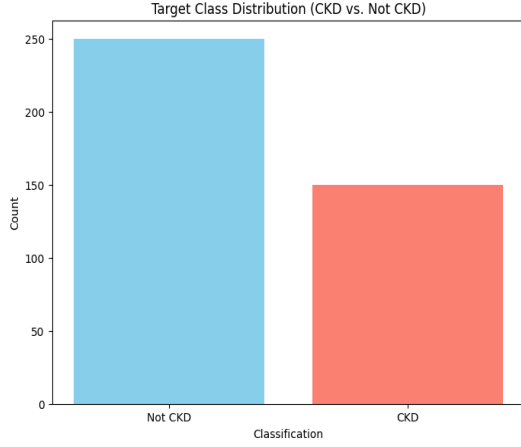
## IV. EXPERIMENTAL DATA

*CKD Dataset*

This project uses Chronic Kidney Disease data from the Kaggle datasets.[10] The dataset is divided into two groups: CKD (Yes) and non-CKD (No), with 24 items and objective variable 1. The CKD dataset contains 25 items, of which 11 are numeric and 14 are nominal items. The total list is 400 articles; 250 have been diagnosed with CKD and the remaining 150 without. These symptoms contribute to the prognosis of CKD by providing an accurate picture of the patient's health status.

---

**Algorithm 1** Principal Component Analysis (PCA)

---

1: **Input:** Data matrix $X \in \mathbb{R}^{n \times m}$, with $n$ samples and $m$ features.
2: **Output:** Principal components of $X$
3: **Step 1:** Standardize the dataset.
   - Center each feature by deducting the mean and scaling by the standard deviation.
4: **Step 2:** Compute the covariance matrix for the standardized data.
   - Covariance matrix $C = \frac{1}{n-1} X^T X$
5: **Step 3:** Determine the covariance matrix C's eigenvalues and eigenvectors.
   - Solve $Cv = \lambda v$, where $\lambda$ represents the eigenvalues and $v$ represents the eigenvectors.
6: **Step 4:** Identify and select the top $k$ eigenvectors by arranging the eigenvalues in descending order and picking those corresponding to the largest eigenvalues.
7: **Step 5:** Construct the projection matrix $W$ from the chosen $k$ eigenvectors.
8: **Step 6:** Project the standardized data into the new $k$-dimensional space.
   - Transform the data using $Z = XW$, where $Z$ is the projected dataset in reduced dimensions.

---

Fig. 5: **Target class distribution**



Target Class Distribution (CKD vs. Not CKD)

*CKD Dataset with PCA*

In order to decrease the dataset's[10] dimensionality while preserving as much information as feasible, principal component analysis, or PCA, is utilized. When working with multi-attribute data sets, PCA is especially helpful because it greatly reduces the number of attributes in the data, making it simpler.The original CKD dataset contains input features 24. By measuring the contribution of each factor to the outcome using PCA, we can build a more accurate and efficient predictive model. The algorithm of PCA [9] is shown below.

## V. RESULTS

The F1-score, accuracy, precision, and recall were the main metrics employed in this research.[3]Table 2 shows how different values for the parameters in each model produced different results.[7] On the original CKD dataset[10], the XG Boost model showed its highest performance, with an accuracy of 0.9833.When PCA[9] was applied to the data set, this accuracy increased to 0.9916. On the original CKD dataset, models such as AdaBoost, RF, SVM, LR, and DT have achieved similar accuracy (0.9833) and after PCA, the XG Boost model achieved higher accuracy 0.9916 (Table 3).Tables with recall, precision, F1-score, training test accuracy, confusion matrix data represent summary of research findings for each model The results of the study showed that many models have impressive automated performance in CKD detection, with an accuracy rate greater than 97% are listed in Table 1.[7] These findings suggest that models developed for this work can accurately predict chronic kidney disease, potentially leading to improvements in biomedical treatment.In previous study, the GB algorithm obtained the best accuracy, by far, 0.990.

However, in this study, our XGBoost model proved to be very accurate, scoring 0.992 after applying PCA and 0.983 in the original dataset.

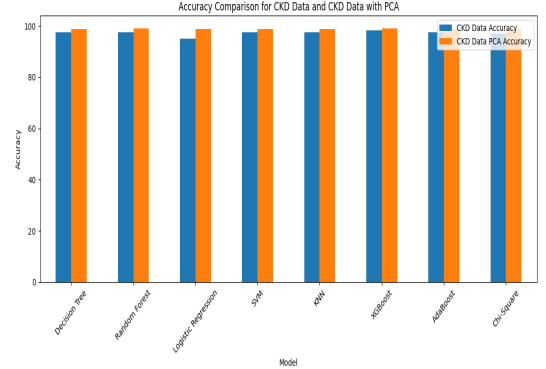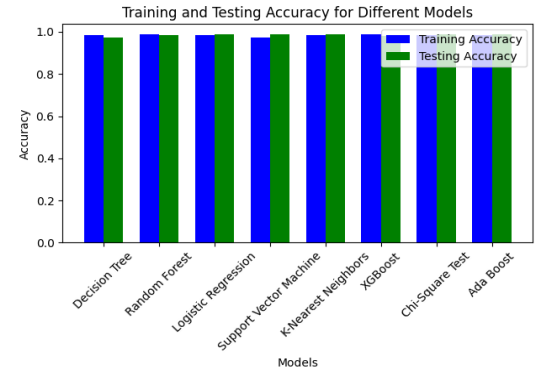Fig. 6: **The specified model's performances**



Accuracy Comparison for CKD Data and CKD Data with PCA

Fig. 7: **Training vs Testing accuracy of different models**



Training and Testing Accuracy for Different Models

## VI. COMPARITIVE ANALYSIS

TABLE II: Algorithm performance on the original CKD dataset [7]

| Classifiers | Accuracy | Precision | Recall | F1-Score |
|---|---|---|---|---|
| AdaBoost | 97.5 | 97.66 | 97.5 | 97.51 |
| Decision Tree | 97.5 | 97.6 | 97.5 | 97.5 |
| XGBoost | 98.75 | 98.8 | 98.75 | 98.75 |
| Random Forest | 97.5 | 97.6 | 97.5 | 97.5 |
| Logistic Regression | 95.0 | 95.62 | 95.0 | 95.06 |
| SVM | 97.5 | 97.66 | 97.5 | 97.51 |
| KNN | 97.5 | 97.6 | 97.5 | 97.5 |
| Chi Square | 98.75 | 98.79 | 98.75 | 98.75 |

Models for chronic kidney disease (CKD) are particularly useful for outcome prediction by identifying high-risk populations.clinical data can be used for this approach for problems seen in routine clinical practice.

TABLE III: Performance of different techniques using PCA on the CKD dataset

| Classifiers | Accuracy | Precision | Recall | F1-Score |
|---|---|---|---|---|
| AdaBoost | 98.75 | 98.79 | 98.75 | 98.75 |
| Decision Tree | 98.75 | 98.77 | 98.75 | 98.74 |
| XGBoost | 99.12 | 99.08 | 99.06 | 99.08 |
| Random Forest | 98.75 | 98.79 | 98.75 | 98.75 |
| Logistic Regression | 98.75 | 98.79 | 98.75 | 98.75 |
| SVM | 98.75 | 98.79 | 98.75 | 98.75 |
| KNN | 98.75 | 98.79 | 98.75 | 98.75 |
| Chi Square | 98.75 | 98.79 | 98.75 | 98.75 |

TABLE IV: Model Training and Testing Parameters for CKD Dataset with PCA

| Model | Training Parameters | Testing Parameters |
|---|---|---|
| Decision Tree | 3200 | 800 |
| Random Forest | 3200 | 800 |
| Logistic Regression | 2880 | 720 |
| SVM | 2880 | 720 |
| KNN | 3200 | 800 |
| XGBoost | 3200 | 800 |
| AdaBoost | 2880 | 720 |

However, building this model is difficult due to limited data.In particular, the dataset [10] used for this study contains only 400 samples, which is small and may compromise the validity of the results.table 3 shows the total number number of trainable and testable parameters used in dataset while pca.

## VII. CONCLUSION & FUTURE ENHANCEMENT

This research developed a machine learning model towards identifying the beginnings of chronic renal failure. Models were trained and validated focusing on identifying and removing irrelevant features to increase prediction accuracy using previously identified data and analyzing the relationship between type reference parameters showed that hemoglobin , albumin, and specific gravity are important prognostic indicators for chronic renal failure.The initial preprocessing of the CKD data set was done to ensure sure that machine learning pattern recognition was followed. We then used principal component analysis (PCA) to identify significant variables associated with prognosis of chronic kidney disease.

## REFERENCES

[1] Ammirati, A. L. (2020). Chronic kidney disease. Revista da Associação Médica Brasileira, 66(Suppl 1), s03-s09.

[2] Aljaaf, A. J., Al-Jumeily, D., Haglan, H. M., Alloghani, M., Baker, T., Hussain, A. J., & Mustafina, J. (2018, July). Early prediction of chronic kidney disease using machine learning supported by predictive analytics. In 2018 IEEE congress on evolutionary computation (CEC) (pp. 1-9). IEEE.

[3] Chittora, P., Chaurasia, S., Chakrabarti, P., Kumawat, G., Chakrabarti, T., Leonowicz, Z., ... & Bolshev, V. (2021). Prediction of chronic kidney disease-a machine learning perspective. IEEE access, 9, 17312-17334.

[4] Gudeti, B., Mishra, S., Malik, S., Fernandez, T. F., Tyagi, A. K., & Kumari, S. (2020, November). A novel approach to predict chronic kidney disease using machine learning algorithms. In 2020 4th International Conference on Electronics, Communication and Aerospace Technology (ICECA) (pp. 1630-1635). IEEE.

[5] Hossain, M. M., Swarna, R. A., Mostafiz, R., Shaha, P., Pinky, L. Y., Rahman, M. M., ... & Iqbal, M. S. (2022). Analysis of the performance of feature optimization techniques for the diagnosis of machine learning-based chronic kidney disease. Machine Learning with Applications, 9, 100330.

[6] Islam, M. A., Akter, S., Hossen, M. S., Keya, S. A., Tisha, S. A., & Hossain, S. (2020, December). Risk factor prediction of chronic kidney disease based on machine learning algorithms. In 2020 3rd international conference on intelligent sustainable systems (ICISS) (pp. 952-957). IEEE.

[7] Islam, M. A., Majumder, M. Z. H., & Hussein, M. A. (2023). Chronic kidney disease prediction based on machine learning algorithms. Journal of pathology informatics, 14, 100189.

[8] Elkholy, S. M. M., Rezk, A., & Saleh, A. A. E. F. (2021). Early prediction of chronic kidney disease using deep belief network. IEEE Access, 9, 135542-135549.

[9] Roweis, S. (1998). Em algorithms for pca and spca. Advances in NeuralInformation Processing Systems.

[10] DatasetLink: https://www.kaggle.com/datasets/mansoordaku/ckdisease.

[11] Venkatrao, K., & Kareemulla, S. (2023). HDLNET: a hybrid deep learning network model with intelligent IoT for detection and classification of chronic kidney disease. IEEE Access.

[12] Elkholy, S. M. M., Rezk, A., & Saleh, A. A. E. F. (2021). Early prediction of chronic kidney disease using deep belief network. IEEE Access, 9, 135542-135549.

[13] Rafi, S., & Das, R. (2021, December). RNN encoder and decoder with teacher forcing attention mechanism for abstractive summarization. In 2021 IEEE 18th India council international conference (INDICON) (pp. 1-7). IEEE.

[14] Debnath, D., Das, R., & Rafi, S. (2022, February). Sentiment-based abstractive text summarization using attention oriented lstm model. In Intelligent Data Engineering and Analytics: Proceedings of the 9th International Conference on Frontiers in Intelligent Computing: Theory and Applications (FICTA 2021) (pp. 199-208). Singapore: Springer Nature Singapore.

[15] Rafi, S., & Das, R. (2021, November). A linear sub-structure with co-variance shift for image captioning. In 2021 8th International Conference on Soft Computing & Machine Intelligence (ISCMI) (pp. 242-246). IEEE.

[16] Rafi, S., & Das, R. (2023). Topic-guided abstractive multimodal summarization with multimodal output. Neural Computing and Applications, 1-16.

[17] Rafi, S., & Das, R. (2023, November). Abstractive Text Summarization Using Multimodal Information. In 2023 10th International Conference on Soft Computing & Machine Intelligence (ISCMI) (pp. 141-145). IEEE.

[18] Rafi, S., & Das, R. (2024). SCT: Summary Caption Technique for Retrieving Relevant Images in Alignment with Multimodal Abstractive Summary. ACM Transactions on Asian and Low-Resource Language Information Processing, 23(3), 1-22.