
Chronic Kidney Disease Prognostication Utilizing Hybrid Ensemble Learning

Akshaya K Shaji¹, Sithara R², Suryamol S³, Adersh P B⁴, Tintu Varghese⁵

^{1,2,3,4} PG - Department of Computer Application, Kristu Jyoti College of Management and Technology, Kottayam Kerala

⁵ Assistant Professor Department of Computer Application, Kristu Jyoti College of Management and Technology, Kottayam Kerala

ABSTRACT

The chronic kidney failure is a extreme fitness issue and if not detected and handled on the early degrees, it can be very lethal. Hence the essential goal of this paper is to increase a dependable system learning version which predicts the CKD with a high accuracy price. The CKD statistics set is downloaded from the famous UCI ML repository but it suffers from a variety of lacking values. To deal with the lacking values KNN Imputation is used. Feature selection is likewise completed with the assist of information benefit because the dataset is massive and consequently the fee of modelling can be very excessive. Various other pre-processing steps like label encoding and Min-max normalization is executed to obtain a clean dataset. After pre-processing, diverse ML algorithms like logistic regression, naïve bayes, synthetic neural network and random wooded area are implemented and their performances are in comparison with the assist of diverse overall performance metrics. A hybrid of Random Forest and Adaboost algorithm is proposed and it achieves a higher accuracy when in comparison to the opposite person element models and subsequently it may be proved that the proposed hybrid version is an awful lot better and accurate in diagnosing CKD. goal of this paper is to increase a dependable system learning version which predicts the CKD with a high accuracy price. The CKD statistics set is downloaded from the famous UCI ML repository but it suffers from a variety of lacking values. To deal with the lacking values KNN Imputation is used. Feature selection is likewise completed with the assist of information benefit because the dataset is massive and consequently the fee of modelling can be very excessive. Various other pre-processing steps like label encoding and Min-max normalization is executed to obtain a clean dataset. After pre-processing, diverse ML algorithms like logistic regression, naïve bayes, synthetic neural network and random wooded area are implemented and their performances are in comparison with the assist of diverse overall performance metrics. A hybrid of Random Forest and Adaboost algorithm is proposed and it achieves

Keywords: Chronic Kidney Disease (CKD), Machine Learning (ML), Integrated/Hybrid Model

1. Introduction

Chronic kidney sickness, also called the persistent kidney failure, describes the sluggish loss of kidney characteristic. The most vital function of the kidneys is to clear out wastes and extra fluids from the blood. CKD manner that the kidney does no longer work as expected and cannot nicely clear out blood. Dangerous tiers of fluid, electrolytes and wastes can build up in one's body when chronic kidney disease reaches the very last ranges and this could be very risky. Chronic kidney disease occurs whilst a disorder or circumstance impairs kidney feature, leading the kidney harm to get worse over several months or years. In the beginning few degrees of continual kidney ailment, one may have just a few signs and symptoms or signs. Chronic kidney disorder might not end up outstanding till the kidney characteristic is substantially impaired. Chronic kidney sickness might progress to stop-level kidney failure, that's fatal without synthetic filtering (dialysis) or a kidney transplant. There are five tiers of CKD and the most dangerous one is degree five because, at this stage, the kidneys are not able to do maximum of their features. The tiers of kidney ailment are dependent on how correct the kidneys can filter waste and additional fluid out of the blood. The kidneys will nevertheless be able to clear out waste from the blood

within the early levels of kidney disorder. In the later levels, the kidneys should work more difficult to dispose of waste and can stop functioning altogether. It is hard to find out the CKD degree of every affected person especially on the early ranges. Glomerular Filtration Rate (GFR) is the pleasant check to degree the extent of kidney feature and to decide the degree of continual kidney disorder of a person. It may be calculated from the outcomes of the patient's blood creatinine, age, race, gender, and other such values. CKD has affected the entire global but specifically for the nations that have a low or medium profits, it has been very disastrous. About 10% of the population international suffers from the persistent kidney disease, and the number of deaths each year is increasing seriously. In the beyond many years as in keeping with the studies and researches executed by means of various health companies, CKD is stated to have triggered plenty of demise and other intense health situations. The range of people who suffer from stop level renal ailment is also turning into high and this considered because the remaining level of CKD and for the patients to survive this, kidney transplantation or normal dialysis have to be executed. A higher accuracy when in comparison to the opposite person element models and subsequently it may be proved that the proposed hybrid version is an awful lot better and accurate in diagnosing CKD.

2. Significance of the Study

In India, there is a huge wide variety of the chronic kidney failure cases pronounced each yr that is surely very alarming. CKD can regularly spoil the whole function of kidney to filter out wastes. If this ailment isn't recognized and dealt with on the early degrees, one would possibly expand a permanent kidney harm. If this disorder keeps progressing, risky electrolytes can be amassed in high levels in someone's blood and motive the person to fall ill. Almost all the parts of a human body can be broken by the progression of CKD and this can be a huge danger. Due to CKD, someone would possibly be afflicted by the trouble of other sicknesses which can be equally risky like hypertension, anemia and so forth. These complications may show up very slowly that someone will no longer be capable of diagnose that he has the continual kidney disorder. Hence early detection of the continual kidney ailment is very essential because it prevents the sickness from getting worse which can be very lethal. CKD does no longer show any symptoms or any ailment related signs in the preliminary few degrees and without taking the check it isn't always possible to tell is a specific character has the CKD or no longer. If detected within the early levels, it'll be very helpful for the patients as they will be capable of get a well-timed treatment and for this reason the development of the disease to the further tiers may be stopped. A man or woman is greater possibly to get the continual kidney disease if he has a own family records of CKD or if he has high blood pressure or diabetes. If Human beings get to recognize about this disease in advance then the remedy manner and the restoration procedure might be a good deal simpler. It is hoped that human beings get to recognise about this ailment quicker with the minimal range of tests feasible and additionally at a low fee. Hence the main motivation behind growing a device learning methodology to diagnose the persistent kidney sickness is the capacity dangers due to the development of CKD to ESRD if no longer detected on the early stages and also the depressingly growing range of cases stated every year.

3. Review Of related Studies

In the past few years there is lots of research accomplished on detecting chronic kidney disease with the useful resource of different types of device gaining knowledge of techniques. Many systems getting to know algorithms like logistic regression, random forest, assist vector device, k-nearest neighbor, naive Bayes classifier and diverse neural networks have been studied and their overall performance was compared with the help of various overall performance metrics and loss which turned into used for the neural community by myself. The dataset for the researches have been acquired from various sources just like the UCI system learning repository, King Fahd University Hospital (KFUH) in Khobar and so

forth. In many research, the most important goal become to come across CKD with the least wide variety of predictions and several statistical assessments had been performed to discard the unwanted attributes. Various characteristic choice strategies have been used as they can be very beneficial in decreasing the charges. Some of the function choice strategies used were Correlation-based Feature Selection (CFS), fruit fly optimization algorithm (FFOA), Density based totally Feature Selection (DFS) and Relief set of rules. Even a Heterogeneous Modified Artificial Neural Network (HMANN) changed into evolved which used ultrasound pics to carry out diverse photo processing steps with the assist of device studying to detect CKD. There were also a few studies on choice tress to diagnose CKD. In one such take a look at it changed into proved that the J48 selection tree and random wooded area completed fantastically correct consequences while as compared to the other sort of tress in gadget learning that did not gain the acceptable bring about detecting CKD. In 2018 Almarashi A, Alghamdi M, and Mechai I delivered a specific kind of detecting approach for CKD by means of using Artificial Neural Network (ANN). The aspect elements of the newly developed device were one enter layer, one hidden layer, and one output layer. There changed into additionally research which proposed a unique methodology based totally on Extreme Gradient Boosting (XGBoost) version. In this model three exclusive function selection methods have been applied. Among the diverse algorithms as compared, ANN and Random Forest finished the first-class performance and SVM became the maximum extensively studied set of rules for diagnosing CKD. Logistic Regression and Naïve Bayes algorithms are fairly new on the subject of diagnosing CKD. It may be seen that maximum of the studies specifically focus at the establishment of fashions and attaining a excessive accuracy however there isn't enough studies on the records pre-processing techniques used. A complete system of dealing with the lacking values is not explained in depth. In maximum of the studies papers the rows with the missing values have been deleted that can result in a lack of crucial records. In different papers the suggest or the median strategies were used which aren't that suitable as they tend to add a variety of undesirable bias and variance to the dataset that could cause a huge problem even as developing the machine getting to know fashions. Most of the prevailing work are about person fashions of system learning algorithms and there is not enough research on Adaboost set of rules and included models that could attain a higher accuracy.

4.Objectives of The Study

The foremost goal of this paper is to diagnose CKD at the early tiers with the least feasible checks and fee and with a high accuracy charge. The paper additionally aims to efficiently deal with the missing values, gift in the CKD records set with the help of KNN Imputation. Feature selection is also carried out with the help of statistics benefit to locate the most critical capabilities that play a essential role in detecting CKD. Various system learning algorithms are carried out and analysed to come across CKD and the first-rate one with the nice performance and accuracy charge is located. Adaboost algorithm is likewise implemented because it boosts the overall performance of the weak classifiers. Finally, the misjudgements generated with the aid of the hooked-up fashions are analysed and an incorporated model is proposed that combines random woodland set of rules and AdaBoost algorithm that may gain a better accuracy and might consequently be an effective and reliable version to come across CKD.

5.Hypotheses of The Study

The hybrid model will achieve a higher accuracy rate when as compared to the person machine learning fashions. The use of KNN Imputation to handle the lacking values inside the dataset and the usage of data benefit as the function choice approach will further boom the accuracy.

6.The Proposed Methodology

The proposed methodology involves the various steps as follows:

6.1.Data Pre-processing

Data pre-processing in Machine Learning is a completely important to convert the raw dataset into a cleaned dataset set that may be desirable to apply variable system mastering algorithms.

6.1.1.Acquire the Dataset

The dataset for predicting the continual kidney disorder is attained from the UCI machine mastering repository which is a well-known source for all of the system studying datasets. The CKD dataset has four hundred affected person information and 25 attributed which might be either the signs and symptoms or other attributes related to the sickness like hypertension, blood stress, unique gravity, albumin etc. Among these 400 affected person records, 250 sufferers have the disease and the other a hundred and fifty of them do not have the ailment.

6.1.2. Import all the important libraries

In order to carry out information pre-processing using Python, diverse predefined Python libraries ought to be imported. All the libraries have a positive mission to do when it comes to gadget mastering. In this diagnosis, numerous libraries including numpy, pandas, matplotlib, sklearn and so forth. Have been imported and most of these libraries have a positive venture to carry out.

6.1.3. Import the dataset

The dataset which has been accumulated for the machine studying assignment is imported. While doing the dataset importing manner, one extra essential factor needs to be finished which involves extracting structured and independent variables. In this dataset, class is the structured variable and all of the different capabilities are impartial variables.

6.1.4 Feature Selection

A Machine Learning version can suffer from the issue of overfitting if the range of features end up equal or big. To prevent this from occurring its miles important to reduce the range of features inside the dataset. Feature selection refers to decreasing the number of features in a dataset while building a device gaining knowledge of version. Another advantage of decreasing the range of input attributes is the decrease inside the cost of constructing the ML model and in a few conditions, it'd even improve the accuracy of the version. There are diverse methods to perform characteristics election and in this paper Information Gain is used. Information benefit is selected over the alternative strategies as it comes underneath filter out method. The characteristic choice the usage of the filter strategies pick out the intrinsic residences of the features measured with the assist of univariate statistics aside from using the cross-validation performance. When compared to the other characteristic choice methods, clear out techniques are swifter and less computationally dearer than wrapper strategies. While dealing with a large dataset with plenty of capabilities, it's miles relatively much less costly and less difficult to make use of the filter strategies. Information gain of each variable with recognize to the target variable is observed to determine the maximum essential capabilities that have a crucial position within the prediction technique. The data benefit of a characteristic can be something from zero to 1. The capabilities with the very best statistics advantage can be retained as they are the most important ones and the capabilities with the least information gain could be discarded. The Information gain of each enter characteristic or the characteristic dataset with respect to the output characteristic.

6.1.5. Identify and manage missing values

Next it is very important to control the missing values in a dataset because it will cause a huge trouble later even as making use of the machine learning algorithms. Hence it is important to handle missing values which can be there in the dataset. The CKD dataset has numerous missing values as a number of the sufferers may forget about or miss to fill in sure values inside the real lifestyles. Hence, KNN imputation is going to be used to deal with missing values because it has confirmed to be effective in experiments. Every lacking price can then be dealt with by means of replacing them with the suggest

value of the k nearest neighbours, to accomplish that in popular Euclidean distance metric is used in default.

6.1.6. Encoding the categorical data

Categorical records in a dataset are those that have certain classes like on this dataset there are specific variables like packed cell quantity, high blood pressure and type. These specific facts might purpose a big assignment while building the model as system gaining knowledge of handiest deals with numbers and mathematics. So it's miles essential to transform those express variables into numbers. There are numerous techniques for this cause and in this paper Label Encoding of sklearn library goes for use.

6.1.7. Feature Scaling

In feature scaling all the input attributes are converted to cost of commonplace variety or distribution in order that it turns into easier to evaluate all of the variables on common grounds to construct a reliable model. In this dataset, it can be observed that the albumin and packed cellular extent columns do not have the identical variety or distribution of values. Packed mobile quantity has a better range when compared to the albumin's range and hence a right end result will now not be done as PCV dominates the albumin. Hence it will become very important to carry out feature scaling to the dataset. In this paper Min-Max Normalization goes for use to perform function scaling. This method converts a characteristic or commentary fee with distribution price inside the variety zero to 1.

6.1.8. Splitting the dataset

All the dataset must be split into the training subset and the checking out subset to proceed with the prediction. The education subset is used to perform schooling whereas the checking out subset is used to carry out testing. One is aware of the consequences within the schooling subset but one is not aware about the prediction end result or output of the trying out subset. The CKD dataset is going to be cut up inside the ratio 70:30.

6.2. Developing individual ML models and evaluating them

The various machine learning algorithms applied for prediction are:

- Logistic Regression
- Naïve Bayes
- Artificial Neural Network
- Random Forest

The performance of those classifiers are analysed primarily based on numerous metrics together with Accuracy, Precision, Recall, F-degree and loss and the first-rate classifier for diagnosing continual kidney disorder is observed.

6.3. Establishing the Integrated/Hybrid Model

Once the individual ML algorithms are analysed and as compared for the misjudgements, a hybrid version is advanced to similarly boom the accuracy. A hybrid of Random Forest and Adaboost algorithm might be proposed. Random Forest is chosen over ANN as it acquires a high accuracy and it's also very well matched with AdaBoost algorithm. ANN also suffers from the problem of unexplained behaviour of community so it is not used for the hybrid model. Boosting is the method of enhancing the performance of susceptible classifiers with the assist of an ensemble approach. This is carried out through developing a new version that maintains solving the issues and issues of the last few fashions. Until a perfect result is attained, this process of making a new model keeps persevering with to create a highly correct version. The final equation for hybrid version is given underneath the vulnerable classifier and θ_m is the load that corresponds to each classifier.

$$F(x) = \text{sign}(\sum_{m=1}^M \theta_m f_m(x))$$

Usually AdaBoost uses Decision Tree Classifier as default vulnerable learner for education motive but any sort of machine getting to know algorithms may be utilised as lengthy as it accepts the parameter called weights. In this research Random Forest is used because the base classifier as it's far tons better

than selection tree. So the fm within the formulation of the included model will be the random wooded area set of rules. The hybrid model will diagnose the CKD greater efficiently and will acquire a higher accuracy when as compared to the person models.

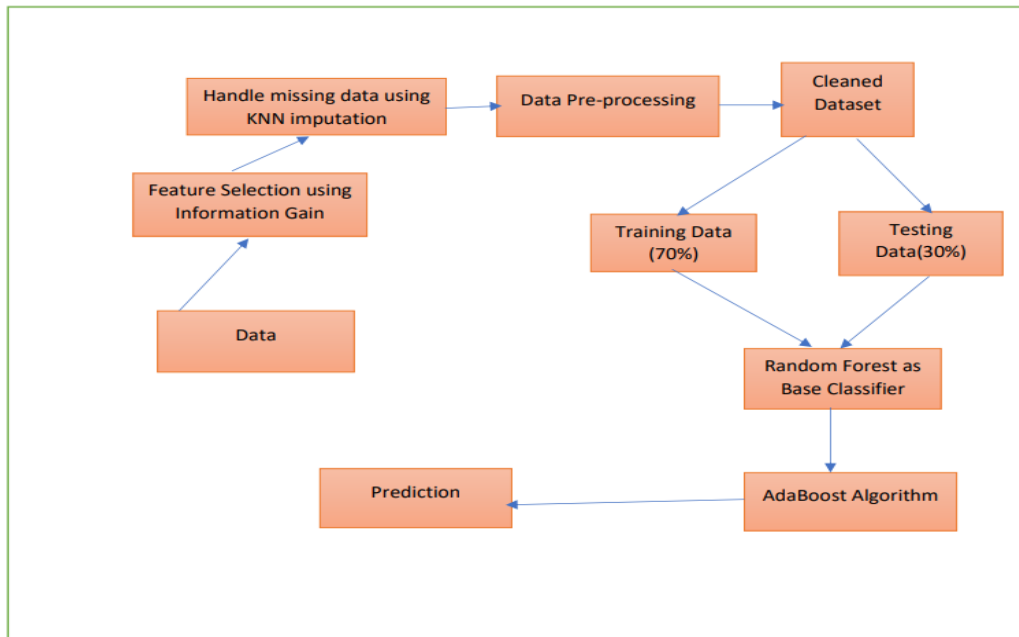


Figure 1: Architecture diagram of Hybrid Model

7.Results and Discussion

The end result of all of the system getting to know algorithms had been in comparison and analyzed with the assist of various overall performance metrics. The experiments had been conducted the usage of Python three.6.7 programming language and it can be performed both the usage of the Jupyter Notebook internet application or Google collab. Many libraries have been used for the implementation and one such vital library is Sciket-study, that is a totally useful library to develop ML models. Various overall performance metrics in the confusion matrix are considered on this studies. The experimental results of all the advanced models is given within the table given below.

Table 1. Performance table

Classifiers	Accuracy	Precision	Recall	F-Measure	Loss
Logistic Regression	0.97701	0.95349	1.0	0.97619	
Naïve Bayes	0.94253	0.89130	1.0	0.94253	
Artificial Neural Network	1.0	1.0	1.0	1.0	0.01262
Random Forest	0.98851	0.97619	1.0	0.98795	
Integrated model	1.0	1.0	1.0	1.0	

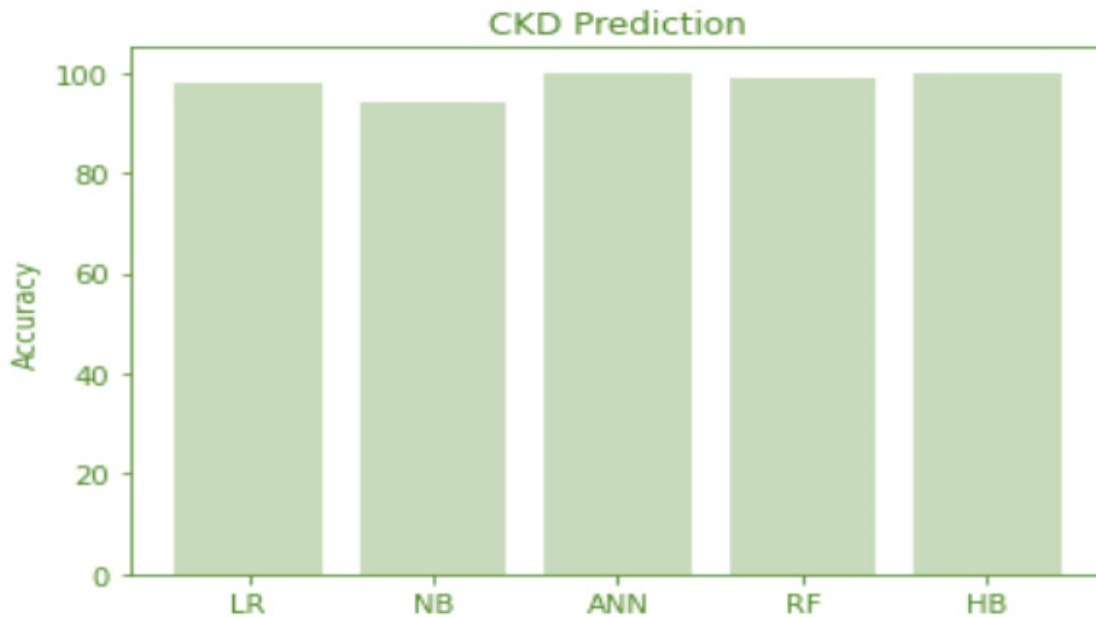


Figure 2: Accuracy graph of the models

The consequences prove the feasibility of the proposed methodology. The function selection approach of Information Gain turned into used to extract the most vital functions which have a vital function in the diagnosis of the persistent kidney ailment. It turned into observed that precise capabilities like particular gravity, albumin, serum creatinine, hemoglobin, packed cell quantity, crimson blood mobile depend and hypertension had been extra vital than the alternative capabilities as they'd a better facts Gain. With the help of KNN Imputation, LOG, NB, ANN and RF had been capable of reap a higher overall performance than the instances when random imputation, deleting rows with missing values or mean and mode Imputation were used. These methods are not superior as they upload a bias and variance to the model. KNN Imputation dealt with the missing values with the aid of replacing them with the suggest of k nearest neighbours. In the experiment imputer optimization was accomplished and it was determined that five become most reliable value for k for the CKD dataset. From the evaluation results, all the fashions have a first-rate performance in opposition to detecting CKD with a terrific accuracy. From the desk it can be seen that Random woodland and Artificial neural community done the first-class overall performance when as compared to the other fashions. To similarly growth the performance and reliability of the individual ML algorithms, an incorporated version is proposed which combines Random Forest and AdaBoost set of rules and it performed a better accuracy. Random Forest is selected over ANN as it extra compatible with AdaBoost set of rules and ANN also suffers from the trouble of unexplained behavior of the network. Whenever ANN produces a probing solution, it does now not give an explanation for the solution as in why and how they're produced. This reduces the notion inside the ANN. Adaboost is normally used to enhance the performance of the choice tree set of rules and if Random Forest is used as opposed to choice tree it in addition improves the overall performance.

8. Conclusion

This research paper has tested the capacity to stumble on CKD using device mastering methodologies. The intention was successfully achieved by means of applying and analysing diverse ML algorithms like logistic regression, random forest, naive Bayes classifier, and synthetic neural network and the performance of those algorithms were compared. KNN imputation became used to handle the lacking

values and Information Gain changed into used as the characteristic selection method. By analyzing the problems and shortcomings of the character ML models, a hybrid version changed into proposed that mixes the AdaBoost and Random Forest algorithms which become capable of achieve a higher accuracy fee. Hence it can be speculated that this hybrid method may be used in the realistic diagnosis of CKD and it could gain a acceptable effect. It can also be mentioned that this system might be useful to the medical facts of the opposite diseases in real scientific analysis. In the destiny, a bigger dataset may be used by attaining a more quantity of patient statistics from numerous hospitals and health companies to improve the accuracy of the prediction. It is was hoping that the performance of the device will be an increasing number of accurate with an growth in the length and high-quality of the dataset.

9. References

1. Qin, J., Chen, L., Liu, Y., Liu, C., Feng, C. and Chen, B., 2019. A machine learning methodology for diagnosing chronic kidney disease. *IEEE Access*, 8, pp.20991-21002.
2. Almasoud, M. and Ward, T.E., 2019. Detection of chronic kidney disease using machine learning algorithms with least number of predictors. *International Journal of Soft Computing and Its Applications*, 10(8).
3. Ahmad, M., Tundjungsari, V., Widiанти, D., Amalia, P. and Rachmawati, U.A., 2017, November. Diagnostic decision support system of chronic kidney disease using support vector machine. In 2017 second international conference on informatics and computing (ICIC) (pp. 1-4). IEEE.
4. Al Imran, A., Amin, M.N. and Johora, F.T., 2018, December. Classification of chronic kidney disease using logistic regression, feedforward neural network and wide & deep learning. In 2018 International Conference on Innovation in Engineering and Technology (ICIET) (pp. 1-6). IEEE.
5. Tekale, S., Shingavi, P., Wandhekar, S. and Chatorikar, A., 2018. Prediction of chronic kidney disease using machine learning algorithm. *International Journal of Advanced Research in Computer and Communication Engineering*, 7(10), pp.92-96.
6. Alassaf, R.A., Alsulaim, K.A., Alroomi, N.Y., Alsharif, N.S., Aljubeir, M.F., Olatunji, S.O., Alahmadi, A.Y., Imran, M., Alzahrani, R.A. and Alturayef, N.S., 2018, November. Preemptive diagnosis of chronic kidney disease using machine learning techniques. In 2018 international conference on innovations in information technology (IIT) (pp. 99-104). IEEE.
7. Ma, F., Sun, T., Liu, L. and Jing, H., 2020. Detection and diagnosis of chronic kidney disease using deep learningbased heterogeneous modified artificial neural network. *Future Generation Computer Systems*, 111, pp.17-26.
8. Amirgaliyev, Y., Shamiluulu, S. and Serek, A., 2018, October. 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) (pp. 1-4). IEEE.
9. Sobrinho, A., Queiroz, A.C.D.S., Da Silva, L.D., Costa, E.D.B., Pinheiro, M.E. and Perkusich, A., 2020. Computer-aided diagnosis of chronic kidney disease in developing countries: A comparative analysis of machine learning techniques. *IEEE Access*, 8, pp.25407-25419.
10. Almansour, N.A., Syed, H.F., Khayat, N.R., Altheeb, R.K., Juri, R.E., Alhiyafi, J., Alrashed, S. and Olatunji, S.O., 2019. Neural network and support vector machine for the prediction of chronic kidney disease: A comparative study. *Computers in biology and medicine*, 109, pp.101-111.
11. Ogunleye, A. and Wang, Q.G., 2018, June. Enhanced XGBoost-based automatic diagnosis system for chronic kidney disease. In 2018 IEEE 14th International Conference on Control and Automation (ICCA) (pp. 805810). IEEE.
12. Jiongming Qin, Lin Chen, et al., "A Machine Learning Methodology for Diagnosing Chronic Kidney Disease", *IEEE Access*, vol 8, pp. 20991-20993, 2020



13. V. Jha, G. Garcia-Garcia, K. Iseki, Z. Li, S. Naicker, B. Plattner, R. Saran, A. Y. Wang, C.W. Yang (2013),” chronic kidney disease: global dimension and perspectives”, The Lancet.
14. Kunwar, V., Chandel, K., Sabitha, A. S., & Bansal, A. (2016). Chronic Kidney Disease analysis using data mining classification techniques. In 2016 6th International Conference-Cloud System and Big Data Engineering (Confluence) (pp. 300-305)
15. A. Salekin and J. Stankovic, (2016) “Detection of chronic kidney disease and selecting important predictive attributes,” IEEE International Conference on Healthcare Informatics.