

ABSTRACT

ABSTRACT

chronic kidney disease refers to the condition of kidneys caused by conditions, diabetes, glomerulonephritis or high blood pressure. These problems may happen gently for a long period of time, often without any symptoms. It may eventually lead to kidney failure requiring dialysis or a kidney transplant to preserve survival time. So the primary detection and treatment can prevent or delay of these complications. The aim of this work is to reduce the diagnosis time and to improve the diagnosis accuracy through classification algorithms. The proposed work deals with classification of different stages in chronic kidney diseases using machine learning algorithms. The experimental results performed on different algorithms like Naive Bayes, Decision Tree, K-Nearest Neighbour and Support Vector Machine. The experimental result shows that the K-Nearest Neighbour algorithm gives better result than the other classification algorithms and produces 98% accuracy. Keywords: Chronic Kidney Disease (CKD), Machine Learning (ML), End-Stage Renal Disease (ESRD), Cardiovascular disease, data mining, machine learning,

1. INTRODUCTION

Data mining is a used for the healthcare industry to enable health systems systematically. It uses data for analytics to identify incompetence and best practices that increase the care and reduce costs. Medical treatment is facing a challenge of knowledge discovery from the growing volume of data. Nowadays huge data are collected continuously through health examination and medical treatment. Classification rules are typically useful for medical problems that have been applied mainly in the area of medical diagnosis. Moreover, various machinelearning (ML) techniques have been applied to the field of medical treatments over the past few years. Chronic kidney disease (CKD) is a worldwide common health problem, with predictable lifetime risk of >50%, higher than that for invasive cancer, diabetes and coronary heart diseases. CKD is a long term disorder caused by damage to both kidneys [1], [2]. There is no single cause and the damage is typically permanent and can lead to ill health. In some cases dialysis or transplantation may become essential. Diabetes mellitus is also becoming more common in one cause of CKD. Chronic kidney disease is become more frequently in older people and consequently is likely to increase in the population as a whole. People with chronic kidney diseases are at higher risk of cardiovascular disease and they should be recognized early so that appropriate preemptive measures can be taken [3,4]. CKD is defined as the presence of kidney damagerevealed by the abnormal

ABSTRACT

albumin excretion or decreased kidney function. The disease is quantified by measured or estimated by Glomerular Filtration Rate (GFR) that persists for more than 3 month of the CKD patients. The glomerular filtration rate (GFR) is the best indicator of how well the kidneys are working. The National Kidney Foundation published treatment guidelines for identified five stages of CKD based on diminishing GFR measurements. The guidelines mention different actions based on the stage of kidney disease [5]. A GFR of 90 or above is considered as normal. Even with a normal GFR, it may be at increased risk for developing CKD if the patients have diabetes, blood pressure in high, or a family history of kidney disease. The risk increases with age over 65 are more than twice as likely to develop CKD as people between the ages of 45 and 65. The remaining paper is organized as follows: Section II deals with literature survey about chronic kidney diseases. In section III methodologies used for classifying chronic kidney diseases are discussed. Section IV deals with experimental and its results. Section V gives prediction of chronic kidney diseases with various performances and its future works.

2. LITERATURE SURVEY

Miguel A. et al. [6] proposed an approach for the management of alarms related to monitoring CKD patients within the eNefro project. The results proof the pragmatism of Data Distribution Services (DDS) for the activation of emergency protocols in terms of alarm ranking and personalization, as well as some observations about security and privacy. Christopher et al. [7] discussed a contextualized method and possibly more interpretable means of communicating risk information on complex patient populations to timeconstrained clinicians. Dataset was collected from American Diabetes Association (ADA) of 22 demographic and clinical variables related to heart attack risk for 588 people with type2 diabetes. The method and tool could be encompasses to other risk-assessment scenarios in healthcare distribution, such as measuring risks to patient safety and clinical recommendation compliance. Srinivasa R. Raghavan et al. [8] explored reviews the literature on clinical decision support system, debates some of the difficulties faced by practitioners in managing chronic kidney failure patients, and sets out the decision provision techniques used in developing a dialysis decision support system. Ricardo T. Ribeiro et al. [9] proposed a method, called clinical based classifier (CBC), discriminates healthy from pathologic conditions. A large multimodal feature database was specifically built for this study. It containing chronic hepatitis, 34 compensated cirrhosis, and 36 decompensated cirrhosis cases, all validated after histopathology examination by liver biopsy. The CBC classification outperformed the nonhierarchical one counter to all scheme, achieving better accuracy. Mitri F.G. et al. [10] presented an ultrasound-based modality

ABSTRACT

complex to stiffness and free from speckle noise and owns some advantages over the conventional ultrasound imaging in terms of the quality. Chih-Yin Ho et al. [11] presented a computer-aided diagnosis tool based on analyzing ultrasonography images and the system could detect and classify various stages of CKD. The dataset was collected thousands of ultrasonic images from patients with kidney diseases, and the selected typical CKD images were applied to be preanalyzed and trained for assessment. The calculated changeover locations are reference indicators could be responsible for physicians an auxiliary and objective computer-aid diagnosis tool for CKD identification and classification. Al-Hyari et al. [12] proposed a new clinical decision support system for identifying patients with CRF. Some data classification algorithms including Artificial Neural Networks, Decision Tree and Naive Bayes are developed and applied to diagnose patients with CRF and determine the evolution stage of the disease.

The dataset containing 102 instances is collected from patients' records and used for this study. The attained results showed that the developed decision tree algorithm is the most accurate CRF classifier (92.2%) when compared to all other algorithms used in this study. Kuo-Su Chen et al. [13] established a detection system based on computer vision and machine learning techniques for simplifying diagnosis of CKD and different stages of CKD. The proposed system required average time of 0.016 seconds for feature extraction and classification of each testing case. The results presented that the system could produce reliable diagnosis based on noninvasive ultrasonography methods and which could be measured as the most proper clinical diagnosis and medical treatment for CKD patients.

Anne Rogers et al. [14] discussed to discover patients' experiences disclosure of CKD in primary care settings. The dataset contains purposive sample of 26 patients, with a mean age of 72 years were cross-examined using constant relative techniques. This study challenges the assumptions characteristic in extensive health policy objectives that are increasingly built on the notion of responsible patients and the ethos of the active support of self-management for pre-conditions. Mohammed Shamim Rahman et al. [15] described the effect of chronic kidney disease (CKD) on morbidity and mortality following Trans catheter aortic valve implantation (TAVI) including patients on hemodialysis, often excluded from randomized trials. There are 118 consecutive patients underwent TAVI 63 were considered as having (CKD) and 55 not having (No-CKD) significant pre-existing CKD. The result shows TAVI is a safe, suitable treatment for patients with pre-existing CKD, though carefulness must be trained, particularly in patients with pre-existing diabetes mellitus and elevated preoperative serum creatinine levels as this confers a greater risk of AKI development, which is associated with increased short term post-operative mortality.

3.DATASET ATTRIBUTES

We have downloaded Chronic Kidney Disease datasets from publically available data

ABSTRACT

from UCI Machine Learning Repository [16]. Table 1 gives a list of all the attributes taken

	ATTRIBUTES	VALUES
	1.Age	Numerical
	2.blood	Numerical
sg(.005,1.010,1.015,1.020,1.025)	3.specific gravity	Nominal
al-(0,1,2,3,4,5)	4.Albumin	Nominal
su-(0,1,2,3,4,5)	5.Sugar	Nominal
rbc-(normal,abnormal)	6.red blood cells	Nominal
pc-(normal,abnormal)	7.pus cells	Nominal
(present,notpresent)	8.pus cell clumps	Nominal pcc-
(present,notpresent)	9.Bacteria	Nominal ba -
	10.blood glucose random	Numerical
	11.blood urea	Numerical
	12.serum creatinine	Numerical
	13.Sodium	Numerical
	14.Potassium	Numerical
	15.Haemoglobin	Numerical
	16.Packed cell volume	Numerical

ABSTRACT

htn - (yes,no)	17.white blood cell count	Numerical
	18. red blood cell count	Numerical
	19.Hypertension	Nominal
dm - (yes,no)	20.diabetes mellitus	Nominal
cad - (yes,no)	21.coronary artery disease	Nominal
	22.Appetite	Nominal
appet - (good,poor)	23.pedal edema	Nominal
pe - (yes,no)	24.Anemia	Nominal
	25.Class	Nominal
class - (ckd,notckd)		

4.CLASSIFICATION ALGORITHMS

A. Naïve Bayes:

The Naive Bayesian classifier is based on Bayes' theorem with independence assumptions among predictors. A Naive Bayesian model is easy to build, with no complex iterative parameter assessment which makes it especially useful for very large datasets. Even though it's simple, the Naive Bayesian classifier often does unexpectedly well and is widely used because it often outperforms more refined classification methods. Bayes theorem delivers a way of calculating the posterior probability, $P(c|x)$, from $P(c)$, $P(x)$ and $P(x|c)$. Naive Bayes classifier assumes that the effect of the value of a predictor (x) on a given class (c) is independent of the values of other predictors. This hypothesis is called class conditional independence [16].

$$P(c/x)=p(x/c) p(c)/p(x)$$

(1) $P(c|x)$ is the posterior probability of given predictor. $P(c)$ is the prior probability of class. $P(x|c)$ is the likelihood which is the ability of predictor

ABSTRACT

given class. $P(x)$ is the prior probability of predictor.

B. Decision Tree:

Decision tree builds classification or regression models in the form of a tree like structure. It breakdowns a dataset into smaller and smaller subsets while at the same time an associated decision tree is incrementally established. The last result is a tree with decision nodes and its leaf nodes. Decision nodes have two or more branches. Leaf node represents a classification or decision. The uppermost decision node in a tree which resembles to the finest predictor called root node. Decision trees can switch both categorical and numerical data values. The core algorithm for building decision trees is called ID3 by J. R. Quinlan which employs a top-down and greedy search over the space of possible branches with no backtracking.

Algorithm

- Start with single node N , with training data D .
- If all the data in D belongs to same class, then N becomes leaf. Otherwise attribute 'A' is selection method based on splitting criterion.
- The instance in 'D' is partitioned accordingly
- Apply algorithm recursively to each subset in 'D' to each subset in 'D' to form decision tree.

The algorithm uses Entropy and Information Gain to construct a decision tree [17].

Entropy ID3 algorithm uses entropy to calculate the similarity of If the sample is completely similar the entropy is zero and if the sample is an equally divided it has entropy of one.

$$E(S) = -\sum p_i \log_2 p_i \quad \text{---(2)}$$

Information Gain

The information gain is based on the decrease in entropy after a dataset is split on an attribute. Creating a decision tree is all about finding attribute that returns the highest information gain (i.e., the most homogeneous branches).

$$\text{GAIN}(T,X) = \text{ENTROPY}(T) - \text{ENTROPY}(T,X) \quad \text{--- (3)}$$

C. K- Nearest Neighbour:

K nearest neighbours is a simple algorithm that stores all available cases and classifies new cases based on a similarity measure (e.g., distance functions). KNN has been used in statistical estimation and pattern recognition already in the beginning of 1970's as a non-parametric technique. A case is classified by a majority vote of its neighbours, in the case being assigned to the class most common between its K nearest neighbours measured by a distance function. If $K = 1$, then the case is simply allocated to the class of its nearest neighbour [18].

Algorithm

- Training set: $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$.
- Assume $X = (x:(1), x:(2), \dots, x:(d))$ is a

ABSTRACT

dimensional feature vector of real numbers for all i

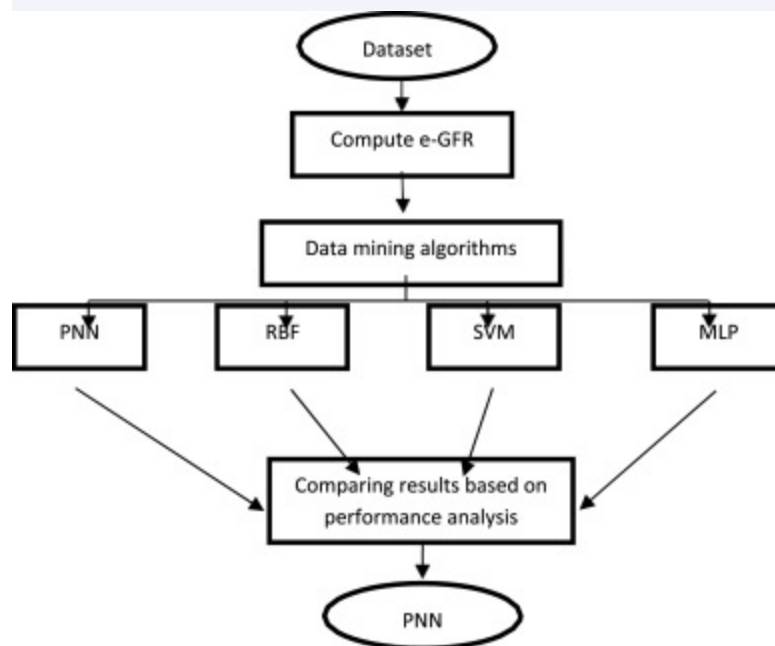
- Y_i is a class label $\{1 \dots C\}$, for all i .
- Find the closest point X_j to X_{new} using distance measures.
- Classify by Y_{knn} = majority vote among the K points.

D. Support Vector Machine A Support Vector Machine (SVM) performs classification by finding the hyper plane that maximizes the margin between the classes. The vectors (cases) that define the hyper plane are the support vectors [19].

Algorithm

- Define an optimal hyper plane: maximize margin.
- Extend the above definition for non-linearly separable Problems: have a penalty term for misclassifications
- Map data to high dimensional space where it is easier to classify with linear decision surfaces: reformulate problem so that data is mapped implicitly to this space. For this type of SVM, training involves the minimization of the error function.

5.METHODOLOGY



components of methodology of chronic kidney disease prediction

1. Data Collection

ABSTRACT

In this research paper we have used Real world data set for predicting CKD status of a patient. The data collected is widely used data and is available at UCI Machine Learning Repository. This Real data belongs to Apollo Hospital in Tamilnadu, India over a period of 2 months. The data set available is specifically used for Chronic Kidney Disease research. It consists of record of 400 people with their respective 25 CKD related attributes. The data consisted of real numbers, Decimal values and Nominal values.

2. Data Pre-Processing

Data pre-processing is a way to convert the noisy and huge data into relevant and clean data, as the data available is Real world data, so it contains inaccurate data, missing values and other Noisy data, for removing this inconsistent data from the Dataset, the proposed system have to clean the raw data.

following are the step involved in data pre processing

2.1. Looking Up For Proper Format: As we have made our model using python, so we need a csv file (comma separated value) for our code. The data downloaded is in the form of RAR file, so we extract the data from the text file available and save it into a csv file so that our python code can read it. This is the first most important step, if the data is not available in requires format then we cannot design the classification model.

2.2. Finding Missing Values: When the data collected is real world data, and then it will contain missing values. This brings more change in the prediction accuracy. Sometimes these missing values can be simply deleted or ignored if they are not large in number. It is the simplest way to handle the missing data but it is not considered healthy for the model as the missing value can be an important attribute contributing to the disease. The missing values can also be replaced by zero this will not bring any change as whole, but this method cannot be much yielding. So an efficient way to handle missing values is to use mean, average of the observed attribute or value. This way we lead to more genuine data and better prediction results.

2.3. Data Transformation: In this step we transform the given real data into required format. The data downloaded consist of Nominal, Real and Decimal values. In this step we convert the Nominal data into numerical data of the form 0 and 1. The positive value is assigned the value of 1 and the negative value is assigned the value of 0. Now the resultant csv file comprises of all the integer .

6.FEATURE SELECTION

ABSTRACT

In this step we select subset of relevant attributes from the total give attributes. This stage helps in reducing the dimensionality and making the model simpler and easy to use, thus leading to short training time and high accuracy.

To obtain highly dependent features for CKD prediction we have used Correlation and dependence method. The term correlation can be defined as mutual relationship between two. In this those attributes are chosen which highly influence the occurrence of Chronic Kidney Disease

. By using the correlation it is found that 5 attributed were highly correlated to the occurrence of CKD from the total of 25 attributes.

The 5 attributes selected from a total of 25 attributes are:

1. specific gravity
2. diabetes mellitus_N
3. albumin
4. packed cell volume
5. red blood cells_N

7. RESULTS

This study is carried to predict whether a patient is suffering from Chronic Kidney Disease or not. This Prediction model is created in Python programming language. In our classification model we have used K-Nearest Neighbour and Naïve Bayes as our classification algorithms; both the classification algorithms were applied to the same data set collected from UCI Repository.

Filtered dataset of 400 people with 5 CKD related attributes from a total of 25 attributes is used. The Filtered attributed are - Specific gravity, diabetes mellitus_N, albumin, packed cell volume and red blood cells_N.

Table2 . Predictive accuracy of classification		
algorithms	Algorithm	Accuracy
	Naïve Bayes Classifier	96.25%
	K-Nearest Neighbour	100%

ABSTRACT

Table 2 represent the prediction accuracy of both Naïve Bayes and K-Nearest Neighbour algorithms. Both the prediction accuracies are compared. Naïve Bayes performed with an accuracy of 96.25% and KNN performed with an accuracy of 100%.

Figure 3 is the graphical representation of both the prediction algorithms, KNN with 100% accuracy and Naïve Bayes with an accuracy of 96.25%. In the above figure x-axis represents the accuracy value and the y-axis represents the algorithm used. The above results highlight that accuracy of KNN algorithm is 3.75% higher than Naïve Bayes classification algorithm. The experimental results show that Chronic Kidney Disease can be better predicted by using K-Nearest Neighbour algorithm with 100% accuracy. The advantage of this research is that it will help Doctors to easily predict CKD with high accuracy and precision in less time period.

8.ADVANTAGES AND LIMITATIONS

There are a number of potential advantages to including KTRs in the CKD classification. Increased recognition of CKD may facilitate implementation of therapeutic strategies to delay progression of kidney function decline or prevent CKD related metabolic complications and CVD. Inclusion of KTRs in a simple severity - based kidney disease classification schema may improve communication between clinicians, enhance public education and facilitate research. Finally, a uniform disease classification and action plan including all patients irrespective of the need or type of renal replacement therapy (i.e. dialysis or transplantation), may enhance the continuity of patient care. These potential advantages should be weighed against the loss of precision inherent to adopting a generic disease classification system.

There are a number of fundamental differences between transplant and nontransplant CKD patients that may limit the applicability of the CKD classification in KTRs. Specific considerations include whether kidney function is an appropriate index of allograft health, whether five CKD stages are needed in KTRs, and whether

ABSTRACT

the clinical action plans associated with each CKD stage are appropriate for KTRs. Additional issues include how and when the CKD classification should be implemented in KTRs, and development of a research agenda to answer the uncertainties regarding the inclusion of KTRs in the CKD classification.

9.APPLICATIONS

The majority of AI use-cases and emerging applications for treating kidney disease appear to fall into two major categories:

- **Patient Monitoring and Prediction Models:** Companies are using machine learning to monitor patients and to predict and prevent the onset of kidney failure.
- **Medical Image Analysis:** Researchers are developing software using machine learning to analyze kidney biopsy images to help prevent disease.

10.CONCLUSIONS

Accurate prediction of chronic kidney disease is one of the emerging topics in medical diagnosis. Even though some approaches using real-time features shows very good performance in terms of accuracy. This work proposes a classification model to predict the chronic kidney disease using various machine learning algorithms.

All the four classification algorithms have been considered for diagnosis of chronic kidney disease. From the above results, the objective is to find the better model for chronic kidney disease. The K-Nearest Neighbour is the better model for diagnosis of chronic kidney disease it attains the accuracy of 98%. It correctly classified the 980 instances from 1000 instances. Thus finally it is observed that KNN is better algorithm for chronic kidney diagnosis.

ABSTRACT

11.FUTURE SCOPE

Early Prediction of Chronic Kidney Disease Using Machine Learning. Predictive analytics for healthcare using machine learning is a challenged task to help doctors decide the exact treatments for saving lives. In this paper, we present machine learning techniques for predicting the chronic kidney disease using clinical data. Four machine learning methods are explored including K-nearest neighbors (KNN), support vector machine (SVM), logistic regression (LR), and decision tree classifiers. These predictive models are constructed from chronic kidney disease dataset and the performance of these models are compared together in order to select the best classifier for predicting the chronic kidney disease. Currently, kidney disease is a major problem. Because there are so many people with this disease. Kidney disease is very dangerous if not immediately treated on time, and may be fatal. If the doctors have a good tool that can identify patients who are likely to have kidney disease in advance, they can heal the patients in time

12.REFERENCES/BIBLIOGRAPHY

- [1] Sinha, Parul, and Poonam Sinha. "Comparative study of chronic kidney disease prediction using KNN and SVM." International Journal of Engineering Research and Technology 4, no. 12 (2015): 608-12.
- [2] Yildirim, Pinar. "Chronic Kidney Disease Prediction on Imbalanced Data by Multilayer Perceptron: Chronic Kidney Disease Prediction." In Computer Software and Applications Conference (COMPSAC), 2017 IEEE 41st Annual, vol. 2, pp. 193-198. IEEE, 2017.
- [3] De Lusignan S, Chan T, Stevens P, O'Donoghue D, Hague N and Dzregah B, et al. "Identifying patients with chronic kidney disease from general practice computer records" ,Oxford Journals of Family Practice,Vol.3, Issue- 22, 2005, pp.234-241.
- [4] Hallan SI, Coresh J, Astor BC, Asberg A, Powe NR and Romundstad S, et al. "International comparison of the relationship of chronic kidney disease prevalence and ESRD risk", Journal American Society of Nephrology,Vol.17, Issue-8, 2006, pp.2275-2284.

ABSTRACT

[5] Levin A, Coresh J, Rossert J, et al. "Definition and classification of chronic kidney disease: a position statement from kidney disease", The New England Journal of Medicine, 2002, pp.36-42.

[6] Miguel A. Estudillo-Valderrama, Alejandro TalaminosBarroso and Laura M. Roa, "A Distributed Approach to Alarm Management in Chronic Kidney Disease", IEEE journal of biomedical and health informatics, Vol.18, Issue-6, 2014, pp. 1796-1803.

[7] Christopher A. Harle, Daniel B. Neill and Rema Padman, "Information Visualization for Chronic Disease Risk Assessment", IEEE Computer Society, 2012, pp.81-85. [8] Srinivasa R. Raghavan, Vladimir Ladik, and Klemens B.

[8] Mitri F.G. et al, "Vibro-acoustography imaging of kidney stones in vitro Vibro-acoustography", IEEE Transactions on Biomedical Engineering 2011.

[9] Chih-Yin Ho, Tun-Wen Pai, Yuan-Chi Peng and ChienHung Lee, "Ultrasonography Image Analysis for Detection and Classification of Chronic Kidney Disease", IEEE conference published on Intelligent and Software Intensive Systems (CISIS), 2012, pp.624 – 629.

[10] Al-Hyari and Al-Tae, "Clinical decision support system for diagnosis and management of Chronic Renal Failure", IEEE Jordan Conference on Applied Electrical Engineering and Computing Technologies (AEECT), 2013, pp.1-6.

[11] Kuo-Su Chen, Yung-Chih Chen and Yang-Ting Chen, "Stage diagnosis for Chronic Kidney Disease based on ultrasonography", IEEE International Conference on Fuzzy Systems and Knowledge Discovery (FSKD), 2014, pp. 525 – 530.

[12] Anne Rogers, Anne Kennedy, Thomas Blakeman and Christian Blickem, "Non-disclosure of chronic kidney disease in primary care and the limits of instrumental rationality in chronic illness self-management", ELSEVIER Social Science & Medicine 131, 2015, pp.31-39.

[13] Mohammed Shamim Rahman, Rajan Sharma and Stephen J.D.

ABSTRACT

Brecker, "Transcatheter aortic valve implantation in patients with pre-existing chronic kidney disease", ELSEVIER International Journal of Cardiology Heart & Vasculature , Vol.5, 2015, pp. 9–18.

[14] https://en.wikipedia.org/wiki/Data_mining.

[15] https://en.wikipedia.org/wiki/Feature_selection.

[16]
<https://machinelearningmastery.com/supervised-and-unsupervised-machine-learning-algorithms/>