# CHRONIC KIDNEY DISEASE PREDICTION BASED ON NAIVE BAYES TECHNIQUE

## Amogh Babu K A[1], Priyanka K[2], Raghavendra Babu T M[3]

[1]B.E. in Computer Science & Engineering, Mandya, Karnataka, India.
[2]Asst. Professor, Dept. of CSE, Nagarjuna College of Engineering & Technology, Bangalore, Karnataka, India.
[3]Asst.Professor, Dept. of CSE, P.E.S. College of Engineering, Mandya, Karnataka, India.

-------------------------------------------------------------------------***-------------------------------------------------------------------------

**Abstract –** *Data Mining has been a recent trend for obtaining a diagnostic result. Huge amount of unmined data is collected by the healthcare department in order to discover the hidden facts for effective diagnosis and also decision making. Data mining is defined as the process of extracting the huge hidden data from a large dataset, categorizing valid and unique patterns in data. There are lot of DM techniques like clustering, classification, association, analysis, regression etc. The main aim of this paper is to predict a YES or NO for Chronic Kidney Disease (CKD) using the classification technique i.e. Naïve Bayes.*

***Key Words***: Naive Bayes, Clustering, User Interface (UI), Data Mining (DM), Chronic Kidney Disease (CKD).

## 1. INTRODUCTION

Data Mining is one among the foremost encouraging areas of analysis with the aim of finding helpful information from voluminous knowledge of datasets. It's been employed in several domains like image mining, opinion mining, web mining, text mining, graph mining etc. Its applications embody anomaly detection, money knowledge analysis, medical knowledge analysis, social network analysis, marketing   research etc. It's become common in health   department   as there's a demand of   analytical methodology for predicting and finding unknown patterns and   obtaining   info   in   health   data.   It   plays a significant role for locating new trends in aid business.

Data Mining is especially helpful in medical field once no handiness of proof favouring a treatment choice is found. Great deal of advanced knowledge is being generated by aid business regarding patients, diseases, hospitals, medical equipments, claims,   treatment price etc. That needs process and analysis for information extraction. Data processing   comes   up   with a   group of   tools   and techniques that once applied to   the   present   processed knowledge, provides information to aid professionals for creating acceptable choices and enhancing the performance of patient management tasks. Patients with similar health problems is sorted along and   effective   treatment plans may well be recommended supported patient's history, physical examination, designation and former treatment patterns. Chronic Kidney Disease (CKD) has   become a world health   issue   and   is   a   locality   of   concern. It's a condition wherever kidneys become broken and can't filter nephrotoxic   wastes   within   the   body.   Our   work preponderantly focuses on police work life threatening diseases   like   chronic   nephrosis   (CKD) victimization Classification algorithms like Naive Bayes.

## 2. LITERATURE SURVEY

At present, health care industry is providing several benefits like fraud detection in health insurance, availability of medical facilities to patients at inexpensive prices, identification of smarter treatment methodologies, and construction of effective healthcare policies, effective hospital resource management, better customer relation, improved patient care and hospital infection control. Disease detection is also one of the significant areas of research in medical. Data mining approaches have become essential for healthcare industry in making decisions based on the analysis of the massive clinical data. Data mining is the process of extracting hidden information from massive dataset. Techniques like classification, clustering, regression and association have been used by in medical field to detect and predict disease progression and to make decision regarding patient's treatment. Classification is a supervised learning approach that assigns objects in a collection to target classes. It is the process which classifies the objects or data into groups, the members of which have one or more characteristic in common.  The techniques of classification are SVM, decision tree, Naive Bayes, ANN etc.
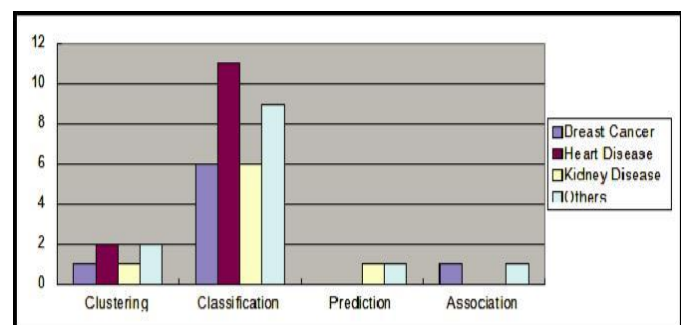


**Fig 2.1: potential use of data mining techniques**

The feasibleness study of employing a distributed approach for the management of alarms from chronic renal   disorder patients. The   key problems relating   to alarm   definition,   classification   and   prioritization consistent with on the market standardization efforts area

unit analyzed for the most situations addressed in dialysis. Then, the middleware projected for alarm management is represented, that follows the publish/subscribe pattern, and supports the OMG DDS (Data Distribution Service) customary. This customary facilitates the period of time observation of the changed info, furthermore because the quantifiability and ability of the answer developed relating to the various stakeholders and resources concerned [1].

The study was to work out the connection between the frequency spectrum of the irregular pulses associated with the stages of the CKD particularly from the chi space Information of the irregular pulse were classified into six stages i.e. 1, 2, 3a, 3b, 4, and five of the CKD patients. The information was collected by the activity pressure throughout beat periods or once blood vessels were in relaxed state. During this amount, the instrumentation records reflections of the heartbeat together with info concerning the amplitudes, frequencies, and pulse wave patterns. Observations were targeted on the part of the signal between amplitude from low to high on pulse patterns i.e. systolic period [3].

## 3. EXISITNG SYSTEM

Nowadays, health care industries are providing several benefits like fraud detection in health insurance, availability of medical facilities to patients at inexpensive prices, identification of smarter treatment methodologies, and construction of effective healthcare policies, effective hospital resource management, better customer relation, improved patient care and hospital infection control. Disease detection is also one of the significant areas of research in medical. There is no automation for chronic kidney disease prediction.

**Limitations of Existing System**

- ➢ Manual Approach
- ➢ Requires Medical Equipments
- ➢ More Expensive
- ➢ Lack of user satisfaction
- ➢ Less Efficient
- ➢ Less Accurate

## 4. PROPOSED SYSTEM

Our Aim is to predict the chronic kidney disease using the machine learning algorithm. Chronic kidney disease (CKD) means your kidneys are damaged and can't filter blood the way they should. The disease is called "chronic" because the damage to your kidneys happens slowly over a long period of time. This damage can cause wastes to build up in your body. CKD can also cause other health problems.10% of the population worldwide is affected by chronic kidney disease (CKD), and millions die each year because the doctors are unable diagnose the disease. The system is automation for predicting the CKD. The system is a Real-world web-based application that can be used by many

hospitals. Naive Bayes is a probabilistic classifier based on Bayes theorem. It assumes variables are independent of each other. The algorithm is easy to build and works well with huge data sets. It has been used because it makes use of small training data to estimate the parameters important for classification. It performs well in multiple class prediction. When assumption of independence holds a Naive Bayes classifier perform better compare to other models like logistic regression and you need less training data.

## 5. METHODOLOGIES

Data Mining is one of the most significant stages of the Knowledge Data Discovery process. The process involves data collection from various sources with pre-processing of the chosen data. The data is then transformed into suitable format for further processing. Data mining technique is applied on the data to extract valuable information and evaluation is done at the end.

### A. Data Collection

The clinical data of 400 records considered for analysis has been taken from UCI Machine Learning Repository. The data obtained after cleaning and removing missing values is 220. There are 25 attributes in the dataset. The numerical attributes include age, blood pressure, blood glucose random, blood urea, serum creatinine, sodium, potassium, hemoglobin, packaged cell volume, WBC count, RBC count. The nominal attributes include specific gravity, albumin and sugar. It also includes RBC, pus cell and pus cell clumps, bacteria, hypertension, diabetes mellitus, coronary artery disease, appetite, pedal edema, anemia and class.
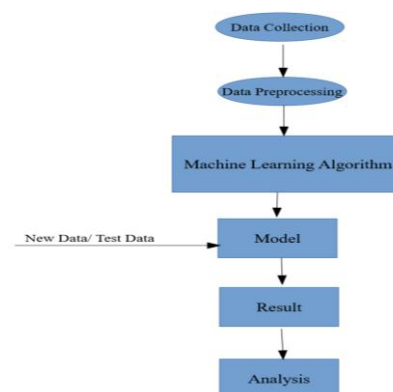
Number of Instances: 400

Number of Attributes: 25

Class: {CKD, NOTCKD}

Missing Attribute Values: yes

Class Distribution: [63% for CKD] [37% for NOTCKD]



**Fig 5.1: Stages of the Knowledge Data Discovery process.**

## B.  Data Pre-processing

Data pre-processing is a data mining technique that involves transforming raw data into an understandable format. The Real-world data is often incomplete, inconsistent, and/or lacking in certain behaviors or trends, and is likely to contain many errors. Data pre-processing is a proven method of resolving such issues. Data reprocessing prepares raw data for further processing. The data from UCI Repository is Pre-processed by removing the noisy data.

## C.  Machine Learning Algorithm

The machine learning algorithm a basic algorithm that we are using is the naïve Bayes algorithm to predict higher accuracy results and classification will always be accurate.

## D.  Dataset

| Attribute Name | Value Range | Description |
|---|---|---|
| age | 2, .., 90 | age |
| bp | 50, ...., 180 | blood pressure |
| sg | 1.005,1.010,1.015,1.020,1.025 | specific gravity |
| al | 0,1,2,3,4,5 | albumin |
| su | 0,1,2,3,4,5 | sugar |
| rbc | 2.1, ...., 8 | red blood cells |
| pc | normal,abnormal | pus cell |
| pcc | present,notpresent | pus cell clumps |
| ba | present,notpresent | bacteria |
| bgr | 22, ...., 490 | blood glucose random |
| bu | 1.5, ...., 391 | blood urea |
| sc | 0.4, ...., 76 | serum creatinine |
| sod | 4.5, ...., 163 | sodium |
| pot | 2.5, ...., 47 | potassium |
| hemo | 3.1, ..., 17.8 | hemoglobin |
| pcv | 9, ..., 54 | packed cell volume |
| wc | 2200,...., 26400 | white blood cell count |
| rc | 2.1,...., 8 | red blood cell count |
| htn | yes, no | hypertension |
| dm | yes, no | diabetes mellitus |
| cad | yes, no | coronary artery disease |
| appet | good,poor | appetite |
| pe | yes, no | pedal edema |
| ane | yes, no | anemia |
| class | ckd,notckd | class |

**Fig: Table Represents the used Data Set.**

## 6. WORKING OF THE SYSTEM

Coming to the working of the proposed system, our main aim is to predict the chronic disease but here we build a web-based application that can be used by hospitals. We have built an application that can be accessed by Admin, Receptionist, Doctor and even the patients. The admin is the person who maintains the entire application and the admin is responsible to add the new parameters or modify the existing parameters for the model. And the receptionist is responsible to add new patients' details and add the data for the new records of new patients and hence helps it to increase the dataset by adding new patient's data. If today we have 400 datasets to train our model next, we will dynamically increase the number of records in the training dataset to train our model; these new datasets of new patients are handled by the receptionist. Next, we have the main aim of our project that is to predict the chronic kidney disease using the naïve Bayes algorithm. Here, the model is

trained using the training dataset and the Naïve Bayes algorithm is executed as follows.



**Step 1:** Scan the dataset (storage servers)

retrieval of required data for mining from the servers such as database, cloud, excel sheet etc.

**Step 2:** Calculate the probability of each attribute value. [n, n_c, m, p]

Here for each attribute we calculate the probability of occurrence using the following formula. (mentioned in the next step). For each class(disease) we should apply the formulae.

**Step 3:** Apply the formulae

$$P(attributevalue(a_i)/subjectvaluev_j) = (n\_c + mp)/(n+m)$$

Where:

n = the number of training examples for which v = vj

n_c = number of examples for which v = vj and a = ai

p = a prior estimate for $P(a_i|v_j)$

m = the equivalent sample size

**Step 4:** Multiply the probabilities by p

for each class, here we multiple the results of each attribute with p and results are used for classification.

**Step 5:** Compare the values and classify the attribute values to one of the predefined sets of class.

**Fig 6.1: Implementation steps of naïve Bayes**

**Query Module-** We can add the query module as a future enhancement to the application where doctor, receptionist and admin of the application can interact with each other.

**Server Deployment-** We can deploy this onto the servers for online chronic kidney disease prediction and even create a wellness application for the users for curing or taking care of the disease.

We tested using other algorithms such as KNN (K-Nearest Neighbor Algorithm), SVM (Support Vector Machines), Decision tree and ANN (Artificial Neural Network) but surprisingly Naïve Bayes gave us amazing results and higher accurate results when compared to other algorithms. But using the J48 algorithm we can get similar accuracy rates like that of the Naïve Bayes itself. In future we can even test J48 algorithm to almost similar results. In future we can provide graphical analysis too, which is user friendly to understand.

**Fig 6.2: Explanation of Naïve Bayes algorithm**

## 7. EXPECTED RESULTS

It is successfully accomplished by applying the Naïve Bayes algorithm for classification. This classification technique comes under data mining technology.



**Fig 7.1: Home Page of the CKD Prediction**



**Fig 7.2: Admin Login Page of the CKD Prediction**



**Fig 7.3: Admin Dashboard**



**Fig 7.4: Admin Module- Parameter Addition**



**Fig 7.5: Receptionist Login Page of CKD Prediction**

**Fig 7.6: Training Dataset submission from Receptionist**



**Fig 7.7: Sample View of Training Dataset**



**Fig 7.8: Doctor Login Page of CKD Prediction**



**Fig 7.9: Uploading Testing Dataset from Doctor**



**Fig 7.10: Single Patient CKD Prediction from Doctor**



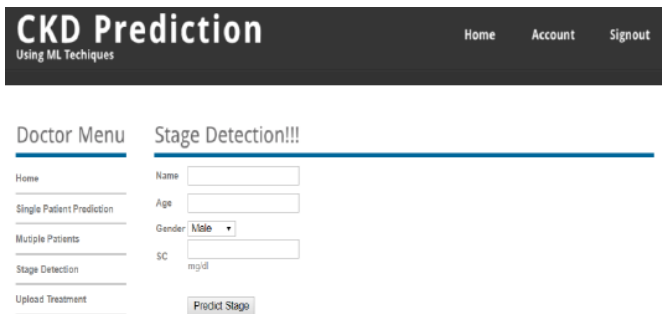**Fig 7.11: Multiple Patient CKD Prediction UI**

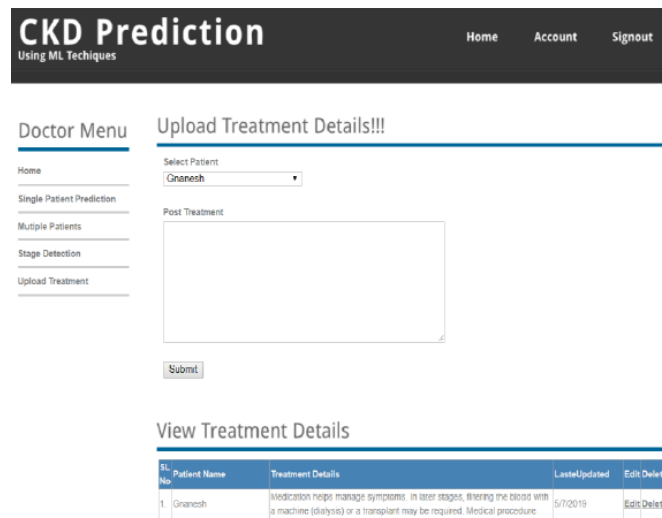**Fig 7.12: Doctor Module – Stage Prediction**



**Fig 7.13: Doctor Module – Upload the treatment Details**


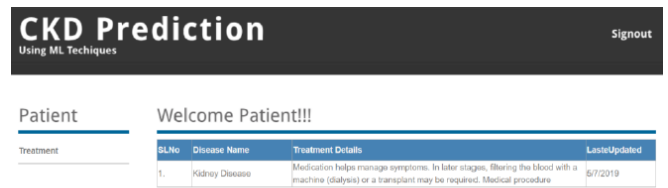
**Fig 7.14: Patient Login Page**



**Fig 7.15: Patient Module- View Treatment Details**

## 8. CONCLUSION

This project is a medical sector application which helps the medical practitioners in predicting the CKD based on the CKD parameters. It is automation for CKD disease prediction and it efficiently and economically speedily identifies the disease, its types and complications from the clinical database. The Accuracy obtained is about 94.6%.
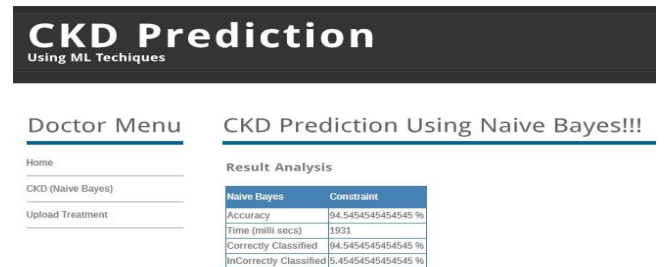


**Fig 8.1: Result analysis**

## 9. FUTURE WORK

We can enhance this problem statement by implementing the below features -

- ➢ Graphical analysis
- ➢ Feature Extraction
- ➢ Stage Prediction

## REFERENCES

[1] Miguel A. Estudillo-Valderrama, Alejandro Talaminos-Barroso, Laura M. Roa, Fellow, IEEE, David Naranjo-Hern´andez, Student Member, IEEE, Javier Reina-Tosina, Senior Member, IEEE, Nuria Arest´e-Fosalba, and Jos´e A. Milan-Martin "A Distributed Approach to Alarm Management in Chronic Kidney Disease", IEEE Journal of Biomedical and Health informatics, VOL. 18, NO.6, November 2014.

[2] ArifahFashaRosmani, UmiHanim Mazlan, Alif Faisal Ibrahim, Dina Shamila Zakaria "iKS:Composition of Chronic Kidney Disease (CKD) Online Informational Self-Care Tool", 2015 IEEE 2015 International Conference on Computer, Communication, and Control Technology (I4CT 2015), April 21 - 23 in Imperial kuching hotel, Kuching, Sarawak, Malaysia.

[3] ErniYudaningtyas, Djoko H. Santjojo, WaruDjuriatno, IndraznoSiradjuddin, Muhammad Rony Hidayatullah,

"Identification of Pulse Frequency Spectrum of Chronic Kidney Disease Patients Measured at TCM Points Using FFT Processing".

[4] Renuka Marutirao Pujari and Mr. Vikas D. Hajare, "Analysis of Ultrasound Images for Identification of Chronic Kidney Disease Stages", 978-14799-3486-7/14 ©2014 IEEE.

[5] Veenita Kunwar, Khushboo Chandel and A. Sai Sabitha "Chronic Kidney Analysis using Data Mining Classification Techniques" 2016 6th InternationalConference - Cloud System and Big Data Engineering (Confluence).

[6]Pinar Yildirim"Chronic Kidney Disease Prediction on Imbalanced Data by Multilayer Perceptron" 2017 IEEE 41st Annual Computer Software and Applications Conference.

[7]. I. H. Witten and E. Frank, "Data Mining Practical Machine Learning Tools and Techniques," 2nd ed., San Francisco/ABD, 2005.

[8].https://archive.ics.uci.edu/ml/datasets/Chronic_Kidney_ Disease (Access Date: 2018 February 7).

[9]. http://www.tbv.com.tr/tr/content/main/page/p/164-kronik-bobrekhastaliginin-    erken-teshisi-ve-korunma-yontemleri (Access Date: 2018 February 7)

[10]. C. Cortes and V. Vapnik, "Support- Vector Networks," Machine Learning, vol. 20, 1995, pp. 273-297.

[11].http://www.datascience.istanbul/2017/07/02/hata-matrisini-confusionmatrix- yorumlama/ (Access Date: 2018 February 7)

**BIOGRAPHIES**

**AMOGH BABU K A**

B.E. in Computer Science and Engg.

http://amoghbabu.xyz

**PRIYANKA K**

Asst. Professor, Dept. of CSE, Nagarjuna College of Engineering and Technology, Bangalore.

**RAGHAVENDRA BABU T M**

Asst. Professor, Dept. of CS&E, P.E.S. College of Engineering, Mandya.