

KIDNEY DISEASES PREDICTION USING MACHINE LEARNING

*Minor project-II report submitted
in partial fulfillment of the requirement for award of the degree of*

**Bachelor of Technology
in
Computer Science & Engineering**

By

K. MOHITH SAI	(21UEDS0038)	(VTU 19278)
THANGA VISHNU VARDHAN REDDY	(21UEDS0062)	(VTU 19696)
KALAVAGUNTA SAKETH	(21UEDS0029)	(VTU 19277)

Under the guidance of
Dr. C. KOTTESWARAN, M.E, PhD
ASSOCIATE PROFESSOR



**DEPARTMENT OF COMPUTER SCIENCE & ENGINEERING
SCHOOL OF COMPUTING**

**VEL TECH RANGARAJAN DR. SAGUNTHALA R&D INSTITUTE OF
SCIENCE & TECHNOLOGY**

(Deemed to be University Estd u/s 3 of UGC Act, 1956)

**Accredited by NAAC with A++ Grade
CHENNAI 600 062, TAMILNADU, INDIA**

May, 2024

KIDNEY DISEASES PREDICTION USING MACHINE LEARNING

*Minor project-II report submitted
in partial fulfillment of the requirement for award of the degree of*

**Bachelor of Technology
in
Computer Science & Engineering**

By

K. MOHITH SAI	(21UEDS0038)	(VTU 19278)
THANGA VISHNU VARDHAN REDDY	(21UEDS0062)	(VTU 19696)
KALAVAGUNTA SAKETH	(21UEDS0029)	(VTU 19277)

*Under the guidance of
Dr. C. KOTTESWARAN, M.E,PhD
ASSOCIATE PROFESSOR*



**DEPARTMENT OF COMPUTER SCIENCE & ENGINEERING
SCHOOL OF COMPUTING**

**VEL TECH RANGARAJAN DR. SAGUNTHALA R&D INSTITUTE OF
SCIENCE & TECHNOLOGY**

(Deemed to be University Estd u/s 3 of UGC Act, 1956)

**Accredited by NAAC with A++ Grade
CHENNAI 600 062, TAMILNADU, INDIA**

May, 2024

CERTIFICATE

It is certified that the work contained in the project report titled "KIDNEY DISEASES PREDICTION USING MACHINE LEARNING " by "K. MOHITH SAI (21UEDS0038) ,THANGA VISHNU VARDHAN REDDY (21UEDS0062), KALAVAGUNTA SAKETH (21UEDS0029)," has been carried out under my supervision and that this work has not been submitted elsewhere for a degree.

Signature of Supervisor

Computer Science & Engineering

School of Computing

Vel Tech Rangarajan Dr. Sagunthala R&D

Institute of Science & Technology

May, 2024

Signature of Professor In-charge

Computer Science & Engineering

School of Computing

Vel Tech Rangarajan Dr. Sagunthala R&D

Institute of Science & Technology

May, 2024

DECLARATION

We declare that this written submission represents our ideas in our own words and where others' ideas or words have been included, we have adequately cited and referenced the original sources. We also declare that we have adhered to all principles of academic honesty and integrity and have not misrepresented or fabricated or falsified any idea/data/fact/source in our submission. We understand that any violation of the above will be cause for disciplinary action by the Institute and can also evoke penal action from the sources which have thus not been properly cited or from whom proper permission has not been taken when needed.

K. MOHITH SAI

Date: / /

THANGA VISHNU VARDHAN REDDY

Date: / /

KALAVAGUNTA SAKETH

Date: / /

APPROVAL SHEET

This project report entitled KIDNEY DISEASES PREDECTION USING MACHINE LEARNING by K. MOHITH SAI (21UEDS0038), THANGA VISHNU VAEDHAN REDDY (21UEDS0062), KALAVGUNTA SAKETH (21UEDS0029) is approved for the degree of B.Tech in Computer Science & Engineering.

Examiners

Supervisor

Dr. C. KOTTESWARAN, M.E,PhD,

Date: / /

Place:

ACKNOWLEDGEMENT

We express our deepest gratitude to our respected **Founder Chancellor and President Col. Prof. Dr. R. RANGARAJAN B.E. (EEE), B.E. (MECH), M.S (AUTO),D.Sc., Foundress President Dr. R. SAGUNTHALA RANGARAJAN M.B.B.S.** Chairperson Managing Trustee and Vice President.

We are very much grateful to our beloved **Vice Chancellor Prof. S. SALIVAHANAN**, for providing us with an environment to complete our project successfully.

We record indebtedness to our **Professor & Dean, Department of Computer Science & Engineering, School of Computing, Dr. V. SRINIVASA RAO, M.Tech., Ph.D.**, for immense care and encouragement towards us throughout the course of this project.

We are thankful to our **Head, Department of Computer Science & Engineering, Dr.M.S. MURALI DHAR, M.E., Ph.D.**, for providing immense support in all our endeavors.

We also take this opportunity to express a deep sense of gratitude to our **Internal Supervisor Dr.C.KOTTESWARAN, M.E.**, for his cordial support, valuable information and guidance, he helped us in completing this project through various stages.

A special thanks to our **Project Coordinators Mr. V. ASHOK KUMAR, M.Tech., Ms. C. SHYAMALA KUMARI, M.E., Mr. SHARAD SHANDHI RAVI, M.Tech.**, for their valuable guidance and support throughout the course of the project.

We thank our department faculty, supporting staff and friends for their help and guidance to complete this project.

K. MOHITH SAI	(21UEDS0038)
THANGA VISHNU VAEDHAN REDDY	(21UEDS0062)
KALAVGUNTA SAKETH	(21UEDS0029)

ABSTRACT

Chronic Kidney Disease is a serious lifelong condition that induced by either kidney pathology or reduced kidney functions. According to the survey in 2021-2022 for every 10 in 3, people are suffering from this chronic kidney disease. It enables us to introduce the optimal subset of parameters to feed machine learning to build a set of predictive models. K-Nearest Neighbor, Random Forest, LBGM and Decision Tree methods are applied for prediction. Among these four methodologies, the proposed model suggests an LBGM is suitable for the early prediction of this kind of disease. Performance measures show that Light Gradient Boosted Machine (LGBM) gives 99.98 procedure concludes that advances in machine learning and predictive analytics, represent a promising model to recognize intelligent solutions, which in turn prove the ability of predication in kidney disease.

Keywords: Chronic Kidney Disease, Decision Tree, K-Nearest Neighbor, Light Gradient Boosted Machine, Machine Learning, Prediction, Random Forest.

LIST OF FIGURES

4.1	CKD-Architecture Diagram	12
4.2	CKD-Data flow Diagram	13
4.3	CKD-Usecase Diagram	14
4.4	CKD-Class Diagram	15
4.5	CKD-Sequence Diagram	16
4.6	CKD-Activity Diagram	17
4.7	CKD Data Set	19
4.8	CKD Data Preprocessing	20
4.9	Code development	21
5.1	Input Data	25
5.2	Output Data	26
5.3	Test result of Unit Testing	28
5.4	Integration testing result	29
5.5	Test Image	31
6.1	Random Forest accuracy	33
6.2	Logistic Regression accuracy	33
6.3	Serum Creatinine Density	36
6.4	Blood pressure Density	37
6.5	White Blood Cell Count Density	37
9.1	Poster	47

LIST OF ACRONYMS AND ABBREVIATIONS

SL.NO	ABBREVIATION	DEFINITION
1	AI	Artificial Intelligence
2	CKD	Chronic Kidney Disease
3	CNN	Convolutional Neural Network
4	DT	Decision Tree
5	ESRD	End Stage Renal Disease
6	FP	False Positive
7	FN	False Negative
8	GBF	Gradient Boosting Framework
9	GFR	Glamorise Filtration Rate
10	KNN	K-Nearest Neighbor
11	LGBM	Light Gradient Boosted Machine
12	RF	Random Forest
13	SVM	Support Vector Machine
14	TP	True Positive
15	TN	True Negative
16	ML	Machine Learning
17	UCI	Unique Client Identifier

TABLE OF CONTENTS

	Page.No
ABSTRACT	v
LIST OF FIGURES	vi
LIST OF ACRONYMS AND ABBREVIATIONS	vii
1 INTRODUCTION	1
1.1 Introduction	1
1.2 Aim of the project	2
1.3 Project Domain	2
1.4 Scope of the Project	3
2 LITERATURE REVIEW	4
3 PROJECT DESCRIPTION	9
3.1 Existing System	9
3.2 Proposed System	9
3.3 Feasibility Study	10
3.3.1 Economic Feasibility	10
3.3.2 Technical Feasibility	10
3.3.3 Social Feasibility	10
3.4 System Specification	11
3.4.1 Hardware Specification	11
3.4.2 Software Specification	11
3.4.3 Standards and Policies	11
4 METHODOLOGY	12
4.1 General Architecture	12
4.2 Design Phase	13
4.2.1 Data Flow Diagram	13
4.2.2 Use Case Diagram	14
4.2.3 Class Diagram	15

4.2.4	Sequence Diagram	16
4.2.5	Activity Diagram	17
4.3	Algorithm & Pseudo Code	17
4.3.1	Algorithm	17
4.3.2	Pseudo Code	18
4.4	Module Description	19
4.4.1	Data Collection and Training	19
4.4.2	Graphical User Interface for the Application	21
4.4.3	Clustering and Classification	22
4.5	Steps to execute/run/implement the project	23
4.5.1	Data Collection	23
4.5.2	Model Selection	23
4.5.3	Model Evaluation and deployment	24
5	IMPLEMENTATION AND TESTING	25
5.1	Input and Output	25
5.1.1	Input Design	25
5.1.2	Output Design	26
5.2	Testing	26
5.3	Types of Testing	27
5.3.1	Unit testing	27
5.3.2	Integration testing	28
5.3.3	System testing	29
5.3.4	Test Result	31
6	RESULTS AND DISCUSSIONS	32
6.1	Efficiency of the Proposed System	32
6.2	Comparison of Existing and Proposed System	33
6.3	Sample Code	34
7	CONCLUSION AND FUTURE ENHANCEMENTS	38
7.1	Conclusion	38
7.2	Future Enhancements	38
8	PLAGIARISM REPORT	39

9	SOURCE CODE & POSTER PRESENTATION	40
9.1	Source Code	40
9.2	Poster Presentation	47
	References	47

Chapter 1

INTRODUCTION

1.1 Introduction

Chronic kidney disease (CKD) is a significant public health problem for world-wide, especially for a low and medium-income countries. CKD means that the kidney does not work and cannot correctly filter the blood. About 10 percent of the population for a worldwide suffering from (CKD), and millions of die each year because of they couldn't get the affordable treatment, with the number increasing in the elderly.

In 2021, a study was conducted by International Society of Numerology(ISN) on global burden disease, they reported that CKD has been raised an important cause of the mortality worldwide with the number of deaths increasing by 82.3 percent in the last two decades. Also, the number of patients reaching End-Stage Renal Disease (ESRD) is increasing, which requires kidney transplantation or dialysis to save the patient's lives.

The worst possible outcome of the chronic kidney disease and the symptoms causing any reduced kidney functioning would lead to kidney failure. When the symptoms are become severe and uncontrollable they can be treated through the dialysis and transplantation. Early detection and treatment of CKD can slow or stop the progression of the kidney disease. But the CKD, in it's early stages, has no symptoms.

The Proper diagnosis and the testing may be the only way to find out whether the patient has been affected by kidney disease. Early detection of CKD in its initial stages can help the patient get the effective treatment and then prohibit the progression to ESRD. CKD is a Long term condition induced the damage done to both kidneys. Kidney damage refers to any kind of kidney pathology that gives the possibility to reduce the capacity of kidney functions, particularly the reduction in Glamorise Filtration Rate (GFR).

Chronic kidney disease (CKD) is a progressive condition that affects the kidneys, which are responsible for filtering waste products and excess fluids from the blood. CKD can be caused by a range of factors including high blood pressure, diabetes, autoimmune diseases, and inherited genetic conditions. As CKD progresses, the kidneys become less able to perform their filtering function, leading to a buildup of waste products and fluids in the body. This can cause a range of symptoms including fatigue, weakness, nausea, and swelling in the legs and feet.

1.2 Aim of the project

The main aim is to identify whether a particular patient is affected by CKD (or) not and it has to be accurate and precise. So, for that we are going to purpose a correlate four pre-existing Machine Learning Algorithms to find the best among all. For this purpose, we gathered a CKD data set from UCI machine learning repository and we examining the correlation between the development of the CKD and predictors using a predictive approach of the analysis. This will help us to reduce the number of required parameters to predict the CKD disease occurrence as well as eliminating the missing, redundant and noisy data. To use certain features to measure its accuracy and predictions.

1.3 Project Domain

Machine learning (ML) is a subfield of artificial intelligence (AI) that focuses on the development of algorithms and models that can learn and make predictions or decisions based on patterns found in data. In the context of the project Kidney Disease Prediction, machine learning techniques is used to analyze large datasets of patient records, and other relevant information to identify patterns that may indicate the presence of Kidney Disease or the likelihood of developing the disease. The algorithms such as K-Nearest Neighbor, Random Forest, Decision Tree, Light Gradient Boosted Machine are used to train models on labeled data to classify patients as either having kidney disease or not. The framework utilizes a machine learning algorithm called Light Gradient Boosted Machine to analyze patient data and make predictions about their risk of developing kidney disease.

Light Gradient Boosted Machine is a powerful algorithm that has been used successfully in several medical domains. It works by building multiple decision trees and combining their predictions to obtain a more accurate result. The project domain includes data collection, processing, and analysis. The framework will be trained on a large dataset of medical images, patient history, and demographic information to learn the patterns and characteristics of kidney disease.

1.4 Scope of the Project

The scope of the project is to Predict whether the patient is suffering from kidney diseases or not by using machine learning algorithms. They can minimize the risk caused by the diseases because “prevention is better than cure”.

Chapter 2

LITERATURE REVIEW

[1] Imesh Udara Ekanayake,et.al.,(2020) , have proposed that the workflow to predict CKD status based on clinical data, incorporating data preprocessing, a missing value handling method with collaborative filtering and attributes selection. The data distribution has properly covered the whole domain in CKD, but the general attributes like appetite, anaemia and pedal oedema are biased towards CKD. It easy to achieve an accurate prediction using the KNN algorithm in machine learning with the data set. It highlights the importance of incorporating the domain knowledge into feature selection when analysing clinical data related to CKD. In the general context, to achieve the accuracy may lead to false positives as observed in Recall.

[2] Yedilkhan Amirgaliyev,et.al.,(2019) , have identified that there are less studies performed in CKD by using SVMs. But there are several studies where SVM has been used abundantly. They used SVMs as a classification model for detecting and diagnosing malignant and benign tumors based on MRI features, ultrasound feature and mammographic features. This result tries to create and test SVM classifier over CKD dataset. To improve the accuracy of CKD classification as positive and negative, the performance of Support Vector Machine (SVM) was evaluated. The SVM is a flexible classifier algorithm that has been proposed as an effective statistical learning method for pattern recognition, which is based on finding optimal hyperplane to separate different classes mapping input data into higher-dimensional feature space. It is automatic classification algorithm for kidney disease diagnosis based on clinical history, physical examinations, and laboratory tests, which are noninvasive, cheap and safe.

[3] Khalid Twarish Alhamazani,et.al.,(2021) , represented the CRISP-DM methodology to the context of the problem so that the different logically organized stages were taken, data collection, pre-processing, learning, evaluation, and selection, which allowed the construction of a model capable of classifying the possibility of a diagnosis of CKD with an accuracy of 93 forests algorithm has obtained quite

optimal results, where predictions of 93 been obtained. The innovation of this work results from the design adjusted to the environment of the health system in Iraq and the pathology of the CKD in our country, with a methodology adapted to the case study and a production architecture proposal for the model with Microsoft Azure tools of form that allows satisfying in the future the scalability of the solution.

[4] Pankaj Chittora,et.al., (2021) , have developed a model to predict CKD disease in patients. The performance of the model was tested on both all attributes and selected features. Among feature selection methods there were Wrapper, Filter and Embedded allowing to select vital features. Classifier algorithms performance was tested on the selected features. The machine learning classifiers such as logistic regression and K- nearest neighbors (KNN) were used for training the model. Logistic and KNN classifiers give satisfactory result and have the negligible difference between precision and recall values. Logistic and KNN algorithms have not enough capacity to distinguish between positive class and negative class as the related score is very low.

[5] Jing Xiao,et.al., (2019) , have proposed the work to predict the severity of CKD using more easily available demographic and blood biochemical features during follow-up. Developed and compared several predictive models using statistical, machine learning and neural network approaches. A Web tool (CKD Prediction System) for clinical practice that can be widely used in the evaluation of proteinuria progress in nephrology has developed. The comparison in this study shows that the linear models performed better. Journal of Translational Medicine Predicting the severity of CKD using machine learning tools. The dialysis machine mixes and monitors the dialysate. Dialysate is the fluid that helps remove the unwanted waste. The incremental feature selection approach with Convolutional Neural Network (CNN) is applied to anticipate disease.

[6] Janakiram Ampolu,et.al., (2019) , represented the causes , symptoms, and complications of CKD in a clear manner such that even common people can easily understand. Once a patient is detected and proved to be affected by CKD then the patient as well as the care takers, including doctors, must follow some constraints. There by it is possible to prevent CKD progression in the patient. Modern methods are needed to prevent the pathogens which are responsible for CKD. With the help

of various engineering techniques one can easily design controllers to assess as well as to prevent CKD permanently. The easiest procedure for identifying CKD is to screen people. Current recommendations suggest screening of individuals with diabetes, hypertension, cardiovascular disease, and family history of kidney diseases in the course of routine health checkups. Much work has been done in medical sciences in the area of CKD, but there is still scope for further research. From the recent studies, advanced tools such as data mining, etc., are considered to be the current trend in the area of CKD.

[7] Salma Shaji,et.al., (2020) , have proposed Chronic Kidney Disease also recognized as Chronic Renal Disease, is an uncharacteristic functioning of kidney or a failure of renal function expanding over a period of months or years. Habitually, chronic kidney disease is detected during the screening of people who are known to be in threat by kidney problems, such as those with high blood pressure or diabetes and those with a blood relative Chronic Kidney Disease(CKD) patients. So the early prediction is necessary in combating the disease and to provide good treatment. This study proposes the use of machine learning techniques for CKD such as Ant Colony Optimization(ACO) technique and Support Vector Machine(SVM) classifier. Final output predicts whether the person is having CKD or not by using minimum number of features. CKD is caused due to diabetes and high blood pressure. Due to Diabetes our many organs get affected and it will be followed by high blood sugar. So it is important to predict the disease as early as possible. This study improvises some of the machine learning techniques to predict the disease.

[8] Revathy Ramesh,et.al., (2020) , proposes a datamining framework for knowledge discovery on the CKD datasets. Large amounts of CKD datasets are collected. Data preperation and preprocessing is done using the traditional methods of data mining process. Three machine learning algorithms namely Decision tree, Random Forest and Support Vector machines are used to predict the early occurence of CKD. The goodness of each algorithm is analysed.The biosciences have advanced to a larger extent and have generated large amounts of information from Electronic Health Records. This have given rise to the acute need of knowledge generation from this enormous amount of data. Data mining methods and machine learning play a major role in this aspect of biosciences. Chronic Kidney Disease(CKD) is a condition in which the kidneys are damaged and cannot filter blood as they always

do. A family history of kidney diseases or failure, high blood pressure, type 2 diabetes may lead to CKD. This is a lasting damage to the kidney and chances of getting worser by time is high. The very common complications that results due to a kidney failure are heart diseases, anemia, bone diseases, high potasium and calcium. The worst case situation leads to complete kidney failure and necessitates kidney transplant to live. An early detection of CKD can improve the quality of life to a greater extent. This calls for good prediction algorithm to predict CKD at an earlier stage . Literature shows a wide range of machine learning algorithms employed for the prediction of CKD. It uses data preprocessing,data transformation and various classifiers to predict CKD and also proposes best Prediction framework for CKD. The results of the framework show promising results of better prediction at an early stage of CKD.

[9] B.Sasi Varna,et.al., (2021) , have proposed several machine learning models in this section, the component models that performed better when diagnosing the data samples were chosen for inclusion. Examination of the component models' mistakes determined the component models' roles as potential biomarkers. Chronic kidney disease (CKD) is a global health issue that causes a high incidence of morbidity and death, as well as the onset of additional illnesses. Because there are no clear symptoms in the early stages of CKD, people frequently miss it. Early identification of CKD allows patients to obtain prompt therapy to slow the disease's development. Due of their rapid and precise identification capabilities, machine learning models can successfully assist doctors in achieving this aim.The CKD data set was collected from the machine learning repository at the University of California, Irvine (UCI). As a result, it will determine whether or not a patient has CKD and, if so, whether or not further drugs should be taken. Six machine learning algorithms (Logistic Regression, AdaBoost, Random Forest, Decision Tree, and Gradient Boosting) were used to establish models.

[10] Dibaba Adeba Debal,et.al., (2022) , represented with both binary and multi classfication for stage prediction have been carried out. The prediction models used include Random Forest (RF), Support Vector Machine (SVM) and Decision Tree (DT). Analysis of variance and recursive feature elimination using cross validation have been applied for feature selection. Evaluation of the models was done using tenfold cross-validation. The results from the experiments indicated that RF based

on recursive feature elimination with cross validation has better performance than SVM and DT. Predictive analysis using machine learning techniques can be helpful through an early detection of CKD for efficient and timely interventions. In this study, Random Forest (RF), Support Vector Machine (SVM) and Decision Tree (DT) have been used to detect CKD. Most of previous researches focused on two classes, which make treatment recommendations difficult because the type of treatment to be given is based on the severity of CKD.

Chapter 3

PROJECT DESCRIPTION

3.1 Existing System

In machine learning, everyone tried to go with the best algorithm which will try to give a more accuracy and play a prominent role in prediction. It uses machine learning algorithms, such as random forests and gradient boosting. It discusses various machine learning approaches used for predicting chronic kidney disease, including decision trees, logistic regression, artificial neural networks, and support vector machines. It provides a comprehensive overview of the existing systems. To predict the outcomes of kidney transplants based on donor and recipient characteristics, Which are available to us, and by comparison, we might be able to notice it's compatibility and it's versatility in the model as well as the system. So, it is always better to look for a alternative approach. The advantages like Early detection and Prevention,Better treatment. Predicting kidney disease is a cost-effective way of identifying individuals who are at risk of developing the disease. Early detection and treatment can help prevent costly and invasive treatments in the future. The prediction of kidney disease can help to raise awareness of the disease and promote healthy lifestyles that can help reduce the risk of developing the disease. This can ultimately lead to improved public health outcomes

3.2 Proposed System

It is better to search for an alternative approach with promising results and we are going to do that. The proposed model compares four very commonly known machine learning algorithms and are going to state the best among them using a case study on Chronic Kidney disease Prediction. As stated, the four machine learning algorithms are: Random Forest, Decision Tree, LGBM, and K- Near Neighbour(KNN). The disadvantages of Chronic kidney disease often has no symptoms in the early stages,

and symptoms may not appear until the disease is advanced. This makes it difficult to diagnose and treat early. CKD can lead to other serious health problems, such as heart disease, high blood pressure, anemia, nerve damage, and bone disease. The advantage is early detection of kidney disease can help healthcare professionals to develop a personalized treatment plan that is tailored to the individual patient's needs. This can help improve treatment outcomes and reduce the risk of complications and by identifying the risk factors associated with kidney disease, individuals can take preventative measures to reduce their risk of developing the disease. These measures may include lifestyle changes such as adopting a healthy diet and regular exercise.

3.3 Feasibility Study

3.3.1 Economic Feasibility

The kidney disease prediction using machine learning project can be economically feasible if the potential benefits outweigh the costs. Proper planning and assessment of the project costs and benefits can help ensure its economic viability.

3.3.2 Technical Feasibility

The kidney disease prediction using machine learning project can be technically feasible if there is access to high-quality data, appropriate machine learning algorithms, powerful computing resources, reliable deployment infrastructure, and ongoing monitoring and maintenance.

3.3.3 Social Feasibility

The social feasibility of a kidney disease prediction using machine learning project depends on the acceptance and adoption of the technology by patients, healthcare providers, and other relevant parties. The project must comply with ethical standards and regulations and be designed to be accessible, equitable, and aligned with the goals of the healthcare system.

3.4 System Specification

3.4.1 Hardware Specification

- Processor : Intel(R) Core(TM) i3-10110U CPU @ 2.10GHz 2.60 GHz
- RAM : 8.00 GB (7.78 GB usable)
- Hard Disc Space : 8GB

3.4.2 Software Specification

- Operating System : Windows 7,8,10 (or) Mac (or) Linux
- Coding Language : Python - version(3.5)
- Platform Requirement: Google Colab
- Framework: Kara's (or) Tensor flow (or) Kaggle

3.4.3 Standards and Policies

General Data Protection Regulation (GDPR):

GDPR is a regulation of the European Union that came into effect in May 2018. It sets out strict guidelines for the processing of personal data, including health data. If you are collecting data from patients in the EU, you must comply with GDPR.

Intellectual Property Rights (IPR):

An IRB is an independent committee that reviews and approves research studies involving human subjects. IRBs ensure that the study is conducted ethically and that the privacy and rights of the patients are protected.

Chapter 4

METHODOLOGY

4.1 General Architecture

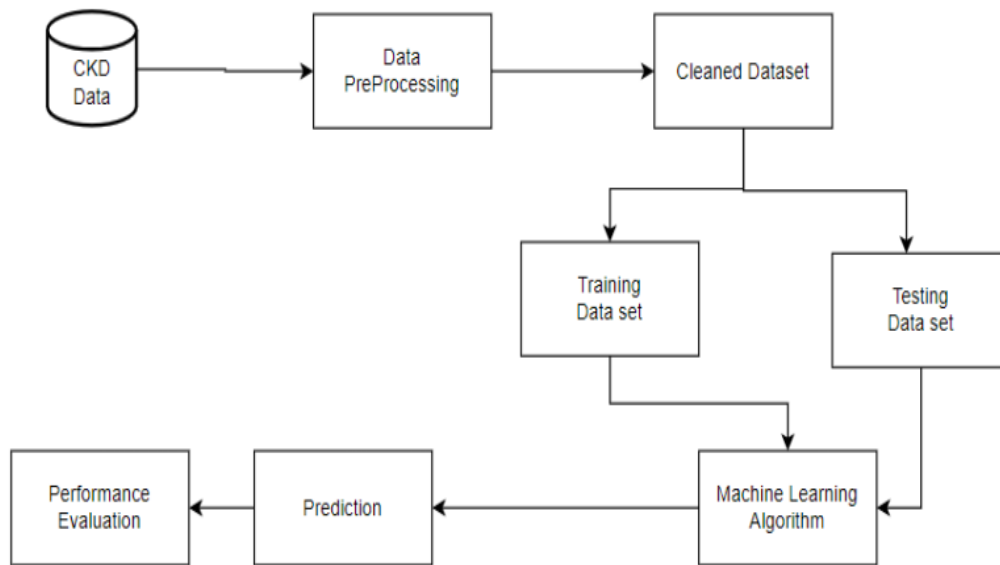


Figure 4.1: CKD-Architecture Diagram

In figure 4.1 show the architecture of the kidney disease prediction. Initially the data set is collected from the University of California, Irvine (UCI). Then that data set is preprocessed to remove the null values and repeated values. Then the output is cleaned dataset then that particular dataset is divided into two different types like training dataset and testing dataset it under goes to the machine learning algorithms like K-Nearest Neighbor, Random Forest, Decision Tree and Light Gradient Boosted Machine. By using these four algorithms the disease is predicted and performance evaluation is done.

4.2 Design Phase

4.2.1 Data Flow Diagram

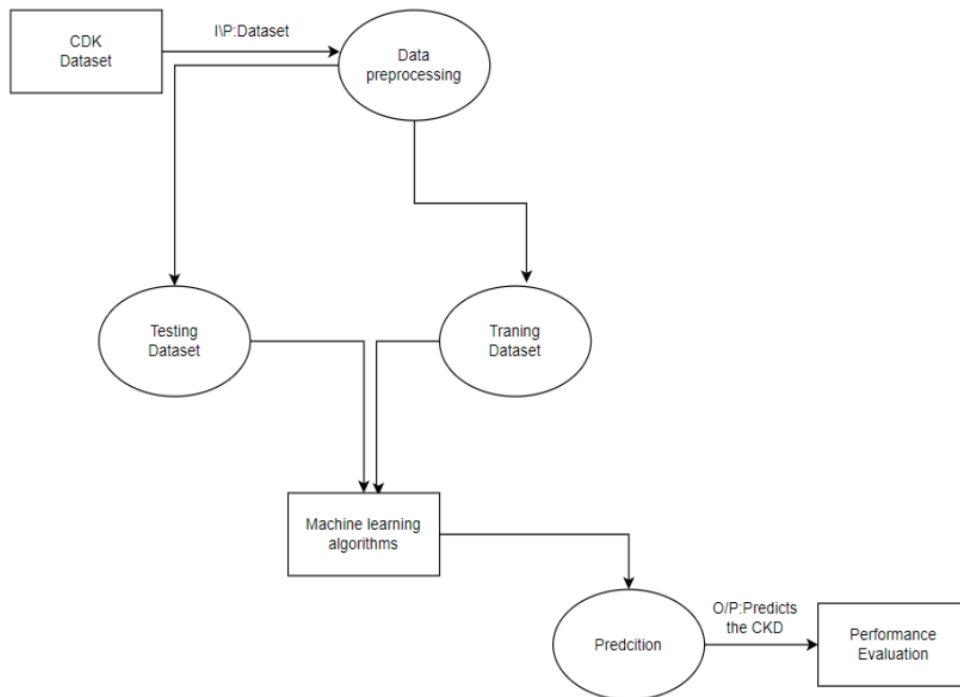


Figure 4.2: CKD-Data flow Diagram

In Figure4.2 Data flow diagram for chronic kidney disease classification using machine learning algorithms. The training dataset and testing dataset are used in the machine learning algorithms for predicting the CKD. Based on the performance evaluation, the kidney disease can be predicted. It has experimented on UCI dataset which contains early stages of 400 CKD patients with 25 attributes.

4.2.2 Use Case Diagram

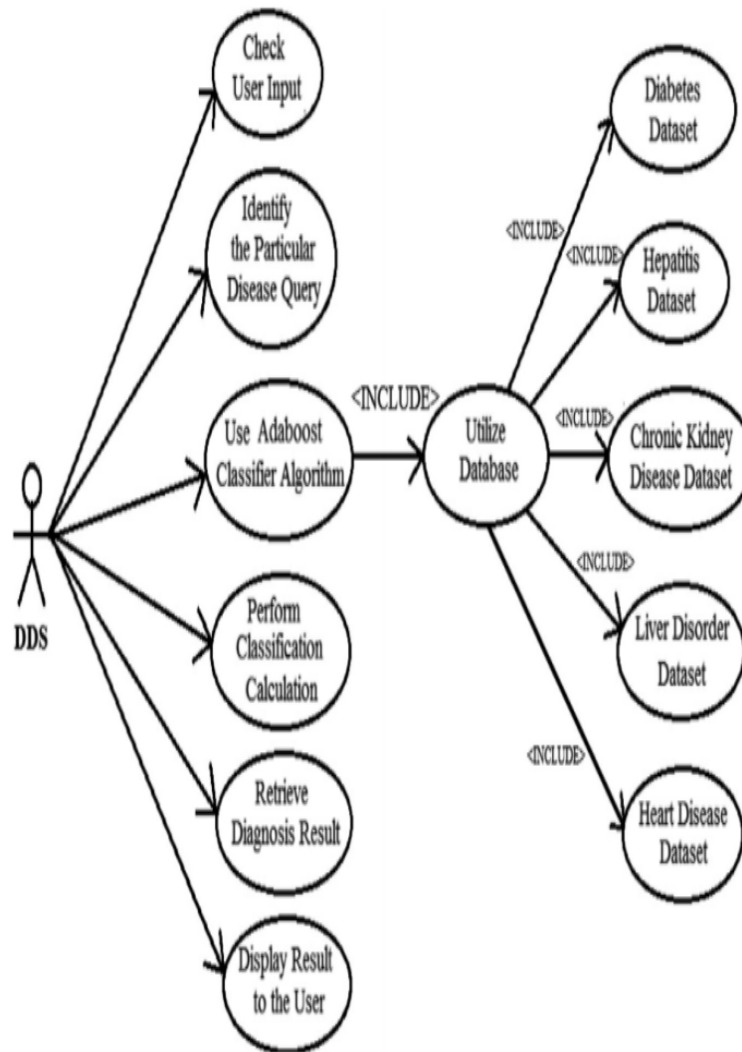


Figure 4.3: CKD-Usecase Diagram

In figure 4.3 Use case diagram is a type of behavioral diagram defined by and created from a Usecase analysis. Its purpose is to present a graphical overview of the functionality provided by a system in terms of actors, their goals (represented as use cases), and any dependencies between those use cases. The main purpose of a use case diagram is to show what systems functions are performed for which actor. The user interacts with the system by providing patient data and viewing the prediction results. The system performs the data preprocessing, feature selection, machine learning model training, and prediction. This use case diagram provides an overview of the functionalities of a kidney disease prediction system using machine learning and how users can interact with it.

4.2.3 Class Diagram

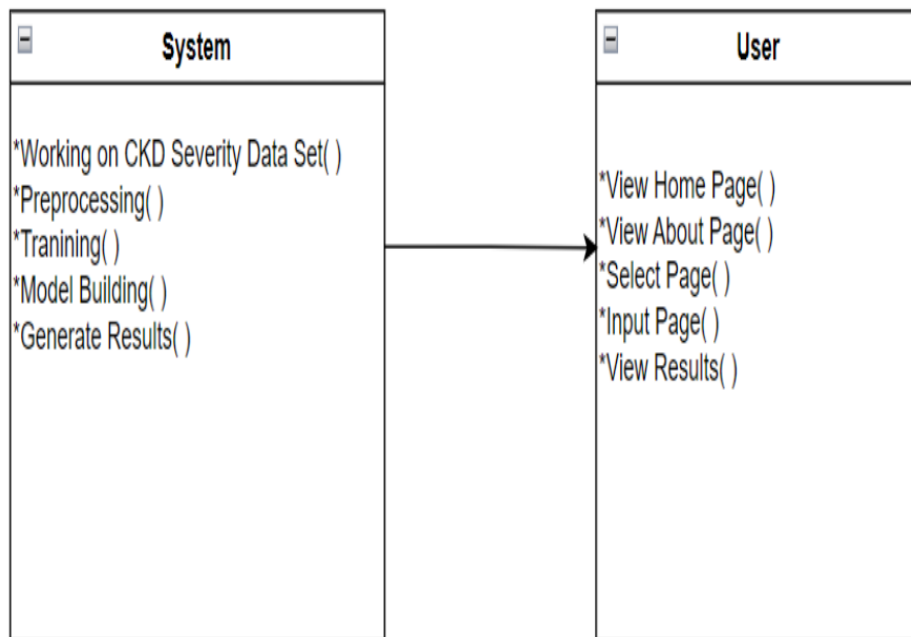


Figure 4.4: CKD-Class Diagram

In figure 4.4 Class Diagram is a type of static structure diagram that describes the structure of a system by showing the system's classes, their attributes, operations (or methods), and The relationships between these classes can be shown using associations, where one class is connected to another class. For example, the User class can have an association with the Prediction class to represent that the user can view the prediction results. Similarly, the Patient class can have an association with the Data Preprocessing class, Feature Selection class, Machine Learning Model class, and Prediction class to represent the flow of data and functionalities. This class diagram provides a static structure of a kidney disease prediction system using machine learning and shows how the different classes are related to each other.

4.2.4 Sequence Diagram

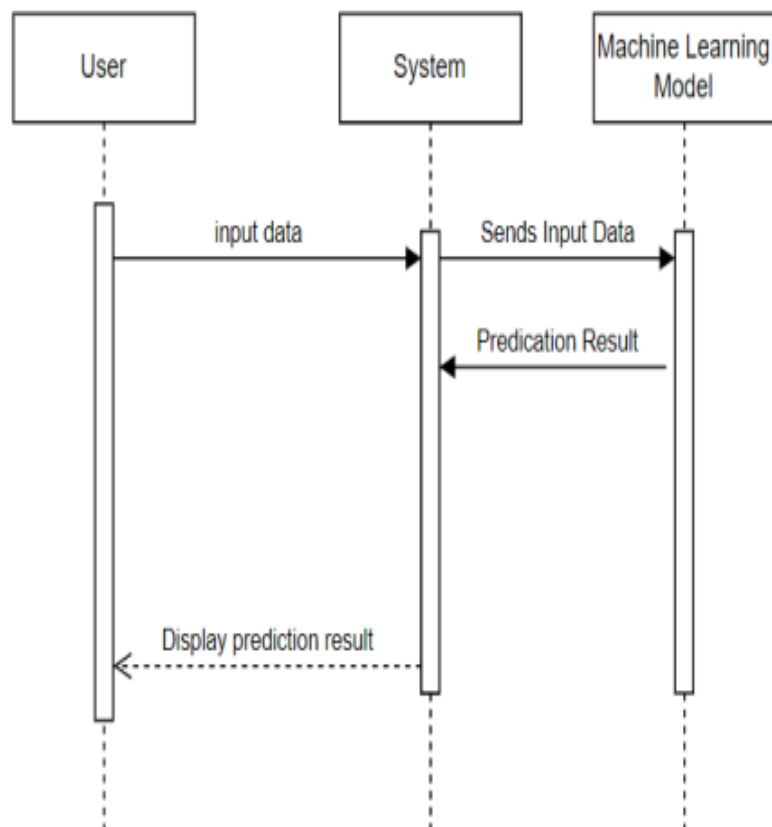


Figure 4.5: CKD-Sequence Diagram

In figure 4.5 Sequence Diagram is a kind of interaction diagram that shows how processes operate with one another and in what order. Sequence diagrams are sometimes called event diagrams, event scenarios, and timing diagrams. The sequence diagram shows the flow of messages between the objects or components involved in the system. The user interacts with the system by providing patient data and viewing the prediction results. The system performs the data preprocessing, feature selection, machine learning model training, and prediction. The sequence diagram shows the order in which these interactions occur. This sequence diagram provides a dynamic view of a kidney disease prediction system using machine learning and shows how the different components interact with each other during the prediction process.

4.2.5 Activity Diagram

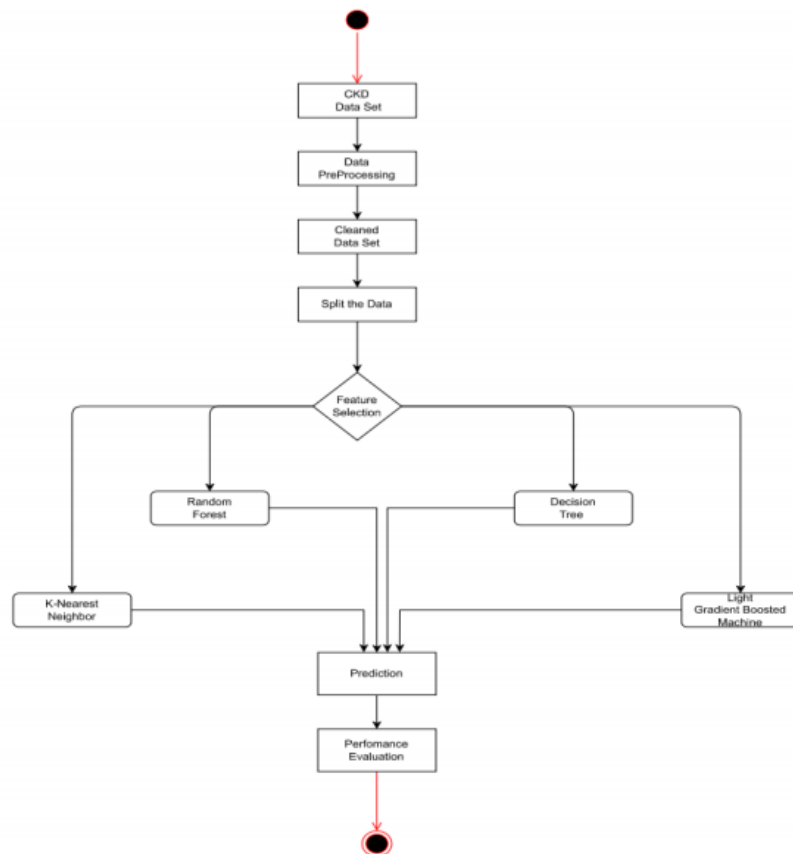


Figure 4.6: CKD-Activity Diagram

In figure 4.6 graphical representations of workflows of stepwise activities and actions with support for choice, iteration and concurrency. In the Unified Modelling Language, activity diagrams can be used to describe the business and operational stepby-step workflows of components in a system. An activity diagram shows the overall flow of control.

4.3 Algorithm & Pseudo Code

4.3.1 Algorithm

1. Import the LGBMClassifier from the lightgbm library
 - from lightgbm import LGBMClassifier
2. Initialize the LGBMClassifier model with a learning rate of 1
 - lgbm = LGBMClassifier(learning rate = 1)

3. Train the model on the training data (x train and y train)
 - `lgbm.fit(x train, y train)`
4. Evaluate the performance of the model on the test data (x test and y test)
 - `lgbm acc = accuracy score(y test,lgbm.predict(x test))`
5. Print the training accuracy, test accuracy, confusion matrix, and classification report of the LGBMClassifier model
 - print the Training Accuracy of LGBM Classifier is `accuracy score(y train, lgbm.predict(x train))`
 - print the Test Accuracy of LGBM Classifier is `lgbm acc`
 - print the confusion matrix(`y test, lgbm.predict(x test)`)
 - print the classification report(`y test, lgbm.predict(x test)`)
6. Create a pandas DataFrame to compare the performance of different models
 - `models = pd.DataFrame('Model' : ['KNN', 'Decision Tree Classifier', 'LGBM Classifier'], 'Score' : [knn acc, dtc acc, lgbm acc])`
7. Sort the models DataFrame by the 'Score' column in descending order
 - `models.sort values(by = 'Score', ascending = False)`

4.3.2 Pseudo Code

```
1
2 # Import necessary libraries
3 import pandas as pd
4 import numpy as np
5 from sklearn.model_selection import train_test_split
6 from sklearn.preprocessing import StandardScaler
7 from sklearn.linear_model import LogisticRegression
8 from sklearn.metrics import accuracy_score
9
10 # Load the dataset
11 data = pd.read_csv('kidneydisease.csv')
12
13 # Data cleaning and preprocessing
14 # ...
15
16 # Split the data into training and testing sets
17 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
18
19 # Feature scaling
20 sc = StandardScaler()
21 X_train = sc.fit_transform(X_train)
22 X_test = sc.transform(X_test)
23
```

```

24 # Train the logistic regression model
25 model = LogisticRegression(random_state=42)
26 model.fit(X_train, y_train)
27
28 # Make predictions on the testing set
29 y_pred = model.predict(X_test)
30
31 # Evaluate the model's accuracy
32 accuracy = accuracy_score(y_test, y_pred)
33 print('Accuracy:', accuracy)

```

4.4 Module Description

4.4.1 Data Collection and Training

1. Data Collection: The Dataset of prediction of chronic kidney disease using machine learning algorithm is downloaded from UCI repository. In that dataset there are 400 patient records are included. Also they include 25 attributes but we take only 14 attributes for building model. Age, Blood pressure, Albumin, Red blood cells, Pus cell, Pus cells clumps, Serum creatinine, Haemoglobin, White blood cell count, Red blood cell count, Anaemia, Classification, Appetite, Packed cell volume all this 14 attributes are used to build model.

FileHomeInsertPage LayoutFormulasDataReviewViewHelp

Calibri11A

BBIU

Font

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard

Clipboard</

Figure 4.7: CKD Data Set

2. Preprocessing: Data Cleaning: Gather open source raw data of CKD patients available on internet. Data obtained from internet does not contains the name of the attribute so first we assigned the names to the attribute. Missing values in the dataset like NA's or blank values are removed by using WEKA function “ReplaceMissingValues” used, which replaces NA's with the mean values of that attribute. Data Reduction: Out of 25 attributes present in the dataset, we have selected 14 important attributes required to build predictive model.

```

import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import plotly.express as px
from sklearn.preprocessing import OneHotEncoder
from sklearn.compose import ColumnTransformer

from sklearn.preprocessing import LabelEncoder
from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestClassifier, AdaBoostClassifier
from sklearn.linear_model import LogisticRegression
from sklearn.neighbors import KNeighborsClassifier
from sklearn.naive_bayes import MultinomialNB
from sklearn.tree import DecisionTreeClassifier
from sklearn.metrics import accuracy_score, classification_report, confusion_matrix, precision_score, f1_score, recall_score
from sklearn.metrics import roc_curve, roc_auc_score
from sklearn.metrics import precision_recall_curve, auc
from sklearn.model_selection import GridSearchCV
from sklearn.datasets import load_iris
from sklearn.datasets import fetch_20newsgroups
from sklearn.feature_extraction.text import TfidfVectorizer

import warnings
warnings.filterwarnings('ignore')
%matplotlib inline
pd.set_option('display.max_columns', 26)

df = pd.read_csv('kidney_disease.csv')
df.head()

```

	id	age	bp	sg	al	sa	rhc	pc	prc	ha	hgr	hs	sc	sod	pot	hemo	pcv	wc	rc	hta	ds	cad	appet	pe	ane	classification
0	48.0	80.0	1020	1.0	0.0	NaN	normal	notpresent	notpresent	121.0	36.0	1.2	NaN	NaN	15.4	44	7800	5.2	yes	no	good	no	no	no	ckd	
1	7.0	50.0	1020	4.0	0.0	NaN	normal	notpresent	notpresent	NaN	16.0	0.8	NaN	NaN	11.3	38	6000	NaN	no	no	no	good	no	no	no	ckd
2	62.0	80.0	1010	2.0	3.0	normal	normal	notpresent	notpresent	423.0	53.0	1.8	NaN	NaN	9.6	31	7500	NaN	no	yes	no	poor	no	yes	ckd	
3	48.0	70.0	1005	4.0	0.0	normal	abnormal	present	notpresent	117.0	56.0	3.8	111.0	2.5	11.2	32	6700	3.9	yes	no	no	poor	yes	yes	ckd	

Figure 4.8: CKD Data Preprocessing

3. Machine Learning Model: The Model Development is a developing the machine learning model for kidney disease prediction requires careful consideration of the data, features, and model selection. It's important to ensure the model is accurate, reliable, and scalable for real-world applications and we apply the four types of machine learning algorithms like K-NN, RF, LGBM, Decision Tree. Finding the more accuracys out of those four algorithms.

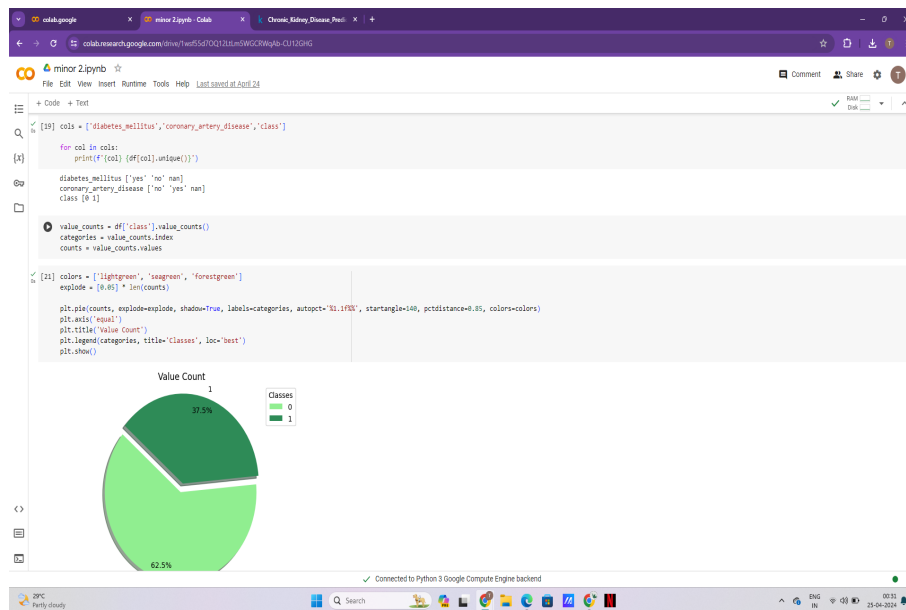


Figure 4.9: Code development

4.4.2 Graphical User Interface for the Application

1. Input Fields:

Age: Textbox or dropdown menu for entering the patient's age.

Gender: Radio buttons or dropdown menu for selecting male/female.

Blood Pressure: Textbox for systolic and diastolic blood pressure readings.

Blood Glucose: Textbox for entering blood glucose levels.

Serum Creatinine: Textbox for entering serum creatinine levels.

Albuminuria: Textbox for entering albuminuria levels.

Family History: Checkbox for indicating family history of CKD.

2. **Predict Button:** Button to trigger the prediction process based on the input data.

3. **Output Display:** Text area or label to display the predicted CKD risk or classification (e.g., low risk, moderate risk, high risk).

4. **Reset Button:** Button to clear all input fields and reset the interface.

5. **Feedback Mechanism:** Provide users with a way to provide feedback or report issues with the application.

6. **Responsive Design:** Ensure that the GUI is responsive and works well on different devices and screen sizes.

7. **Branding:** Incorporate branding elements such as logos, colors, and fonts to maintain consistency with your application or organization's branding.

4.4.3 Clustering and Classification

GLCM Algorithm

The Feature Extraction is a method of capturing visual content of images for indexing retrieval. Grey Level Co-occurrence Matrix (GLCM) method is a way of extracting second order statistical texture features. The GLCM functions characterize the texture of an image by calculating with specific values and in a specified spatial relationship occur in an image, creating a GLCM, and then extracting statistical measures from this matrix. The image segmentation is the process of partitioning a digital image into multiple segments (sets of pixels, also known as image objects). The goal of segmentation is to simplify and/or change the representation of an image into something that is more meaningful and easier to analyze.

K Means clustering

The k-means algorithm is a clustering algorithm that can be used for plant disease detection. The algorithm works by grouping similar data points into clusters, with each cluster being represented by its centroid. Initialize the system and input data which involves setting up the detection system and capturing images of the plants. K-means clustering is the extracted features are then used as input to the k-means algorithm. The algorithm groups the features into a specified number of clusters (k), with each cluster being represented by its centroid. The number of clusters is usually determined by trial and error, or by using techniques such as the elbow method.

Feature extraction:

The next step is to extract relevant features from the images. This could involve techniques such as color-based feature extraction or texture analysis. Classify images based on cluster labels, Once the clustering is complete, each image can be classified based on the cluster it belongs to. If the cluster contains images of healthy plants, then any new image that belongs to that cluster can be classified as healthy. Conversely, if the cluster contains images of diseased plants, then any new image that belongs to that cluster can be classified as diseased.

4.5 Steps to execute/run/implement the project

4.5.1 Data Collection

- **Patient Demographics:** Basic information about the patient, including age, gender, ethnicity, occupation, and lifestyle habits like smoking and alcohol consumption, which can help in understanding potential risk factors.
- **Medical History:** Past medical conditions and family history can provide insights into the patient's overall health and predisposition to certain diseases, including CKD. Conditions like hypertension, diabetes, and cardiovascular diseases are known risk factors for CKD.
- **Labeling Images:** Each image needs to be labeled with information such as the plant species, the disease present (if any), and the severity of the disease. This labeling process helps in training the machine learning model.
- **Imaging Studies:** Renal ultrasound and other imaging modalities can provide visual information about the size, shape, and structure of the kidneys, aiding in the diagnosis and monitoring of CKD.
- **Data Source:** Data can be sourced from electronic health records, laboratory databases, patient interviews, or remote monitoring devices, depending on availability and accessibility.
- **Data Preprocessing:** Before building predictive models, it's essential to clean and preprocess the data, handling missing values, outliers, and encoding categorical variables appropriately.

4.5.2 Model Selection

- **Understand the Data:** Begin by thoroughly understanding your dataset. Identify the features available for prediction and their relevance to CKD. Ensure the dataset is clean and properly preprocessed.
- **Define Performance Metrics:** Decide on the evaluation metrics based on the nature of your problem. For CKD prediction, common metrics include accuracy, precision, recall, F1-score, and area under the ROC curve (AUC-ROC).
- **Baseline Models:** Start with simple baseline models to establish a benchmark performance. These could include logistic regression, decision trees, or naive Bayes classifiers.
- **Regularization Techniques:** Consider using regularization techniques to prevent

overfitting in complex models. Common regularization methods include L1 and L2 regularization for linear models and dropout for neural networks.

- **Cross-Validation:** Employ cross-validation techniques such as k-fold cross-validation to assess the generalization performance of your models.
- **Model Evaluation:** Evaluate the performance of each model using the chosen metrics on both the training and validation datasets. Ensure that the selected model generalizes well to unseen data.

4.5.3 Model Evaluation and deployment

- **Data Collection and Preprocessing:** Collect relevant data on patients' demographics, medical history, laboratory tests, and other pertinent information related to CKD. Preprocess the data to handle missing values, outliers, and ensure consistency and quality.
- **Feature Selection and Engineering:** Identify the most relevant features for predicting CKD. Engineer new features if necessary to improve the predictive power of the model.
- **Model Training:** Select an appropriate machine learning algorithm (e.g., logistic regression, decision trees, random forests, support vector machines) or deep learning model.
- **Model Evaluation:** Evaluate the model's performance using appropriate metrics such as accuracy, precision, recall, F1-score, and area under the receiver operating characteristic curve (ROC-AUC).
- **Monitoring and Maintenance:** Implement monitoring systems to track the model's performance in real-time. Regularly update the model as new data becomes available or as the underlying patterns change.
- **Documentation:** Document the entire process, including data collection, preprocessing steps, model architecture, evaluation results, and deployment procedures.

Chapter 5

IMPLEMENTATION AND TESTING

5.1 Input and Output

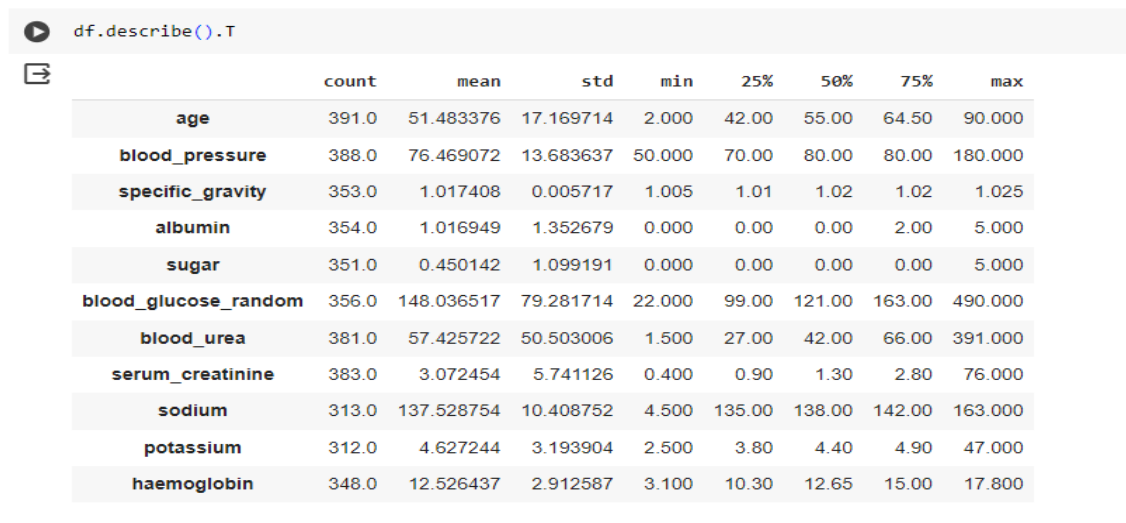
5.1.1 Input Design

J	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z	AA
1	id	age	sex	bp	sg	al	su	dbc	pcr	pot	lba	hgr	lbu	sc	lod	pcr	hemo	pcr	lvs	tc	lba	cho	cad	appet	pe	ana	classification
2	1	40	80	1.02	1	0	normal	notpresent/notpresent	121	36	1.2					15.4	44	7800	5.2	yes	yes	no	good	no	no	no	did
3	1	7	90	1.02	4	0	normal	notpresent/notpresent		18	9.8					11.3	36	6000		no	no	no	good	no	no	no	did
4	2	82	80	1.01	2	3	normal	normal	notpresent/notpresent	425	59	1.8				9.6	31	7500		no	yes	no	poor	yes	yes	did	
5	3	40	70	1.025	4	0	normal	abnormal	present	notpresent	117	56	3.8	111	2.5	12.2	32	6700	1.9	yes	no	no	poor	yes	yes	did	
6	4	51	80	1.01	2	0	normal	normal	notpresent/notpresent	306	26	1.4				11.8	35	7800	4.6	no	no	no	good	no	no	no	did
7	5	60	90	1.015	3	0		notpresent/notpresent	74	25	1.1	141	3.2	12.2	29	7800	4.4	yes	yes	yes	no	good	yes	no	no	did	
8	6	68	70	1.01	8	0	normal	notpresent/notpresent	330	54	24	104	4	12.4	36		no	no	no	no	good	no	no	no	no	did	
9	7	24		1.015	2	4	normal	abnormal	notpresent/notpresent	410	11	1.1				12.4	44	6900	5	no	yes	no	good	yes	no	no	did
10	8	52	100	1.015	3	0	normal	abnormal	present	notpresent	136	60	1.9			10.8	33	9600	4	yes	yes	no	good	no	yes	did	
11	9	53	90	1.02	2	0	abnormal	abnormal	present	notpresent	70	107	7.2	114	3.7	9.5	29	11300	3.7	yes	yes	no	poor	no	yes	did	
12	10	56	60	1.01	2	4	abnormal	abnormal	present	notpresent	490	55	4			9.4	28		yes	yes	no	good	no	yes	did		
13	11	63	70	1.01	3	0	abnormal	abnormal	present	notpresent	180	60	2.7	111	4.2	10.8	32	4500	1.8	yes	yes	no	poor	yes	no	did	
14	12	68	70	1.015	3	1	normal	present	notpresent	208	72	2.1	138	3.8	9.7	38	11300	1.8	yes	yes	yes	poor	yes	no	did		
15	13	68	70					notpresent/notpresent	90	86	4.6	135	3.4	9.8			yes	yes	yes	yes	poor	yes	no	no	did		
16	14	68	80	1.01	3	2	normal	abnormal	present	present	157	90	4.1	136	6.4	5.6	36	11000	1.6	yes	yes	yes	poor	yes	no	did	
17	15	40	80	1.015	3	0	normal	notpresent/notpresent	70	162	3.6	141	4.9	7.6	24	3800	1.6	yes	no	no	no	good	no	yes	did		
18	16	47	70	1.015	2	0	normal	notpresent/notpresent	90	86	2.2	138	4.1	12.6			no	no	no	no	good	no	no	no	no	did	
19	17	47	80					notpresent/notpresent	114	87	5.2	139	3.7	12.1			yes	no	no	no	poor	no	no	no	no	did	
20	18	66	100	1.025	8	3	normal	notpresent/notpresent	360	27	1.3	135	4.3	12.7	37	11400	4.3	yes	yes	yes	yes	good	no	no	no	did	
21	19	82	60	1.015	1	0	abnormal	abnormal	present	notpresent	300	11	1.6			10.3	36	5300	1.7	yes	no	yes	good	no	no	did	
22	20	81	80	1.015	2	0	abnormal	abnormal	notpresent/notpresent	170	148	3.8	131	3.2	7.7	24	5300	1.2	yes	yes	yes	poor	yes	yes	did		
23	21	60	90					notpresent/notpresent	180	76	4.1			10.9	32	6300	1.6	yes	yes	yes	yes	good	no	no	no	did	
24	22	48	80	1.025	4	0	normal	abnormal	notpresent/notpresent	95	169	7.7	136	3.8	9.8	32	6900	1.4	yes	no	no	good	no	yes	did		
25	23	21	70	1.01	8	0	normal		notpresent/notpresent								no	no	no	no	poor	no	yes	did			
26	24	42	100	1.015	4	0	normal	abnormal	notpresent/present		50	1.4	129	4	11.1	29	8300	4.8	yes	no	no	poor	no	no	no	did	
27	25	61	90	1.025	8	0	normal	notpresent/notpresent	336	75	1.9	141	5.2	9.9	29	8400	1.7	yes	yes	no	good	no	yes	did			
28	26	75	80	1.015	8	0	normal	notpresent/notpresent	156	45	2.4	140	3.4	11.6	35	10300	4	yes	yes	no	poor	no	no	no	did		
29	27	65	70	1.01	3	4	normal	abnormal	notpresent/notpresent	264	87	2.7	130	4	12.5	37	9600	4.1	yes	yes	yes	good	yes	no	no	did	
30	28	75	70		1	3		notpresent/notpresent	120	11	1.4						no	yes	no	no	good	no	no	no	no	did	
31	29	68	70	1.005	1	0	abnormal	abnormal	present	notpresent	26	1.4				12.9	38		no	no	yes	good	no	no	no	did	

Figure 5.1: Input Data

Figure 5.1 Here we are using the 2nd approach because this is an easy way and we can also increase the instances along with it. Currently, we are having 400 instances and a predictive model will work just fine with more instances.

5.1.2 Output Design



The image shows a Jupyter Notebook interface. At the top, there is a code cell with the command `df.describe().T`. Below it, the output is displayed as a table. The table has 10 columns: **age**, **blood_pressure**, **specific_gravity**, **albumin**, **sugar**, **blood_glucose_random**, **blood_urea**, **serum_creatinine**, **sodium**, **potassium**, and **haemoglobin**. Each column contains statistical data: count, mean, std, min, 25%, 50%, 75%, and max.

	count	mean	std	min	25%	50%	75%	max
age	391.0	51.483376	17.169714	2.000	42.00	55.00	64.50	90.000
blood_pressure	388.0	76.469072	13.683637	50.000	70.00	80.00	80.00	180.000
specific_gravity	353.0	1.017408	0.005717	1.005	1.01	1.02	1.02	1.025
albumin	354.0	1.016949	1.352679	0.000	0.00	0.00	2.00	5.000
sugar	351.0	0.450142	1.099191	0.000	0.00	0.00	0.00	5.000
blood_glucose_random	356.0	148.036517	79.281714	22.000	99.00	121.00	163.00	490.000
blood_urea	381.0	57.425722	50.503006	1.500	27.00	42.00	66.00	391.000
serum_creatinine	383.0	3.072454	5.741126	0.400	0.90	1.30	2.80	76.000
sodium	313.0	137.528754	10.408752	4.500	135.00	138.00	142.00	163.000
potassium	312.0	4.627244	3.193904	2.500	3.80	4.40	4.90	47.000
haemoglobin	348.0	12.526437	2.912587	3.100	10.30	12.65	15.00	17.800

Figure 5.2: Output Data

Figure 5.2 The output of this predictive is represented numerically ranging between 0 to 1 because we are using logistic regression here and this initiative gives the output in the range 0 and 1. And this is common for every model, the output will be in the range of 0 and 1.

5.2 Testing

A strategy in this testing is integrating the systems test cases along with the design techniques into a well planned series of steps with the results in the successful construction of software. This testing strategy should work along with the test planning, test execution, test case design, and the data collection resultant and evaluation. To accommodate a low-level test is a strategy for software testing. Software testing is a critical element of software quality assurance which represents the ultimate review of specification design and coding. Interesting anomalies from the software are being to represented by using Testing. Thus, the period of testing is performed for the proposed system and there by is ready for user acceptance testing.

5.3 Types of Testing

5.3.1 Unit testing

Chronic kidney disease is a medical condition characterized by the gradual loss of kidney function over time. It can have various causes and stages, and its management often involves medical interventions, lifestyle changes, and sometimes software tools for tracking and managing patient data.

Input

```
1 models = {
2     "Random Forest": RandomForestClassifier(),
3
4 }
5
6 for model_name, model in models.items():
7     model.fit(X_train, y_train)
8
9     # Make predictions on the test set
10    y_pred = model.predict(X_test)
11
12    # Evaluate the model
13    accuracy = accuracy_score(y_test, y_pred)
14    print(f'{model_name} Accuracy: {accuracy:.2f}')
15
16    # Generate a classification report
17    print(f'{model_name} Classification Report:')
18    print(classification_report(y_test, y_pred))
19
20    cm = confusion_matrix(y_test, y_pred)
21    print(f'{model_name} Confusion Matrix:')
22    print(cm)
23
24    print('-' * 55)
```

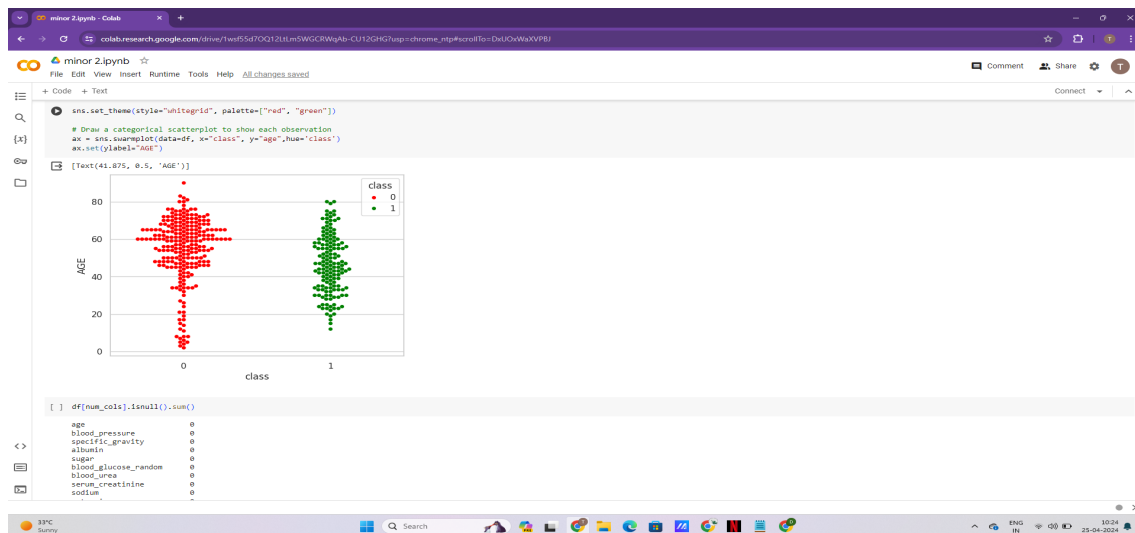


Figure 5.3: Test result of Unit Testing

Test result : PASSED

5.3.2 Integration testing

Integration testing for chronic kidney disease (CKD) involves examining how different components of CKD management interact to ensure the entire system works effectively. CKD management typically involves various healthcare professionals, medical interventions, medications, and lifestyle modifications.

Input

```

1 df.sample(10)
2 df.info()
3 df.describe().T
4 df['packed_cell_volume'] = pd.to_numeric(df['packed_cell_volume'], errors='coerce')
5 df['white_blood_cell_count'] = pd.to_numeric(df['white_blood_cell_count'], errors='coerce')
6 df['red_blood_cell_count'] = pd.to_numeric(df['red_blood_cell_count'], errors='coerce')
7 cat_cols = [col for col in df.columns if df[col].dtype == 'object']
8 cat_cols
9 num_cols = [col for col in df.columns if df[col].dtype != 'object']
10 num_cols
11 for col in cat_cols:
12     print(f'{col} {df[col].unique()}')
13 df['diabetes_mellitus'].replace({'\tno': 'no', '\tyes': 'yes', ' yes': 'yes'}, inplace=True)
14 df['coronary_artery_disease'] = df['coronary_artery_disease'].str.replace('\tno', 'no')
15 df['class'].replace({'ckd\t': 'ckd', 'notckd': 'not ckd'}, inplace=True)
16 df['class'] = df['class'].map({'ckd': 0, 'not ckd': 1})
17 df['class'] = pd.to_numeric(df['class'], errors='coerce')
18 cols = ['diabetes_mellitus', 'coronary_artery_disease', 'class']

```



```

19
20 for col in cols:
21     print(f'{col} {df[col].unique()}')
22 value_counts = df['class'].value_counts()
23 categories = value_counts.index
24 counts = value_counts.values
25 colors = ['lightgreen', 'seagreen', 'forestgreen']
26 explode = [0.05] * len(counts)

```

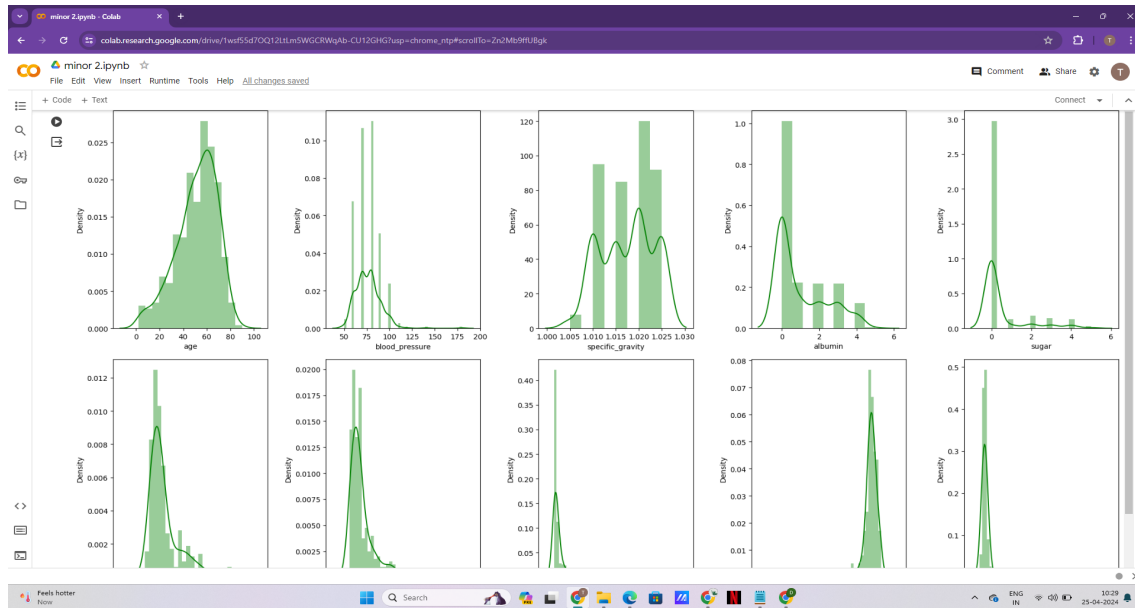


Figure 5.4: Integration testing result

Test result : PASSED

5.3.3 System testing

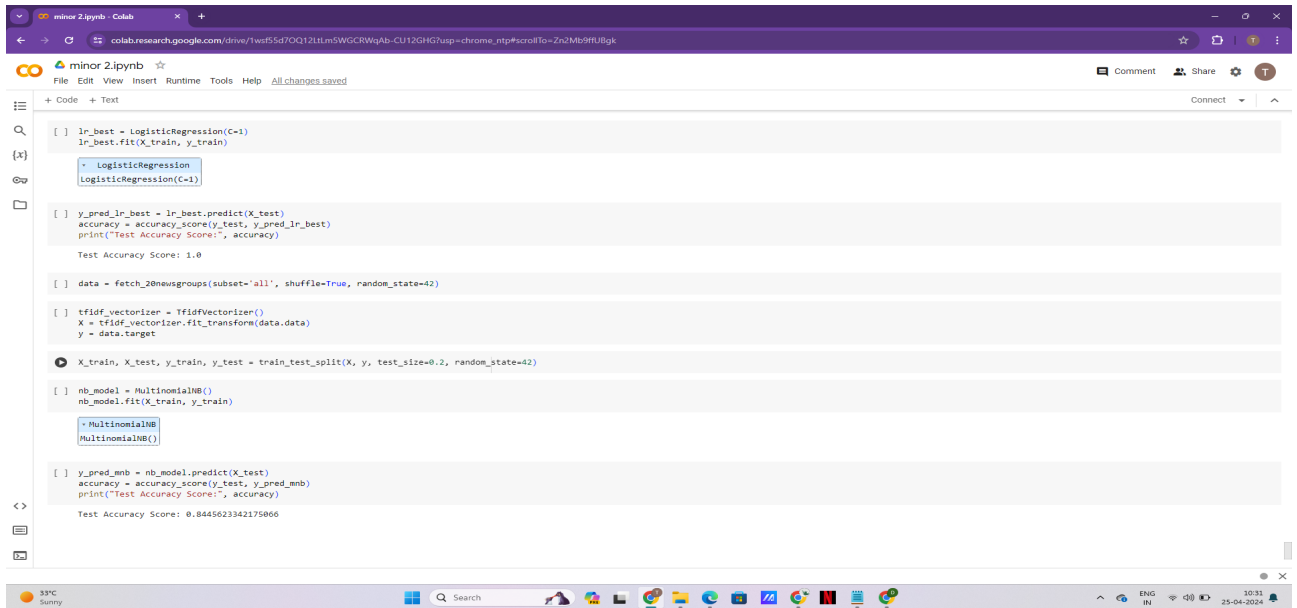
A strategy in this system testing is integrating the systems test cases along with the design techniques into a well planned series of steps with the results in the successful construction of software. This testing strategy should work along with the test planning, test execution, test case design, and the data collection resultant and evaluation. To accommodate a low-level test is a strategy for software testing. Software testing is a critical element of software quality assurance which represents the ultimate review of specification design and coding. Interesting anomalies from the software are being to represented by using Testing. Thus, the period of testing is performed for the proposed system and there by is ready for user acceptance testing.

Input

```
1 plt.title('Blood Pressure Density')
2 plt.xlabel('Blood Pressure')
3 plt.ylabel('Density')
4 plt.legend()
5 plt.show()
6 normal_wbc = df[(df['white_blood_cell_count'] >= 4500) & (df['white_blood_cell_count'] <= 11000)][ '
    white_blood_cell_count']
7 low_wbc = df[df['white_blood_cell_count'] < 4500]['white_blood_cell_count']
8 high_wbc = df[df['white_blood_cell_count'] > 11000]['white_blood_cell_count']
9
10 plt.figure(figsize=(12, 6))
11 sns.kdeplot(data=normal_wbc, color='green', shade=True, label='Normal White Blood Cell Count (4500 -
    11000)')
12 sns.kdeplot(data=low_wbc, color='blue', shade=True, label='Low White Blood Cell Count (< 4500)')
13 sns.kdeplot(data=high_wbc, color='red', shade=True, label='High White Blood Cell Count (> 11000)')
14
15 plt.title('White Blood Cell Count Density')
16 plt.xlabel('White Blood Cell Count')
17 plt.ylabel('Density')
18 plt.legend()
19 plt.show()
20 normal_rbc = df[((df['red_blood_cell_count'] >= 4.35) & (df['red_blood_cell_count'] <= 5.65)) | ((df
    ['red_blood_cell_count'] >= 3.92) & (df['red_blood_cell_count'] <= 5.13))]['red_blood_cell_count
    ']
21 low_rbc = df[(df['red_blood_cell_count'] < 4.35) | (df['red_blood_cell_count'] < 3.92)][ '
    red_blood_cell_count']
22 high_rbc = df[(df['red_blood_cell_count'] > 5.65) | (df['red_blood_cell_count'] > 5.13)][ '
    red_blood_cell_count']
23
24 plt.figure(figsize=(12, 6))
25 sns.kdeplot(data=normal_rbc, color='green', shade=True, label='Normal Red Blood Cell Count (4.35 -
    5.65)')
26 sns.kdeplot(data=low_rbc, color='blue', shade=True, label='Low Red Blood Cell Count (< 4.35 for men,
    < 3.92 for women)')
27 sns.kdeplot(data=high_rbc, color='red', shade=True, label='High Red Blood Cell Count (> 5.65 for men
    , > 5.13 for women)')
28
29 plt.title('Red Blood Cell Count Density')
30 plt.xlabel('Red Blood Cell Count')
31 plt.ylabel('Density')
32 plt.legend()
33 plt.show()
34 normal_hb_male = df[(df['haemoglobin'] >= 14) & (df['haemoglobin'] <= 18)][ 'haemoglobin']
35 normal_hb_female = df[(df['haemoglobin'] >= 12) & (df['haemoglobin'] <= 16)][ 'haemoglobin']
36 low_hb = df[(df['haemoglobin'] < 12) | (df['haemoglobin'] < 14)][ 'haemoglobin']
37 high_hb = df[(df['haemoglobin'] > 16) | (df['haemoglobin'] > 18)][ 'haemoglobin']
```

Test Result : PASSED

5.3.4 Test Result



The screenshot shows a Jupyter Notebook titled 'minor 2.ipynb' in a web browser. The notebook contains the following code cells:

```
[ ] lr_best = LogisticRegression(C=1)
lr_best.fit(X_train, y_train)

[ ] y_pred_lr_best = lr_best.predict(X_test)
accuracy = accuracy_score(y_test, y_pred_lr_best)
print("Test Accuracy Score:", accuracy)

Test Accuracy Score: 1.0

[ ] data = fetch_20newsgroups(subset='all', shuffle=True, random_state=42)

[ ] tfidf_vectorizer = TfidfVectorizer()
X = tfidf_vectorizer.fit_transform(data.data)
y = data.target

[ ] X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

[ ] nb_model = MultinomialNB()
nb_model.fit(X_train, y_train)

[ ] y_pred_mnb = nb_model.predict(X_test)
accuracy = accuracy_score(y_test, y_pred_mnb)
print("Test Accuracy Score:", accuracy)

Test Accuracy Score: 0.8445623342175866
```

The interface includes a left sidebar with icons for file explorer, search, and other notebook functions. The bottom status bar shows the system clock as 10:31 on 25-04-2024.

Figure 5.5: Test Image

Chapter 6

RESULTS AND DISCUSSIONS

6.1 Efficiency of the Proposed System

As we have mentioned the earlier that it is better to search for an alternative approach with promising results and we are going to do that. So, in this study, we are going to compare four very commonly known machine learning algorithms and we are going to state the best among them using a case study on Chronic Kidney disease Prediction. As stated, the four machine learning algorithms are: Random Forest Decision Tree LGBM K- Near Neighbour (KNN) Why we are going with only four algorithms is that we don't want to overload and confuse the learning procedure by including more machine learning algorithms. So, we thought that four would be an optimal choice and the machine would work just fine. Using confusion matrix we received the accuracy of all four methodologies such as K-Nearest Neighbor, Random Forest, Decision Tree and Light Gradient Boosted Machine. We receive the accuracy of KNN is 94 percentage and Random Forest, Decision Tree, Light Gradient Boosted Machine accuracy are 98 percentage. Compared to Light Gradient Boosted Machine, Decision Tree and Random Forest we receive less accuracy in K-Nearest Neighbor.

6.2 Comparison of Existing and Proposed System

```
➔ Random Forest Accuracy: 1.00
Random Forest Classification Report:
              precision    recall  f1-score   support

     0       1.00      1.00      1.00        52
     1       1.00      1.00      1.00        28

   accuracy          1.00          1.00          1.00        80
  macro avg          1.00          1.00          1.00        80
weighted avg          1.00          1.00          1.00        80

Random Forest Confusion Matrix:
[[52  0]
 [ 0 28]]
```

Figure 6.1: **Random Forest accuracy**

```
Logistic Regression Accuracy: 0.93
Logistic Regression Classification Report:
              precision    recall  f1-score   support

     0       0.94      0.94      0.94        52
     1       0.89      0.89      0.89        28

   accuracy          0.93          0.93          0.93        80
  macro avg          0.92          0.92          0.92        80
weighted avg          0.93          0.93          0.93        80

Logistic Regression Confusion Matrix:
[[49  3]
 [ 3 25]]
```

Figure 6.2: **Logistic Regression accuracy**

6.3 Sample Code

```
1 # Importing libraries
2 import numpy as np
3 import pandas as pd
4 import matplotlib.pyplot as plt
5 from sklearn.model_selection import train_test_split
6 from sklearn.preprocessing import LabelEncoder, MinMaxScaler
7 from sklearn.metrics import accuracy_score, confusion_matrix, classification_report
8 from sklearn.neighbors import KNeighborsClassifier
9 from sklearn.tree import DecisionTreeClassifier
10 from sklearn.ensemble import RandomForestClassifier
11 from keras.models import Sequential
12 from keras.layers import Dense
13 from keras.initializers import RandomNormal
14
15 # Read CSV file
16 df = pd.read_csv('chronickidney.csv')
17
18 # Selecting columns to retain
19 columns_to_retain = ['sg', 'al', 'sc', 'hemo', 'pcv', 'wbcc', 'rbcc', 'htn', 'classification']
20 df = df[columns_to_retain].dropna()
21
22 # Encoding categorical variables
23 for column in df.columns:
24     if df[column].dtype != np.number:
25         df[column] = LabelEncoder().fit_transform(df[column])
26
27 # Scaling features
28 x = df.drop('classification', axis=1)
29 y = df['classification']
30 scaler = MinMaxScaler()
31 x_scaled = scaler.fit_transform(x)
32
33 # Splitting data into training and testing sets
34 x_train, x_test, y_train, y_test = train_test_split(x_scaled, y, test_size=0.2, shuffle=True)
35
36 # Building a neural network model
37 model = Sequential()
38 model.add(Dense(256, input_dim=len(x.columns), kernel_initializer=RandomNormal(seed=13), activation=
    'relu'))
39 model.add(Dense(1, activation='hard_sigmoid'))
40 model.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy'])
41
42 # Training the model
43 history = model.fit(x_train, y_train, epochs=399, batch_size=x_train.shape[0])
44
45 # Saving the model
46 model.save('ckd.model')
```

```

47
48 # Plotting accuracy and loss
49 plt.plot(history.history['accuracy'])
50 plt.plot(history.history['loss'])
51 plt.title('Model Accuracy & Loss')
52 plt.ylabel('Accuracy and Loss')
53 plt.xlabel('Epoch')
54 plt.legend(['Accuracy', 'Loss'], loc='upper right')
55 plt.show()
56
57 # Printing shapes of training and testing data
58 print('Shape of training data:', x_train.shape)
59 print('Shape of testing data:', x_test.shape)
60
61 # Making predictions
62 pred = model.predict(x_test)
63 pred = [1 if y >= 0.5 else 0 for y in pred]
64
65 # Printing original and predicted values
66 print('Original:', ", ".join(str(x) for x in y_test))
67 print('Predicted:', ", ".join(str(x) for x in pred))
68
69 # KNN Classifier
70 knn = KNeighborsClassifier()
71 knn.fit(x_train, y_train)
72 knn_acc = accuracy_score(y_test, knn.predict(x_test))
73 print(f"Training Accuracy of KNN is {accuracy_score(y_train, knn.predict(x_train))}")
74 print(f"Test Accuracy of KNN is {knn_acc}")
75 print(f"Confusion Matrix:\n{confusion_matrix(y_test, knn.predict(x_test))}\n")
76 print(f"Classification Report:\n{classification_report(y_test, knn.predict(x_test))}")
77
78 # Decision Tree Classifier
79 dtc = DecisionTreeClassifier()
80 dtc.fit(x_train, y_train)
81 dtc_acc = accuracy_score(y_test, dtc.predict(x_test))
82 print(f"Training Accuracy of Decision Tree Classifier is {accuracy_score(y_train, dtc.predict(
      x_train))}")
83 print(f"Test Accuracy of Decision Tree Classifier is {dtc_acc}")
84 print(f"Confusion Matrix:\n{confusion_matrix(y_test, dtc.predict(x_test))}\n")
85 print(f"Classification Report:\n{classification_report(y_test, dtc.predict(x_test))}")
86
87 # Random Forest Classifier
88 rdcf = RandomForestClassifier(criterion='entropy', max_depth=11, max_features='auto',
      min_samples_leaf=2,
89                               min_samples_split=3, n_estimators=130)
90 rdcf.fit(x_train, y_train)
91 rdcf_acc = accuracy_score(y_test, rdcf.predict(x_test))
92 print(f"Training Accuracy of Random Forest Classifier is {accuracy_score(y_train, rdcf.predict(
      x_train))}")
93 print(f"Test Accuracy of Random Forest Classifier is {rdcf_acc}")

```

```

94 print(f"Confusion Matrix:\n{confusion_matrix(y_test, rdcf.predict(x_test))}\n")
95 print(f"Classification Report:\n{classification_report(y_test, rdcf.predict(x_test))}")
96
97 # Sorting models based on accuracy
98 models = pd.DataFrame({
99     'Model': ['KNN', 'Decision Tree Classifier', 'Random Forest Classifier'],
100     'Score': [knn_acc, dtc_acc, rdcf_acc]
101 })
102 models = models.sort_values(by='Score', ascending=False)
103
104 print(models)

```

Output

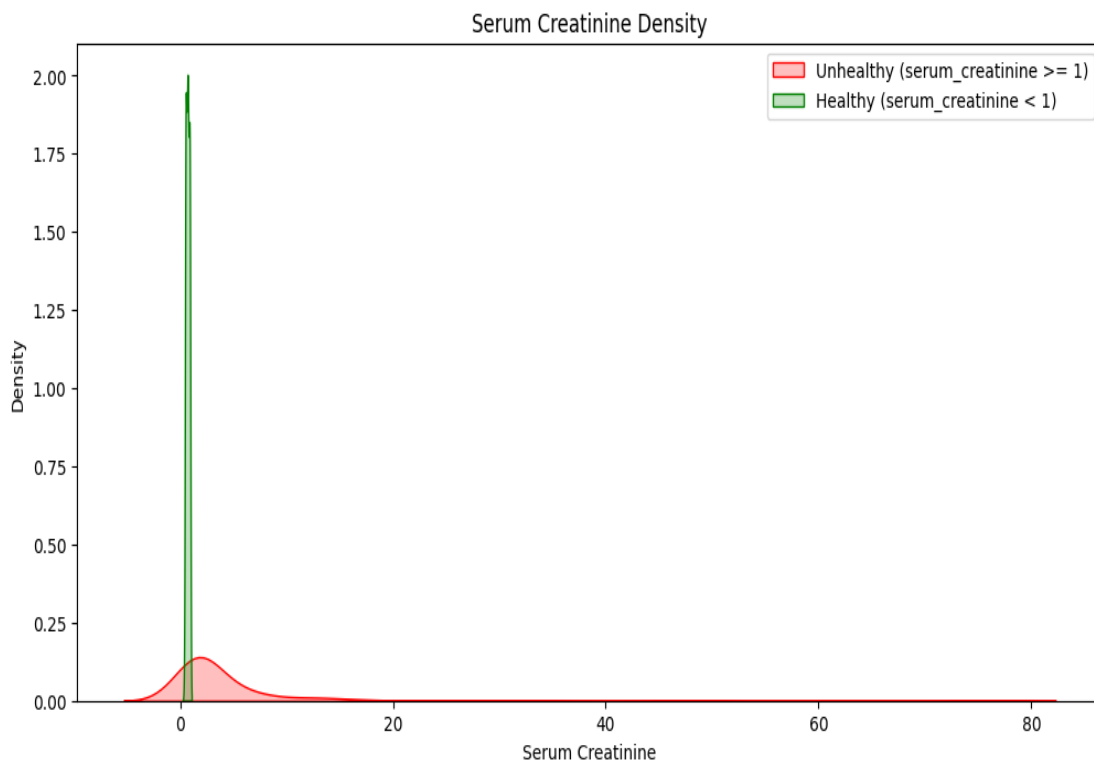


Figure 6.3: Serum Creatinine Density

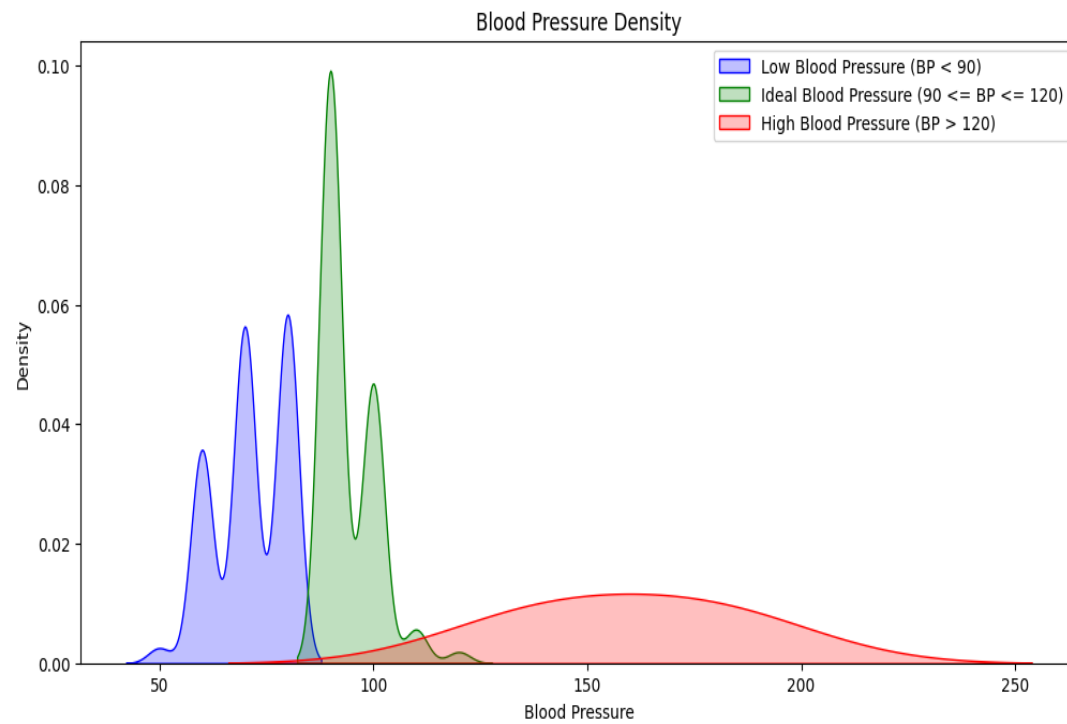


Figure 6.4: **Blood pressure Density**

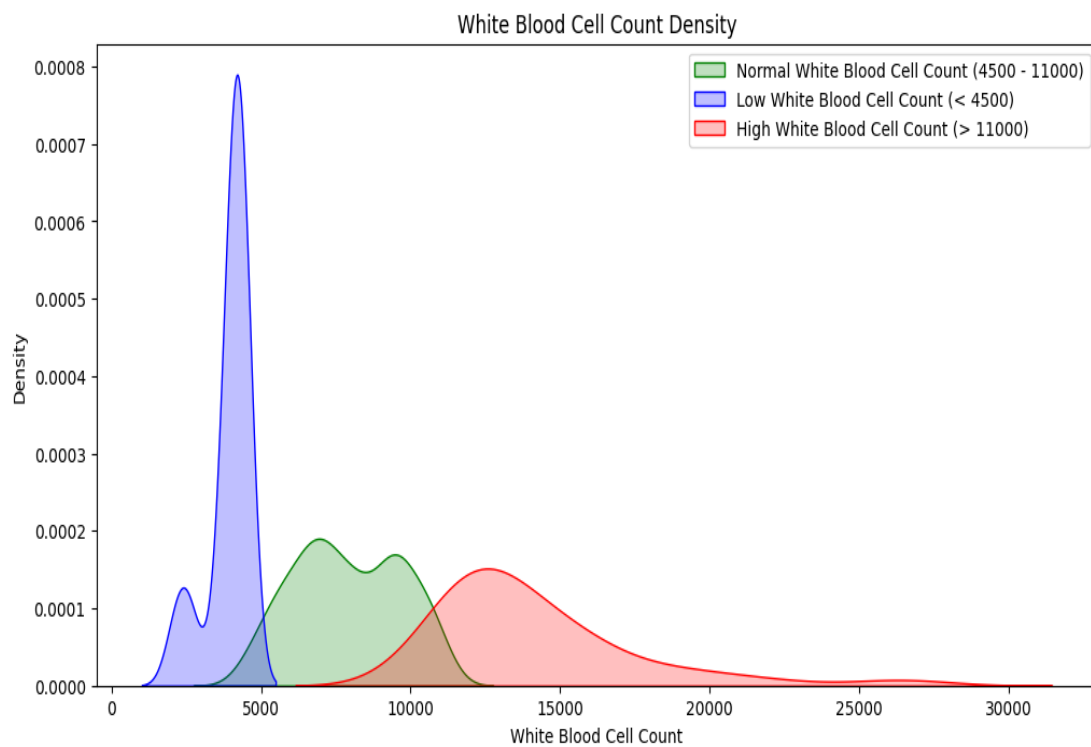


Figure 6.5: **White Blood Cell Count Density**

Chapter 7

CONCLUSION AND FUTURE ENHANCEMENTS

7.1 Conclusion

CKD means that the kidney does not work and cannot correctly filter the blood. About 10 percent of the population for a worldwide suffering from (CKD), and millions of die each year because of they couldn't get the affordable treatment, with the number increasing in the elderly. In 2010, a study was conducted by International Society of Numerology (ISN) on global burden disease, they reported that CKD has been raised an important cause of the mortality worldwide with the number of deaths increasing by 82.3 percent in the last two decades. Also, the number of patients reaching End-Stage Renal Disease (ESRD) is increasing, which requires kidney transplantation or dialysis to save the patient's lives.

7.2 Future Enhancements

As a complementary solution, we could build a mobile application, to make it available to everyone to detect whether they are CKD or not and as of now if we have a correct data set we could use it to predict many more diseases and also if the application collaborates with a highly reputed company then we can allow making more advancements to the existing model. There would be disadvantages because of we need to have higher graphic card compatibility to run the model and the solution for that would be, the involvement of a framework such as Kara's or Tensor-Flow at the back end.

Chapter 8

PLAGIARISM REPORT

The screenshot displays a web browser window with the URL `plagiarismdetector.net`. The page title is "Plagiarism Scan Report". At the top, there are three buttons: "Check Grammar", "Make it Unique", and "AI Detector". Below these, statistics are shown: Characters: 1024, Words: 154, Sentences: 9, and Speak Time: 2 Min. A sidebar on the right shows two circular progress indicators: a red one for "11% Plagiarized" and a green one for "89% Unique". Below these is a blue bar labeled "100%". A section titled "View Plagiarized Sources" shows "11% Plagiarized Content". The main content area contains a paragraph about Chronic Kidney Disease and machine learning, with some text highlighted in pink. At the bottom, there is an "Ad closed by Google" message and a Windows taskbar with various icons and a system clock showing 2:49 PM on 5/2/2024.

Plagiarism Scan Report

Check Grammar Make it Unique AI Detector

Characters: 1024 Words: 154 Sentences: 9 Speak Time: 2 Min

GO PRO Deep Search NO ADS SUPPORT Accurate Reports! Go Pro

Chronic Kidney Disease is a serious lifelong condition that induced by either kid- ney pathology or reduced kidney functions. According to the survey in 2021-2022 for every 10 in 3, people are suffering from this chronic kidney disease. It enables us to introduce the optimal subset of parameters to feed machine learning to build a set of predictive models. K-Nearest Neighbor, Random Forest, LGBM and Decision Tree methods are applied for prediction. Among these four methodologies, the pro- posed model suggests an LGBM is suitable for the early prediction of this kind of disease. Performance measures show that Light Gradient Boosted Machine (LGBM) gives 99 98 procedure concludes that advances in machine learning and predictive analytics, represent a promising model to recognize intelligent solutions, which in turn, gives the ability of predication in kidney Neighbor, Light Gradient Boo

11% Plagiarized 89% Unique

100%

View Plagiarized Sources

11% Plagiarized Content

Ad closed by Google

100°F Sunny Search ENG 2:49 PM 5/2/2024

Chapter 9

SOURCE CODE & POSTER PRESENTATION

9.1 Source Code

```
1
2 import pandas as pd
3 import numpy as np
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6 import plotly.express as px
7 from sklearn.preprocessing import OneHotEncoder
8 from sklearn.compose import ColumnTransformer
9
10 from sklearn.preprocessing import LabelEncoder
11 from sklearn.model_selection import train_test_split
12 from sklearn.ensemble import RandomForestClassifier, AdaBoostClassifier
13 from sklearn.linear_model import LogisticRegression
14 from sklearn.neighbors import KNeighborsClassifier
15 from sklearn.naive_bayes import MultinomialNB
16 from sklearn.tree import DecisionTreeClassifier
17 from sklearn.metrics import accuracy_score, classification_report, confusion_matrix, precision_score,
    f1_score, recall_score
18 from sklearn.metrics import roc_curve, roc_auc_score
19 from sklearn.metrics import precision_recall_curve, auc
20 from sklearn.model_selection import GridSearchCV
21 from sklearn.datasets import load_iris
22 from sklearn.datasets import fetch_20newsgroups
23 from sklearn.feature_extraction.text import TfidfVectorizer
24
25
26 import warnings
27 warnings.filterwarnings('ignore')
28 %matplotlib inline
29 pd.set_option('display.max_columns', 26)
30 df = pd.read_csv('kidney_disease.csv')
31 df.head()
32 df.drop('id', axis=1, inplace=True)
33 df.head()
```

```

34 df.columns = ['age', 'blood_pressure', 'specific_gravity', 'albumin', 'sugar', 'red_blood_cells', 'pus_cell',
35               'pus_cell_clumps', 'bacteria', 'blood_glucose_random', 'blood_urea', 'serum_creatinine',
36               'sodium', 'potassium', 'haemoglobin', 'packed_cell_volume', 'white_blood_cell_count', 'red_blood_cell_count',
37               'hypertension', 'diabetes_mellitus', 'coronary_artery_disease', 'appetite', 'peda_edema',
38               'aanemia', 'class']
39 df.sample(10)
40 df.info()
41 df.describe().T
42 df['packed_cell_volume'] = pd.to_numeric(df['packed_cell_volume'], errors='coerce')
43 df['white_blood_cell_count'] = pd.to_numeric(df['white_blood_cell_count'], errors='coerce')
44 df['red_blood_cell_count'] = pd.to_numeric(df['red_blood_cell_count'], errors='coerce')
45 cat_cols = [col for col in df.columns if df[col].dtype == 'object']
46 cat_cols
47 num_cols = [col for col in df.columns if df[col].dtype != 'object']
48 num_cols
49 for col in cat_cols:
50     print(f'{col} {df[col].unique()}')
51 df['diabetes_mellitus'].replace({'\tno': 'no', '\tyes': 'yes', ' yes': 'yes'}, inplace=True)
52 df['coronary_artery_disease'] = df['coronary_artery_disease'].str.replace('\tno', 'no')
53 df['class'].replace({'ckd\t': 'ckd', 'notckd': 'not_ckd'}, inplace=True)
54 df['class'] = df['class'].map({'ckd': 0, 'not_ckd': 1})
55 df['class'] = pd.to_numeric(df['class'], errors='coerce')
56 cols = ['diabetes_mellitus', 'coronary_artery_disease', 'class']
57
58 for col in cols:
59     print(f'{col} {df[col].unique()}')
60 value_counts = df['class'].value_counts()
61 categories = value_counts.index
62 counts = value_counts.values
63 colors = ['lightgreen', 'seagreen', 'forestgreen']
64 explode = [0.05] * len(counts)
65
66 plt.pie(counts, explode=explode, shadow=True, labels=categories, autopct='%1.1f%%', startangle=140,
67         pctdistance=0.85, colors=colors)
68 plt.axis('equal')
69 plt.title('Value Count')
70 plt.legend(categories, title='Classes', loc='best')
71 plt.show()
72 sns.pairplot(df);
73 plt.figure(figsize=(20, 15))
74 for i, column in enumerate(num_cols, 1):
75     plt.subplot(3, 5, i)
76     sns.distplot(df[column], color='green')
77     plt.xlabel(column)
78 plt.tight_layout()

```

```

79 plt.show()
80 plt.figure(figsize=(20, 15))
81 for i, column in enumerate(cat_cols, 1):
82     plt.subplot(3, 4, i)
83     sns.countplot(data=df, x=column, palette='husl')
84     plt.xlabel(column)
85     plt.ylabel('Count')
86
87 plt.tight_layout()
88 plt.show()
89 for column in cat_cols:
90     print(f"--- {column} ---")
91     total_count = len(df[column])
92     unique_values = df[column].value_counts()
93     for value, count in unique_values.items():
94         percentage = (count / total_count) * 100
95         print(f"{value}: {percentage:.1f}%")
96     print()
97 unhealthy = df[df['serum_creatinine'] >= 1]['serum_creatinine']
98 healthy = df[df['serum_creatinine'] < 1]['serum_creatinine']
99
100 plt.figure(figsize = (12,6))
101 sns.kdeplot(data=unhealthy, color='red', shade=True, label='Unhealthy (serum_creatinine >= 1)')
102 sns.kdeplot(data=healthy, color='green', shade=True, label='Healthy (serum_creatinine < 1)')
103
104 plt.title('Serum Creatinine Density')
105 plt.xlabel('Serum Creatinine')
106 plt.ylabel('Density')
107 plt.legend()
108 plt.show()
109 low_bp = df[df['blood_pressure'] < 90]['blood_pressure']
110 ideal_bp = df[(df['blood_pressure'] >= 90) & (df['blood_pressure'] <= 120)]['blood_pressure']
111 high_bp = df[df['blood_pressure'] > 120]['blood_pressure']
112
113 plt.figure(figsize=(12, 6))
114 sns.kdeplot(data=low_bp, color='blue', shade=True, label='Low Blood Pressure (BP < 90)')
115 sns.kdeplot(data=ideal_bp, color='green', shade=True, label='Ideal Blood Pressure (90 <= BP <= 120)')
116 sns.kdeplot(data=high_bp, color='red', shade=True, label='High Blood Pressure (BP > 120)')
117
118 plt.title('Blood Pressure Density')
119 plt.xlabel('Blood Pressure')
120 plt.ylabel('Density')
121 plt.legend()
122 plt.show()
123 normal_wbc = df[(df['white_blood_cell_count'] >= 4500) & (df['white_blood_cell_count'] <= 11000)]['white_blood_cell_count']
124 low_wbc = df[df['white_blood_cell_count'] < 4500]['white_blood_cell_count']
125 high_wbc = df[df['white_blood_cell_count'] > 11000]['white_blood_cell_count']
126

```

```

127 plt.figure(figsize=(12, 6))
128 sns.kdeplot(data=normal_wbc, color='green', shade=True, label='Normal White Blood Cell Count (4500 -
    11000)')
129 sns.kdeplot(data=low_wbc, color='blue', shade=True, label='Low White Blood Cell Count (< 4500)')
130 sns.kdeplot(data=high_wbc, color='red', shade=True, label='High White Blood Cell Count (> 11000)')
131
132 plt.title('White Blood Cell Count Density')
133 plt.xlabel('White Blood Cell Count')
134 plt.ylabel('Density')
135 plt.legend()
136 plt.show()
137 normal_rbc = df[((df['red_blood_cell_count'] >= 4.35) & (df['red_blood_cell_count'] <= 5.65)) | ((df
    ['red_blood_cell_count'] >= 3.92) & (df['red_blood_cell_count'] <= 5.13))]['red_blood_cell_count
    ']
138 low_rbc = df[(df['red_blood_cell_count'] < 4.35) | (df['red_blood_cell_count'] < 3.92)][ '
    red_blood_cell_count']
139 high_rbc = df[(df['red_blood_cell_count'] > 5.65) | (df['red_blood_cell_count'] > 5.13)][ '
    red_blood_cell_count']
140
141 plt.figure(figsize=(12, 6))
142 sns.kdeplot(data=normal_rbc, color='green', shade=True, label='Normal Red Blood Cell Count (4.35 -
    5.65)')
143 sns.kdeplot(data=low_rbc, color='blue', shade=True, label='Low Red Blood Cell Count (< 4.35 for men,
    < 3.92 for women)')
144 sns.kdeplot(data=high_rbc, color='red', shade=True, label='High Red Blood Cell Count (> 5.65 for men
    , > 5.13 for women)')
145
146 plt.title('Red Blood Cell Count Density')
147 plt.xlabel('Red Blood Cell Count')
148 plt.ylabel('Density')
149 plt.legend()
150 plt.show()
151 normal_hb_male = df[(df['haemoglobin'] >= 14) & (df['haemoglobin'] <= 18)][ 'haemoglobin']
152 normal_hb_female = df[(df['haemoglobin'] >= 12) & (df['haemoglobin'] <= 16)][ 'haemoglobin']
153 low_hb = df[(df['haemoglobin'] < 12) | (df['haemoglobin'] < 14)][ 'haemoglobin']
154 high_hb = df[(df['haemoglobin'] > 16) | (df['haemoglobin'] > 18)][ 'haemoglobin']
155
156 # KDE grafi ini izme
157 plt.figure(figsize=(12, 6))
158 sns.kdeplot(data=normal_hb_male, color='green', shade=True, label='Normal Haemoglobin (14 - 18) for
    males')
159 sns.kdeplot(data=normal_hb_female, color='blue', shade=True, label='Normal haemoglobin (12 - 16) for
    females')
160 sns.kdeplot(data=low_hb, color='red', shade=True, label='Low Haemoglobin (< 12 for both)')
161 sns.kdeplot(data=high_hb, color='purple', shade=True, label='High Haemoglobin (> 16 for females, >
    18 for males)')
162
163 plt.title('Hemoglobin Level Density')
164 plt.xlabel('Hemoglobin Level (g/dl)')
165 plt.ylabel('Density')

```

```

166 plt.legend()
167 plt.show()
168 normal_sodium = df[(df['sodium'] >= 135) & (df['sodium'] <= 145)][ 'sodium' ]
169 low_sodium = df[df['sodium'] < 135][ 'sodium' ]
170 high_sodium = df[df['sodium'] > 145][ 'sodium' ]
171
172 plt.figure(figsize=(12, 6))
173 sns.kdeplot(data=normal_sodium, color='green', shade=True, label='Normal Sodium (135 - 145)')
174 sns.kdeplot(data=low_sodium, color='blue', shade=True, label='Low Sodium (< 135)')
175 sns.kdeplot(data=high_sodium, color='red', shade=True, label='High Sodium (> 145)')
176
177 plt.title('Sodium Level Density')
178 plt.xlabel('Sodium Level (mEq/L)')
179 plt.ylabel('Density')
180 plt.legend()
181 plt.show()
182 normal_urea = df[(df['blood_urea'] >= 5) & (df['blood_urea'] <= 20)][ 'blood_urea' ]
183 low_urea = df[df['blood_urea'] < 5][ 'blood_urea' ]
184 high_urea = df[df['blood_urea'] > 20][ 'blood_urea' ]
185
186 plt.figure(figsize=(12, 6))
187 sns.kdeplot(data=normal_urea, color='green', shade=True, label='Normal Urea (5 - 20)')
188 sns.kdeplot(data=low_urea, color='blue', shade=True, label='Low Urea (< 5)')
189 sns.kdeplot(data=high_urea, color='red', shade=True, label='High Urea (> 20)')
190
191 plt.title('Blood Urea Level Density')
192 plt.xlabel('Blood Urea Level (mg/dl)')
193 plt.ylabel('Density')
194 plt.legend()
195 plt.show()
196 sns.set_theme(style="whitegrid", palette=["red", "green"])
197
198 # Draw a categorical scatterplot to show each observation
199 ax = sns.swarmplot(data=df, x="class", y="age", hue='class')
200 ax.set(ylabel="AGE")
201 df[num_cols].isnull().sum()
202 median_values = df[num_cols].median()
203 df[num_cols] = df[num_cols].fillna(median_values)
204 df[num_cols].isnull().sum()
205 df[cat_cols].isnull().sum()
206 df[cat_cols].isnull().sum()
207 df[cat_cols].isnull().sum()
208 LabelEncoder = LabelEncoder()
209
210 for col in cat_cols:
211     df[col] = LabelEncoder.fit_transform(df[col])
212 df.head()
213 y = df['class']
214 X = df.drop(['class'], axis=1)
215 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.2, random_state = 42)

```



```

216 models = {
217     "Random Forest": RandomForestClassifier(),
218     "Logistic Regression": LogisticRegression(),
219     "K-Nearest Neighbors": KNeighborsClassifier(),
220     "MultinomialNB": MultinomialNB(),
221     "Decision Tree": DecisionTreeClassifier(),
222     "AdaBoost": AdaBoostClassifier(),
223 }
224
225 for model_name, model in models.items():
226     model.fit(X_train, y_train)
227
228     # Make predictions on the test set
229     y_pred = model.predict(X_test)
230
231     # Evaluate the model
232     accuracy = accuracy_score(y_test, y_pred)
233     print(f'{model_name} Accuracy: {accuracy:.2f}')
234
235     # Generate a classification report
236     print(f'{model_name} Classification Report:')
237     print(classification_report(y_test, y_pred))
238
239     cm = confusion_matrix(y_test, y_pred)
240     print(f'{model_name} Confusion Matrix:')
241     print(cm)
242
243     print('-' * 55)
244 lr = LogisticRegression()
245 lr.fit(X_train, y_train)
246 y_pred_lr = lr.predict(X_test)
247 print('Accuracy', '%.3f' % accuracy_score(y_test, y_pred_lr))
248 print('Precision', '%.3f' % precision_score(y_test, y_pred_lr))
249 print('Recall', '%.3f' % recall_score(y_test, y_pred_lr))
250 print('F1 Score', '%.3f' % f1_score(y_test, y_pred_lr))
251 iris = load_iris()
252 X = iris.data
253 y = iris.target
254 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
255 lr = LogisticRegression()
256 param_grid = {'C': [0.001, 0.01, 0.1, 1, 10, 100, 1000]}
257 grid_search = GridSearchCV(lr, param_grid, cv=5, scoring='accuracy')
258 grid_search.fit(X_train, y_train)
259
260 # En iyi hiperparametreleri ve do ruluk skorunu yazd rma
261 print("Best Params:", grid_search.best_params_)
262 print("Best accuracy score:", grid_search.best_score_)
263 lr_best = LogisticRegression(C=1)
264 lr_best.fit(X_train, y_train)
265 y_pred_lr_best = lr_best.predict(X_test)

```

```
266 accuracy = accuracy_score(y_test , y_pred_lr_best)
267 print("Test Accuracy Score:", accuracy)
268 data = fetch_20newsgroups(subset='all' , shuffle=True , random_state=42)
269 tfidf_vectorizer = TfidfVectorizer()
270 X = tfidf_vectorizer.fit_transform(data.data)
271 y = data.target
272 X_train , X_test , y_train , y_test = train_test_split(X, y, test_size=0.2, random_state=42)
273 nb_model = MultinomialNB()
274 nb_model.fit(X_train , y_train)
275 y_pred_mnb = nb_model.predict(X_test)
276 accuracy = accuracy_score(y_test , y_pred_mnb)
277 print("Test Accuracy Score:", accuracy)
```

9.2 Poster Presentation

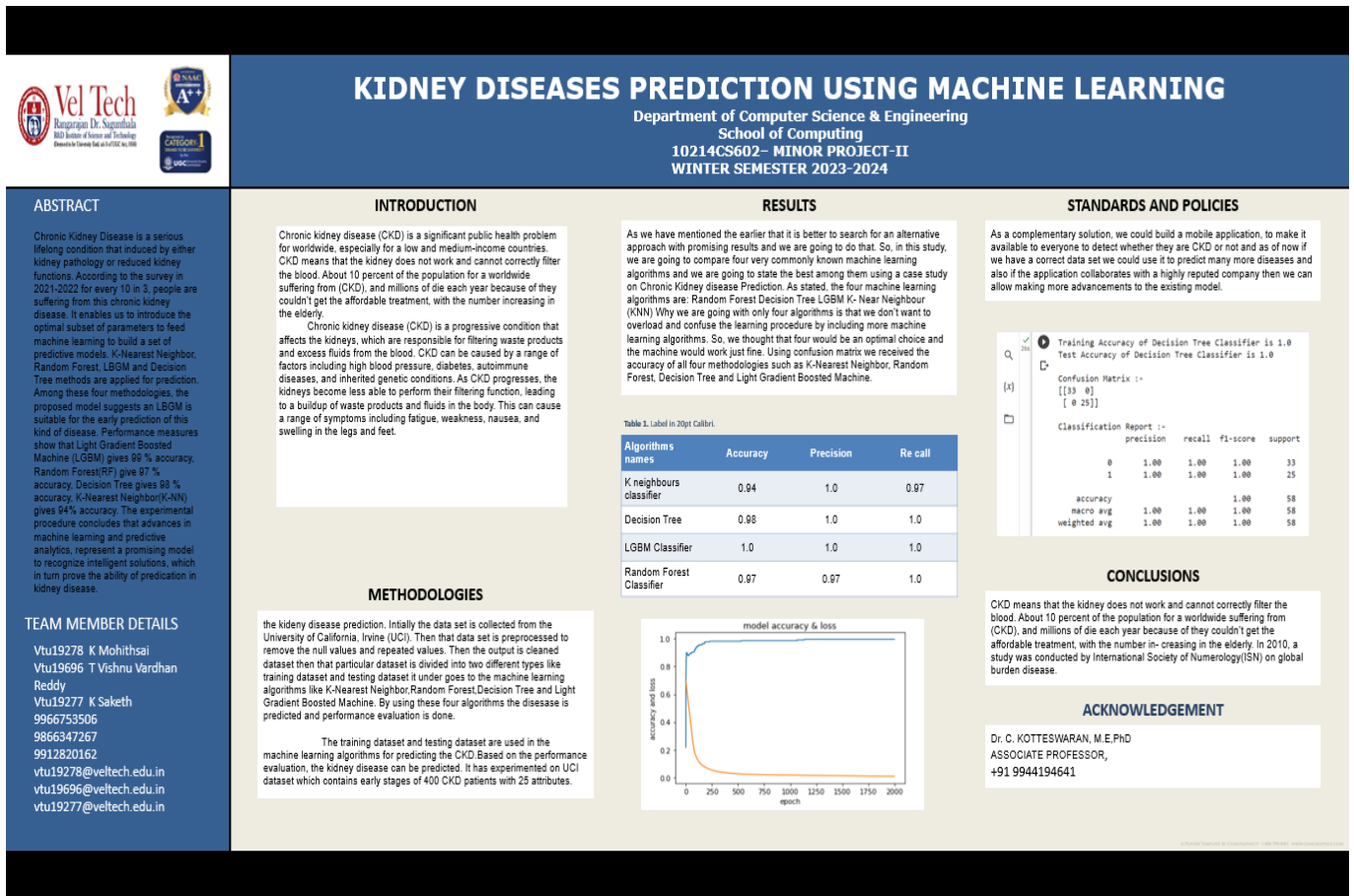


Figure 9.1: Poster

References

- [1] Imesh Udara Ekanayake Damayanthi Herath , "Chronic Kidney Disease Prediction Using Machine Learning Methods," in 2020 IEEE Moratuwa Engineering Research Conference, 2020 .
- [2] Yedilkhan Amirgaliyev Shahriar Shamiluulu Azamat Serek, "Analysis of Chronic Kidney Disease Dataset by Applying Machine Learning Methods",International Journal of Computing and Business Research (IJCBR), vol. 6, no. 2, (2019).
- [3] T Shaikhina, Torgyn, et al. "Implementation of Machine Learning Models for the Prevention of Kidney Diseases (CKD) or Their Derivatives" Hindawi Computational Intelligence and Neuroscience Volume 2021, Article ID 3941978, 8 pages <https://doi.org/10.1155/2021/3941978> (2021).
- [4] PANKAJ CHITTORA ZBIGNIEW LEONOWICZ SANDEEP CHAURASIA, "Prediction of Chronic Kidney Disease - A Machine Learning Perspective",Article in IEEE Access. D
- [5] Jing Xiao1 Ruifeng Ding Xiulin Xu Haochen Guan , et al. "Comparison and development of machine learning tools in the prediction of chronic kidney disease progression." Xiao et al. J Transl Med 17:119 <https://doi.org/10.1186/s12967-019-1860-0> (2019).
- [6] Yesubabu Kakitapallia Janakiram Ampolua Satya Dinesh Madasub, et al., "Detailed Review of Chronic Kidney Disease" International Publishing, Kidney Dis DOI: 10.1159/000504622. (2019).
- [7] Reshma S Salma Shaji S R Ajina Vishnu Priya S R,"Chronic Kidney Disease Prediction using Machine Learning",International Journal of Engineering Research Technology (IJERT) ISSN: 2278-0181 IJERTV9IS070092 Vol. 9 Issue 07, July-2020.
- [8] S.Revathy B.Bharathi et al, "Chronic Kidney Disease Prediction using Machine Learning Models ," International Journal of Engineering and Advanced Technology (IJEAT) ISSN: 2249 – 8958, Volume-9 Issue-1, Oct.

- [9] B.Sasi Varna G.Pravallika, "CHRONIC KIDNEY DISEASE METHODOLOGY BY USING MACHINE LEARNING," IJRTI Volume 7, Issue 6 ISSN: 2456-3315(2021).
- [10] Dibaba Adeba Debal and Tilahun Melak Sitote, "Chronic kidney disease prediction using machine learning techniques," Informatics in Medicine Unlocked, Debal and Sitote Journal of Big Data 9:109 <https://doi.org/10.1186/s40537-022-00657-5> (2022).