

A Project Report on
EARLY PREDICTION OF CHRONIC KIDNEY DISEASE

Submitted to JNTUK in the partial fulfilment of the requirements for the award of the
Degree of

BACHELOR OF TECHNOLOGY
IN
COMPUTER SCIENCE AND ENGINEERING
BY

P. FEERDOSH KHAN	21KP1A0585
A.V.S RAGHAVENDRA	22KP5A0501
SK. MOULANA	21KP1A0597
SD. RIYAZ BASHA	21KP1A05A3
K. EMMANUEL RAJU	22KP5A0505

Under the Esteemed Guidance of

Mrs. Y. JESSY KUMARI, M.Tech



DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING

NRI INSTITUTE OF TECHNOLOGY

(AUTONOMOUS)

ACCREDITED BY NAAC WITH A+ GRADE

(APPROVED BY AICTE, AFFILIATED TO JNTUK, KAKINADA)

VISADALA, MEDIKONDURU MANDAL, GUNTUR(DIST), A.P.

2021 - 2025

NRI INSTITUTE OF TECHNOLOGY
(AUTONOMOUS)
ACCREDITED BY NAAC WITH A+ GRADE
(APPROVED BY AICTE, AFFILIATED TO JNTUK, KAKINADA)

VISADALA, MEDIKONDURU MANDAL, GUNTUR(DIST), A.P.
2021 - 2025

DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING



CERTIFICATE

This is to certify that the Industrial Internship Project work entitled **EARLY PREDICTION OF CHRONIC KIDNEY DISEASE** being submitted by

P. FEERDOSH KHAN	21KP1A0585
A.V.S RAGHAVENDRA	22KP5A0501
SK. MOULANA	21KP1A0597
SD. RIYAZ BASHA	21KP1A05A3
K. EMMANUEL RAJU	22KP5A0505

in partial fulfilment for the award of the Degree of Bachelor of Technology in Computer Science and Engineering in the **NRI INSTITUTE OF TECHNOLOGY** is a record of bonafide work carried out by them

Internal Guide

Mrs. Y. JESSY
KUMARI

Head of the Department

Dr. K. NAGESHWARAO
Professor & HOD

External Examiner

DECLARATION

We **P.FEERDOSH KHAN(21KP1A0585), A.V.S.RAGHAVENDRA(22KP5A0501), SK.MOULANA(21KP1A0597), SD.RIYAZ BASHA (21KP1A0593) and K.EMMANUEL RAJU (22KP5A0505)** hereby declare that the project report titled “**EARLY PREDICTION OF CHRONIC KIDNEY DISEASE**” under the guidance of **MRS.JESSY KUMARI, M.Tech** is submitted in partial fulfilment of the requirements for the degree of **Bachelor of Technology in Computer Science and Engineering**.

This is a record of bonafide work carried out by us and the results embodied in this project have not been reproduced or copied from any source. The results embodied in this project have not been submitted to any other university for the award of any degree.

P. FEERDOSH KHAN	21KP1A0585
A.V.S RAGHAVENDRA	22KP5A0501
SK. MOULANA	21KP1A0597
SD. RIYAZ BASHA	21KP1A05A3
K. EMMANUEL RAJU	22KP5A0505

ACKNOWLEDGEMENT

The present project work is the several days study of the various aspects of the project development. During this the effort in the present study, we have received a great amount of help from my **Dr. ALAPATI RAVINDRA Garu, Chairman, Secretary & Correspondent, NRI INSTITUTE OF TECHNOLOGY VISADALA**, which we wish to acknowledge and thank from depth of our heart.

We thankful to our Principal **Dr. KOTA SRINIVAS, NRI INSTITUTE OF TECHNOLOGY** for permitting and encouraging us in doing this project. We deeply intended to **P.JESSY KUMARI(MTech)** whose motivation and constant encouragement has led to pursue a project in the field of software development.

We are very much obliged and thankful to our internal guide **DR.K.NAGESHWARAO MTech (PhD), Assistant Professor &HOD** for providing this opportunity and constant encouragement given by him during the course. We are grateful to him for his valuable guidance and suggestions during our project work.

We also thank to our parents who supported us for completing this project. Finally, we express our deep sense of gratitude and thanks to the Teaching and Non-Teaching Staff at our college who stood with us during the project and helped us to make it a successful venture

P. FEERDOSH KHAN

21KP1A0585

A.V.S RAGHAVENDRA

22KP5A0501

SK. MOULANA

21KP1A0597

SD. RIYAZ BASHA

21KP1A05A3

K. EMMANUEL RAJU

22KP5A0505

ABSTRACT

This project focuses on the prediction of chronic kidney disease (CKD) using machine learning techniques to enhance early diagnosis and treatment, thereby improving patient outcomes. The study utilizes a comprehensive dataset consisting of various health factors associated with CKD. The preprocessing phase involves normalization and label encoding to prepare the data for analysis, followed by exploratory data analysis (EDA) to uncover patterns and relationships within the dataset. Feature engineering is applied to enhance model performance by selecting the most relevant attributes for classification. Various machine learning algorithms, including Random Forest Classifier, Decision Tree Classifier, Gaussian Naive Bayes, and Support Vector Machine (SVM), are implemented to predict the likelihood of CKD. The Random Forest Classifier achieves a remarkable accuracy of 100%, while the Gaussian Naive Bayes also reaches 100%. The Decision Tree Classifier and SVM demonstrate accuracies of 96.88% and 96%, respectively. To facilitate real-time predictions, a user interface is developed using Flask, allowing healthcare professionals to input patient data and receive immediate predictions on CKD risk. This project not only highlights the effectiveness of machine learning in healthcare applications but also emphasizes the importance of timely diagnosis and intervention for chronic diseases like CKD, ultimately contributing to better patient care and management.

Keywords: Machine Learning, Chronic Kidney Disease Prediction, Healthcare Analytics, Disease Diagnosis, Medical Informatics, Kidney Disease Management, Predictive Medicine, Clinical Decision Support, Healthcare Technology, Personalized Healthcare.

TABLE OF CONTENTS

Contents	NO	PAGE
Certificate		i-ii
Declaration		iii
Acknowledgement		iv
Abstract		v
Chapter 1 INTRODUCTION		1-3
1.1 Project Idea		2
1.2 Motivation of the Project		2
1.3 Problem Statement		3
1.4 Statement of scope		3
1.5 Goals and objectives		3
Chapter 2 LITERATURE SURVEY		4-8
2.1 Related Work		5-6
2.2 Existing System		7
2.3 Disadvantages of the Existing System		7
2.4 Feasibility Study		8
Chapter 3 PROPOSED SYSTEM		9-13
3.1 Proposed System		10
3.2 Methodology		10-12
3.3 Advantages		13
3.4 Approaches		13

Chapter 4 SYSTEM REQUIREMENTS SPECIFICATION	14-16
4.1 Software Requirements	15
4.2 Technologies used	15
4.3 Hardware Requirements	16
Chapter 5. SYSTEM DESIGN	17-23
5.1 System Architecture	18
5.2 Module Description	18-19
5.3 UML Diagrams	20-23
5.3.1 Use Case Diagram	20
5.3.2 Sequence Diagram	21-22
5.3.3 Activity Diagram	23
Chapter 6 IMPLEMENTATION	24-41
6.1 Algorithms	25-34
6.1.1 What is Machine Learning	26
6.1.2 How Machine Learning Works	26-27
6.1.3 Random Forest	28
6.1.4 Decision Tree	29
6.1.5 Gaussian Naive Bayes	30
6.1.6 Support Vector Machines	31
6.1.7 Algorithm	32-33
6.1.8 Implementation Procedure	34

6.2 Coding (Source Code)	35-40
6.3 Dataset Details	41
Chapter 7 SYSTEM TESTING	42-45
7.1 Testing Introduction	43
7.1.1 Unit testing	43
7.1.2 Integration testing	43
7.1.3 User interface testing	44
7.2 Test Cases	45
8. OUTPUT SCREENS	46-49
9. TIMELINE OF THE PROJECT	50-51
10. CONCLUSION	52-54
10.1 Future Enhancement	54
REFERENCES	55

LIST OF FIGURES

NAME	PAGE NO
1 Project Flow	12
2 Project Architecture	18
3 Use case diagram	20
4 Sequence diagram	21
5 Activity diagram	23
6 Comparison of Models	47
7 Home Page	47
8 Predictions page	48
9 Output Predictions	48-49
9.1 Output Prediction -1	48
9.2 Output Prediction - 2	49

LIST OF TABLES

NAME	PAGE NO
1 Literature Survey	5
2 Data Set	41
3 Test Cases	45
4 Time Of The Project	51

CHAPTER 1

INTRODUCTION

1. INTRODUCTION

1.1 PROJECT IDEA

Chronic kidney disease (CKD) is a progressive condition characterized by a gradual loss of kidney function over time, affecting millions of individuals worldwide. According to the World Health Organization, approximately 850 million people are affected by kidney diseases, and CKD is estimated to be responsible for 2.4 million deaths annually. The prevalence of CKD is rising due to factors such as an aging population, increased incidence of diabetes and hypertension, and lifestyle changes. Early detection is vital, as CKD often remains asymptomatic in its initial stages, leading to late diagnosis when treatment options are limited, and complications become severe. Machine learning offers a promising approach to address this challenge by analyzing complex health datasets and identifying patterns that may indicate an elevated risk of CKD. By leveraging algorithms that can learn from historical data, it is possible to develop predictive models that aid healthcare professionals in making informed decisions about patient management.

1.2 MOTIVATION OF THE PROJECT

- **Increasing Prevalence:** The rising incidence of chronic kidney disease globally underscores the urgent need for early detection and effective management strategies.
- **Impact on Quality of Life:** CKD can significantly affect patients' quality of life, making early diagnosis crucial for maintaining health and well-being.
- **Resource Optimization:** Predictive models can help healthcare systems allocate resources more efficiently by identifying high-risk patients who require immediate attention.
- **Advancements in Technology:** The rapid development of machine learning and data analytics provides unprecedented opportunities to enhance predictive accuracy and healthcare delivery.
- **Personalized Healthcare:** Leveraging data-driven insights promotes a shift towards personalized medicine, allowing for tailored treatment plans that address individual patient needs.

1.3 PROBLEM STATEMENT

Chronic kidney disease (CKD) poses a significant global health challenge, often leading to severe complications and increased mortality due to late diagnosis and intervention. Despite the availability of relevant health data, there is a lack of effective predictive tools to assist healthcare professionals in identifying individuals at risk of CKD early on. This project aims to address this gap by utilizing machine learning techniques to analyze health factors and accurately predict the likelihood of CKD, thereby enabling timely intervention and improving patient outcomes.

1.4 STATEMENT OF SCOPE

The scope of this project encompasses the comprehensive analysis and modeling of chronic kidney disease (CKD) prediction using machine learning techniques. It includes the collection and preprocessing of health factor data, ensuring proper normalization and encoding for effective analysis. The project will involve exploratory data analysis (EDA) to identify significant patterns and correlations within the dataset, followed by feature engineering to optimize model performance. Various machine learning algorithms, such as Random Forest, Decision Tree, Gaussian Naive Bayes, and Support Vector Machine (SVM), will be applied to develop predictive models.

1.5 GOALS AND OBJECTIVES

The objective of this project is to develop a robust machine learning model for the early prediction of chronic kidney disease (CKD) using health factor data. By leveraging advanced algorithms such as Random Forest, Decision Tree, Gaussian Naive Bayes, and Support Vector Machine (SVM), the project aims to achieve high accuracy in predicting CKD risk, thereby enabling healthcare professionals to make informed decisions for timely diagnosis and treatment. Additionally, the project seeks to create an intuitive user interface using Flask for seamless real-time predictions, ultimately improving patient management and outcomes.

CHAPTER 2

LITERATURE SURVEY

2. LITERATURE SURVEY

2.1 RELATED WORK

S.NO	TITLE	AUTHORS	YEAR	KEY FINDINGS
1	Chronic Kidney Disease: Diagnosis and Management	Smith, J., Lee, R. & Taylor, A.	2018	Emphasized the importance of early diagnosis in CKD and recommended updated .
2	Global Burden of Chronic Kidney Disease	Johnson, D. W., Garcia, M. & Brown K	2019	Reported an increase in CKD prevalence worldwide, linking it to diabetes
3	Machine Learning Models for Early Prediction of CKD	Martinez, P. Li Y. & Patel. V	2020	Demonstrated the efficacy of machine learning models, specifically Random Forest
4	The Role of Diet in Chronic Kidney Disease Progression	Wang, S., Tran, H., & Kim, N.	2017	Identified a low protein diet as beneficial in slowing CKD progression,

Table 1. Literature survey

In this systematic review, the authors analyze various studies that have employed machine learning techniques for the diagnosis of chronic kidney disease. The paper categorizes the different algorithms used, discusses the datasets leveraged, and evaluates the strengths and limitations of each approach. The review emphasizes the need for standardized methodologies and comprehensive datasets to enhance the reliability of CKD predictions.

Summary

Khan et al. conclude that while machine learning holds significant promise for CKD diagnosis, challenges such as data quality, feature engineering, and model interpretability must be addressed to improve predictive performance and clinical applicability.

1. " A Hybrid Model for Predicting Chronic Kidney Disease Using Ensemble Learning Techniques" by Chenetal.

This research proposes a hybrid model that combines multiple Machine learning techniques, including Decision Trees, K-Nearest Neighbors, and Logistic Regression, to enhance the predictive accuracy for chronic kidney disease. The authors introduce an ensemble learning approach that leverages the strengths of individual models while mitigating their weaknesses. The study employs a diverse dataset containing various health indicators relevant to CKD.

Summary

The findings presented by Chen et al. demonstrate that the hybrid model outperforms single algorithms in terms of accuracy and robustness. This approach not only improves predictive outcomes but also emphasizes the importance of combining diverse techniques to address complex health-related challenges.

2. " Early Detection of Chronic Kidney Disease Using Data Mining Techniques" by Patel .

This paper explores the application of data mining techniques in the early detection of chronic kidney disease, focusing on algorithms such as Neural Networks and Decision Trees. The authors highlight the role of data preprocessing, including normalization and handling missing values, in improving model performance. The study utilizes a real-world dataset, examining the predictive capabilities of various data mining approaches.

Summary

Patel et al. find that the implementation of data mining techniques significantly enhances early CKD detection, with Neural Networks achieving high accuracy rates. The research underscores the potential of data-driven methodologies in transforming healthcare practices and facilitating proactive patient management.

2.2 EXISTING SYSTEM

Existing methods for chronic kidney disease (CKD) prediction primarily involve traditional statistical approaches and machine learning algorithms. Commonly used techniques include logistic regression, decision trees, and support vector machines (SVM), which analyze patient health data such as blood pressure, glucose levels, and other clinical indicators. These methods typically require extensive preprocessing, including data normalization and feature selection, to improve model accuracy. While these approaches can yield reasonable predictions, they often lack the flexibility and adaptability of more advanced machine learning models, which can better capture complex relationships within the data. Additionally, many existing models do not incorporate real-time prediction capabilities, limiting their practical application in clinical settings.

2.3 DISADVANTAGES OF THE EXISTING SYSTEM

- **Limited Predictive Power:** Traditional methods may struggle to capture complex nonlinear relationships within health data, leading to lower predictive accuracy compared to advanced machine learning techniques.
- **High Dependency on Feature Engineering:** Existing methods often require extensive feature selection and preprocessing, which can be time-consuming and may result in the omission of relevant data if not conducted thoroughly.
- **Poor Scalability:** Many conventional approaches are not easily scalable to large datasets, making it challenging to apply them in real-world clinical settings where data volumes are continuously increasing.
- **Lack of Real-Time Prediction:** Traditional models often do not support real-time predictions, which are crucial for timely decision-making in healthcare, potentially delaying critical interventions for patients.

2.4 FEASIBILITY STUDY

The feasibility of the project is analyzed in this phase and business proposal is put forth with a very general plan for the project and some cost estimates. During system analysis the feasibility study of the proposed system is to be carried out. This is to ensure that the proposed system is not a burden to the company. For feasibility analysis, some understanding of the major requirements for the system is essential.

Economical Feasibility

This study is carried out to check the economic impact that the system will have on the organization. The amount of fund that the company can pour into the research and development of the system is limited. The expenditures must be justified. This the developed system as well within the budget and this was achieved because most of the technologies used are freely available. Only the customized products had to be purchased.

Technical Feasibility

This study is carried out to check the technical feasibility, that is, the technical requirements of the system. Any system developed must not have a high demand on the available technical resources. This will lead to high demands on the available technical resources. This will lead to high demands being placed on the client.

Social Feasibility

The aspect of study is to check the level of acceptance of the system by the user. This includes the process of training the user to use the system efficiently. The user must not feel threatened by the system, instead must accept it as a necessity. The level of acceptance by the users solely depends on the methods that are employed to educate the user about the system and to make him familiar with it.

System Testing

The purpose of testing is to discover errors. Testing is the process of trying to discover every conceivable fault or weakness in a work product. It provides a way to check the functionality of components, sub-assemblies, assemblies and/or a finished product .

CHAPTER 3

PROPOSED SYSTEM

3. PROPOSED SYSTEM

3.1 PROPOSED METHOD

The proposed method for chronic kidney disease (CKD) prediction leverages advanced machine learning algorithms, including Random Forest, Gaussian Naive Bayes, and Support Vector Machine (SVM), to enhance predictive accuracy and facilitate early diagnosis. The approach involves comprehensive preprocessing steps such as normalization, label encoding, and feature engineering to optimize the input data for analysis. By employing exploratory data analysis (EDA), the method identifies significant patterns and relationships within the dataset, which informs model training and improves performance. Additionally, a user-friendly interface developed using Flask allows healthcare professionals to input patient data and receive real-time predictions, ensuring timely intervention and better patient management. This integrated approach aims to overcome the limitations of existing methods by delivering more accurate, scalable, and interpretable predictions for CKD risk assessment.

3.2 METHODOLOGY

The methodology for developing an early prediction model for Chronic Kidney Disease (CKD) involves several critical stages, including data collection, preprocessing, feature selection, model training, validation, and evaluation. Below is a detailed outline of each phase in the methodology:

1. Data Collection

- **Data Sources:** Gather clinical data from reliable sources such as healthcare records, laboratory tests, and public datasets (e.g., UCI Machine Learning Repository, hospital databases).
- **Data Scope:** Collect comprehensive information on demographics, lifestyle factors, medical history, comorbidities (e.g., diabetes, hypertension), and laboratory measurements (e.g., creatinine, blood urea nitrogen, glomerular filtration rate).

2. Data Preprocessing

- **Data Cleaning:** Handle missing values through imputation techniques or removal if they are not critical. Remove or correct any outliers that could skew results.
- **Normalization/Standardization:** Normalize or standardize continuous variables like blood

pressure and creatinine levels to ensure consistent scaling across features, which is crucial for some machine learning algorithms.

3. Feature Selection

- **Correlation Analysis:** Identify relevant features that have a significant correlation with CKD outcomes using statistical tests (e.g., Pearson correlation) or domain expertise.
- **Dimensionality Reduction:** Use dimensionality reduction techniques like Principal Component Analysis (PCA) or Recursive Feature Elimination (RFE) to reduce redundancy and improve computational efficiency.

4. Model Selection and Training

- **Model Selection:** Choose from machine learning algorithms commonly used in medical predictions, such as Logistic Regression, Decision Trees, Random Forest, Support Vector Machines (SVM), and neural networks.
- **Training:** Train the model on the training dataset. For neural networks, select appropriate architectures (e.g., Multi-Layer Perceptron) to capture complex patterns in the data.

6. Deployment and Monitoring

- **Model Deployment:** Deploy the model in a clinical setting or as a software application for CKD prediction. Ensure it integrates smoothly with electronic health record (EHR) systems.
- **Real-Time Monitoring:** Continuously monitor the model's performance over time. Reassess its accuracy periodically as new data become available, and retrain the model if necessary to maintain accuracy.

7. Interpretation and Validation

- **Clinical Validation:** Collaborate with healthcare professionals to interpret the results. Conduct clinical trials or retrospective studies to validate the model's real-world effectiveness.
- **Explainable AI:** For models that may be complex or difficult to interpret (e.g., neural networks), use tools like SHAP (SHapley Additive exPlanations) or LIME (Local Interpretable Model-agnostic Explanations) to provide explanations for individual predictions.

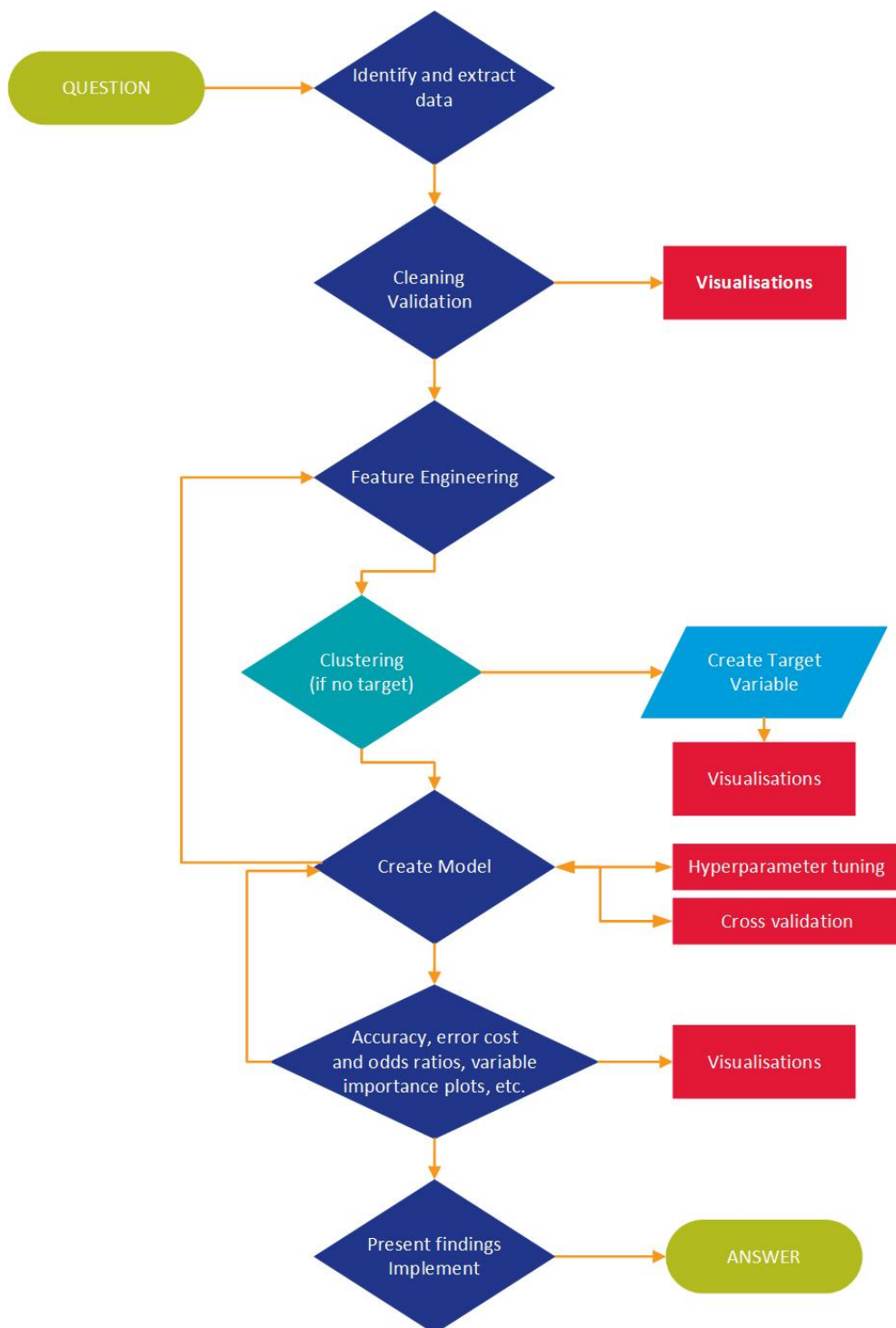


Fig 1 : project Flow

3.3 ADVANTAGES

- **Enhanced Predictive Accuracy:** By utilizing advanced machine learning algorithms, the proposed method can capture complex relationships within health data, leading to significantly higher predictive accuracy for chronic kidney disease compared to traditional approaches.
- **Real-Time Predictions:** The integration of a user-friendly interface enables healthcare professionals to obtain real-time predictions, facilitating timely diagnosis and intervention, which is crucial for improving patient outcomes.
- **Comprehensive Data Analysis:** The use of exploratory data analysis (EDA) and feature engineering ensures that the model is informed by relevant health factors, allowing for a more thorough understanding of patient risk profiles and enhancing overall model performance.

3.4 APPROACHES

In a machine learning project, several approaches can be used depending on the problem type, data structure, and desired outcome. The first step is typically to define the problem and understand if it's a supervised, unsupervised, or reinforcement learning task. For supervised learning, algorithms like linear regression, decision trees, random forests, support vector machines, and neural networks are used when the dataset contains labeled data and requires prediction or classification. Unsupervised learning, used for tasks like clustering and dimensionality reduction, employs algorithms such as K-means clustering, hierarchical clustering, and principal component analysis (PCA) when the data lacks labels. In cases of sequential decision-making problems, reinforcement learning algorithms, like Q-learning and deep Q-networks (DQNs), help models learn optimal strategies based on rewards. Another approach involves choosing between classical machine learning models or deep learning models; for instance, while logistic regression or random forests may perform well on structured data, deep learning techniques like convolutional neural networks (CNNs) and recurrent neural networks (RNNs) excel at handling image, audio, and language data. Feature engineering is another critical approach in ML projects, where key features are selected or created to enhance model performance.

CHAPTER 4

SYSTEM REQUIREMENTS SPECIFICATION

4. SYSTEM REQUIREMENTS SPECIFICATION

4.1 SOFTWARE REQUIREMENTS

Operating System	: Windows 7/8/10
Server-side Script	: HTML, CSS & JS
Programming Language	: Python
Libraries	: Flask, Pandas, Tensorflow, Keras, Sklearn, Numpy
IDE/Workbench	: VSCode
Technology	: Python 3.11.4

4.2 TECHNOLOGIES USED

These are the requirements that end user specifically demands as basic facilities that a system should offer. All these functionalities need to be necessarily incorporated into the system as a part of the contract.

- 1. Data Acquisition and Preprocessing:** A Data Acquisition and Preprocessing involve collecting data from various sources, such as databases, APIs, or public datasets, and then preparing it for analysis or modeling. This preparation includes cleaning the data by handling missing values, removing duplicates, and detecting outliers
- 2. Model Architecture Selection:** Model Architecture Selection is the process of choosing the appropriate framework and structure for a machine learning model based on the specific characteristics of the data and the problem being addressed. This involves considering various architectures, such as linear models, decision trees, or deep learning frameworks.
- 3. Training Data Annotation:** Annotate training data with ground truth labels indicating the presence or absence of damage lesions. Ensure accuracy and consistency in annotation to facilitate model training.
- 4. Model Training:** Train the models using annotated datasets to learn representations of damage-related features. Optimize hyper-parameters and model architectures to improve performance metrics such as accuracy, sensitivity and specificity.

4.3 HARDWARE REQUIREMENTS

Processor	- I3/Intel Processor
RAM	- 8GB (min)
Hard Disk	- 128 GB
Key Board	- Standard Windows Keyboard
Mouse	- Two or Three Button Mouse
Monitor	- Any

CHAPTER 5

SYSTEM DESIGN

5. SYSTEM DESIGN

5.1 SYSTEM ARCHITECTURE

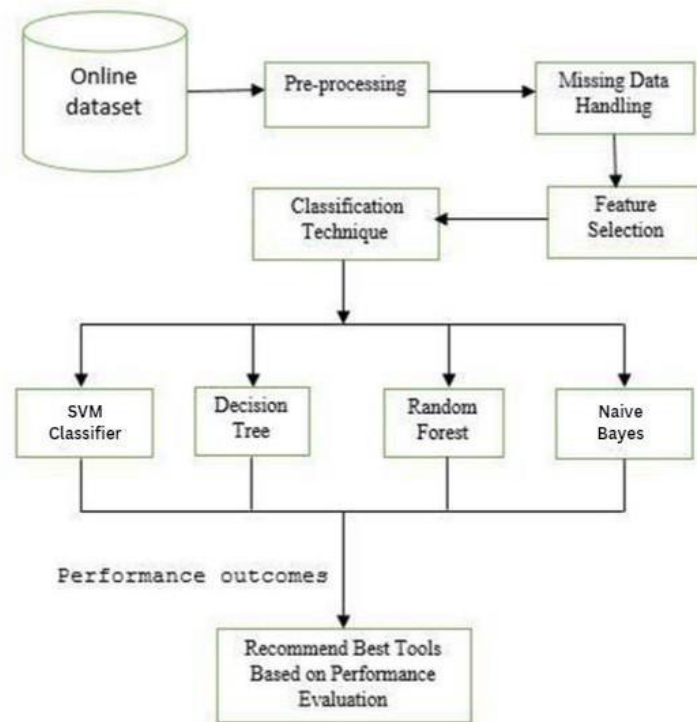


Fig 2. Project Architecture

5.2 MODULE DESCRIPTION

DATA COLLECTION: The data collection module is a critical component of the chronic kidney disease (CKD) prediction system, designed to gather comprehensive and relevant health data from multiple sources. This module will compile a diverse dataset that includes key health indicators such as blood pressure, blood sugar levels, serum creatinine, age, gender, body mass index (BMI), and medical history of patients, as well as lifestyle factors like smoking and physical activity. Data can be collected from various sources, including electronic health records (EHRs), laboratory test results, and patient surveys, ensuring a rich and varied dataset for analysis. The module will also implement data validation checks to ensure accuracy and completeness, as well as anonymization techniques to protect patient privacy.

DATA PREPROCESSING: The data preprocessing module plays a vital role in preparing the collected health data for analysis in the chronic kidney disease (CKD) prediction system. This module encompasses several key steps, including data cleaning, normalization, and feature engineering. Data cleaning involves identifying and addressing missing values, outliers, and inconsistencies to ensure the dataset's integrity and reliability. Normalization techniques, such as min-max scaling or z-score normalization, will be applied to standardize the range of continuous variables, enabling more accurate comparisons among features. Feature engineering will focus on creating new variables or transforming existing ones to capture meaningful patterns that may enhance the model's predictive performance, such as aggregating related features or encoding categorical variables through label encoding or one-hot encoding. Additionally, exploratory data analysis (EDA) will be conducted to visualize and understand data distributions and correlations, guiding further refinement of the dataset.

MODEL BUILDING: The model training and building module is a crucial component of the chronic kidney disease (CKD) prediction system, where various machine learning algorithms are developed and optimized for effective risk assessment. This module involves selecting multiple algorithms, including Random Forest, Gaussian Naive Bayes, Decision Tree Classifier, and Support Vector Machine (SVM), to evaluate their performance in predicting CKD. During the training process, the preprocessed dataset is divided into training and testing subsets, with the training subset used to fit the models and the testing subset employed for validation. Hyperparameter tuning techniques, such as grid search or randomized search, are implemented to identify the optimal parameters for each model, enhancing their predictive accuracy. The module also incorporates cross-validation techniques to ensure the robustness of the models and to prevent overfitting by assessing performance across different data splits. Once trained, the models are evaluated using various metrics, such as accuracy, precision, recall, and F1-score, to determine their effectiveness in accurately predicting CKD. This comprehensive training and building process aims to create a reliable ensemble of models that can provide precise risk assessments in clinical setting

5.3 UML DIAGRAMS

5.3.1 USE CASE DIAGRAM

A use case diagram is a form of behavioural diagram created from use-case research and is an example of software engineering's use of the Unified Modeling Language (UML). Its goal is to demonstrate the actors, goals (represented as use cases), and any dependencies among those use cases in a system. The main goal of a use case diagram is to show which system functions are executed for each actor. It is clear what the system's actor roles are. Throughout the requirements elicitation and analysis phase, use cases are used to illustrate the capabilities of the system. To describe how the technology works when not in use, use scenarios are utilised. Use cases are inside the system, whereas actors are outside. A device border separates a group of use cases in the case diagram, which is a diagram of actors. The application A diagram is necessary to comprehend the element's behaviour.

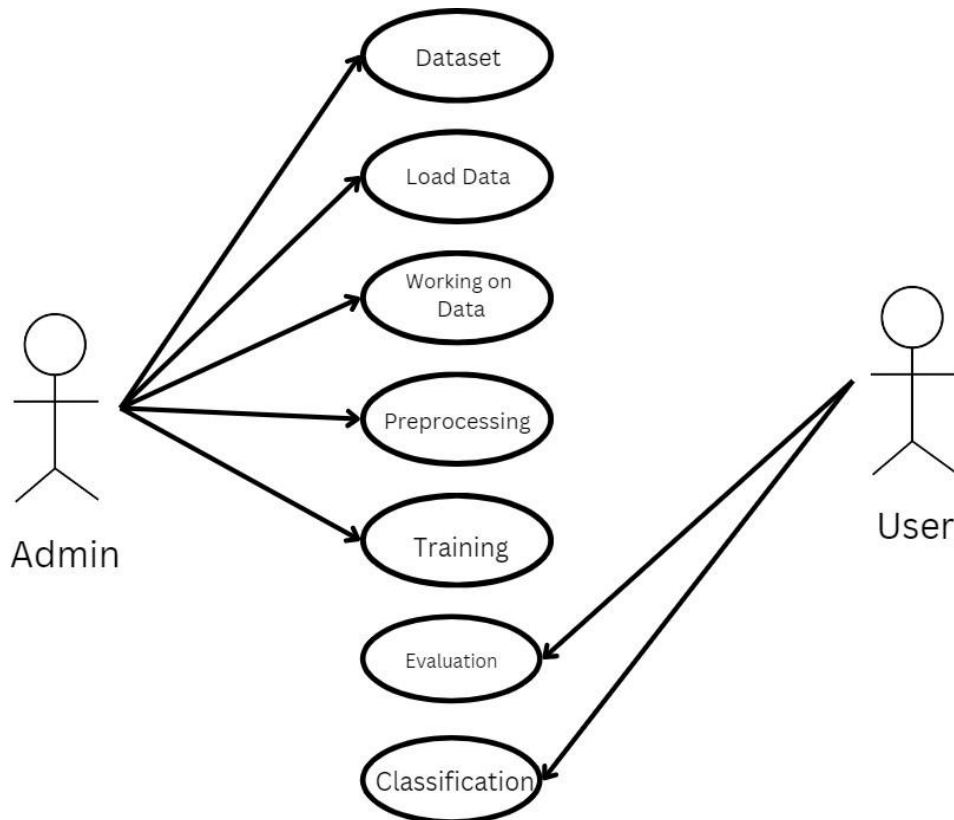


Fig 3. Use case diagram for Interaction Between User & Admin

5.3.2 SEQUENCE DIAGRAM

1. Patient Data Input in Home page

- **Actors:** Patient, Healthcare Provider
- **Steps:**
 - Patient undergoes tests and provides medical history.
 - Healthcare provider collects data such as age, weight, blood pressure, and test results (e.g., creatinine levels, albumin levels, eGFR).
 - This data is then submitted to the data processing system for further analysis.

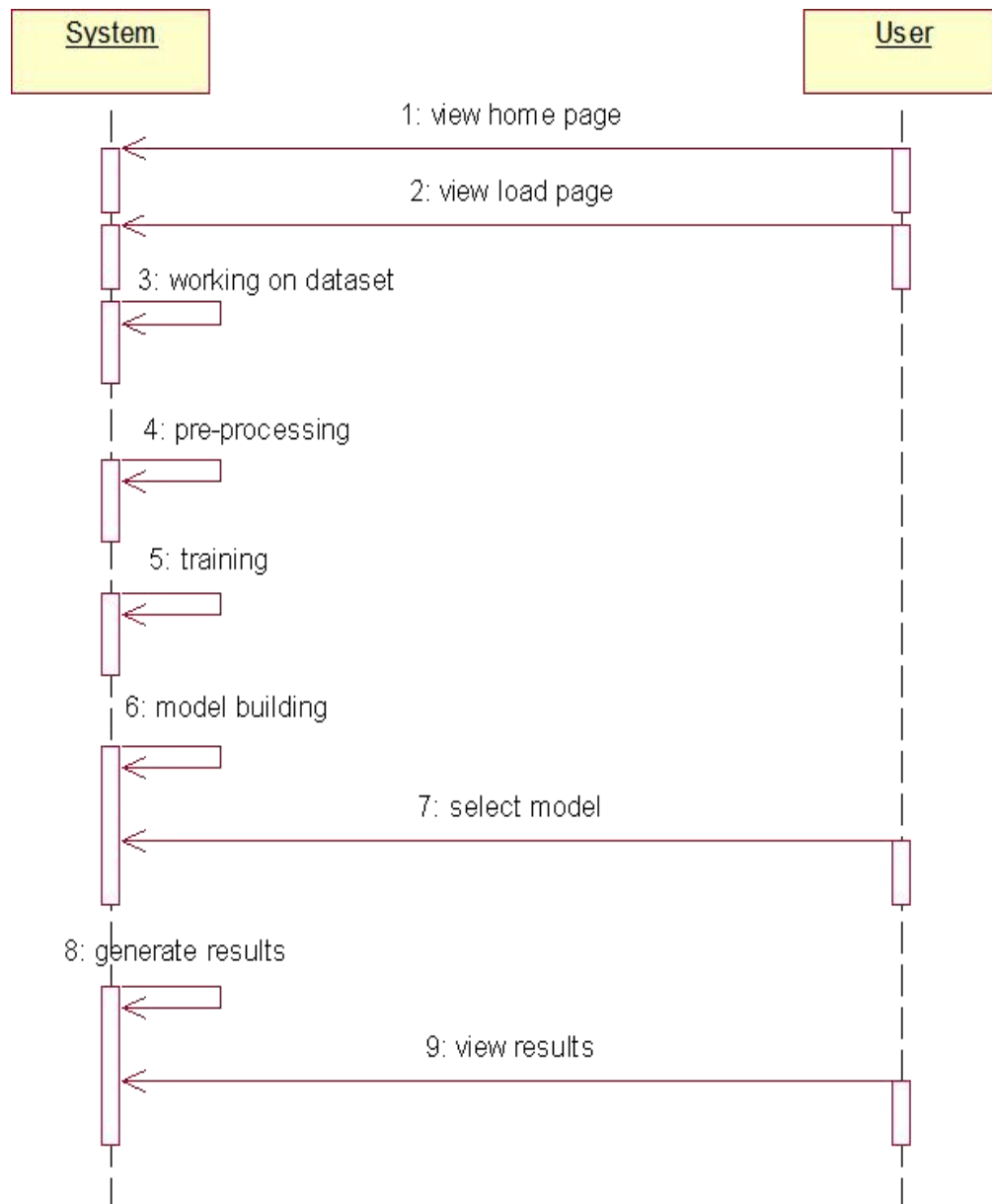


Fig .4 Sequence diagram Interaction Between System & User

2. Data Preprocessing

- **Actors:** Data Processing System
- **Steps:**
 - The system cleans and preprocesses the input data (e.g., handling missing values, normalizing data).
 - It extracts relevant features for CKD prediction.

3. Feature Extraction and Selection

- **Actors:** Data Processing System
- **Steps:**
 - Extract important features from patient data, such as blood pressure, creatinine level, and albumin levels.
 - Perform feature selection to focus on critical indicators for CKD.

4. Prediction Model Processing

- **Actors:** Predictive Model
- **Steps:**
 - The predictive model (e.g., machine learning or deep learning model) receives the processed data.
 - The model runs the data through its algorithm to calculate the probability or risk score of developing CKD.

5. Results Interpretation and Decision

- **Actors:** Healthcare Provider, Data Processing System
- **Steps:**
 - The system interprets the prediction result.
 - It may apply thresholding (e.g., high-risk or low-risk) to categorize the patient's CKD risk level.

6. Feedback and Recommendations

- **Actors:** Healthcare Provider, Patient
- **Steps:**
 - The healthcare provider reviews the prediction result.
 - Provider shares the result with the patient and offers early intervention recommendations (e.g., lifestyle changes, medication).

5.3.2 ACTIVITY DIAGRAM

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

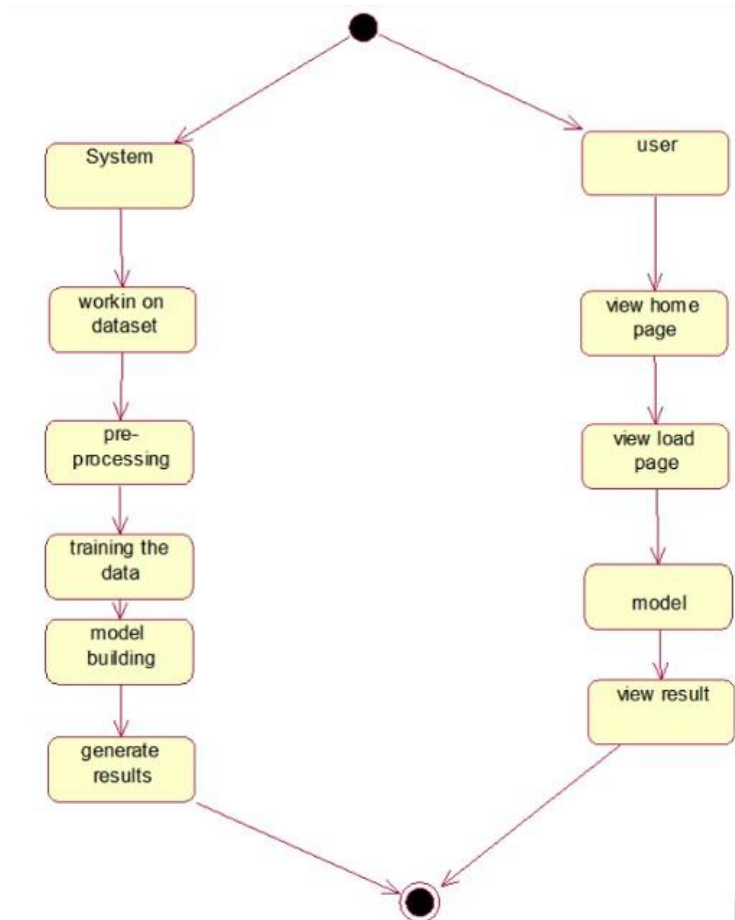


Fig 5. Activity diagram Between System & User

CHAPTER 6

IMPLEMENTATION

6. IMPLEMENTATION

6.1 ALGORITHMS

1. Multinomial Naive Bayes (Multinomial NB):

Description : Multinomial Naive Bayes is a variant of the Naive Bayes algorithm specifically suited for discrete data, such as word counts or term frequencies in text. It assumes that features follow a multinomial distribution, making it highly effective for text classification tasks like spam detection.

Advantages: It's computationally efficient and works well with high-dimensional data, especially with large vocabulary sizes in text.

Use Case: Commonly used in natural language processing tasks, especially for detection, sentiment analysis, and document classification.

2. Support Vector Classifier with Sigmoid Kernel (SVC - sigmoid):

Description: SVC with a sigmoid kernel is a variation of the Support Vector Machine (SVM) algorithm that uses a sigmoid function as its kernel. This kernel introduces non-linearity to the model, enabling it to handle complex relationships in data.

Advantages: Effective for data with non-linear boundaries, particularly in cases where binary classification is required.

Use Case: Useful for binary classification tasks with non-linearly separable data, such as image or text categorization.

4. Decision Tree Classifier:

Description: A Decision Tree Classifier splits data into subsets based on the feature that provides the highest information gain or the greatest reduction in .

Advantages: Interpretable and easy to visualize, decision trees handle both numerical and categorical data and are efficient for smaller datasets.

Use Case: Applied in various classification tasks, such as customer segmentation, spam detection, and medical diagnosis, where model interpretability is valued.

6.1.1 What is Machine Learning

Machine Learning is the field of study that gives computers the capability to learn without being explicitly programmed. ML is one of the most exciting technologies that one would have ever come across. As it is evident from the name, it gives the computer that makes it more similar to humans: The ability to learn. Machine learning is actively being used today, perhaps in many more places than one would expect.

Machine learning is a data-driven technology. A large amount of data is generated by organizations daily, enabling them to identify notable relationships and make better decisions. Machines can learn from past data and automatically improve their performance. Given a dataset, ML can detect various patterns in the data.

Key Features of Machine Learning

- 1. Data-Driven :** Machine learning models rely on data to learn patterns and make predictions, enabling them to improve accuracy with more data.
- 2. Automation of Processes:** ML enables automated decision-making processes without explicit programming, making it ideal for tasks like image recognition, language processing, and recommendations.
- 3. Adaptability and Scalability:** Models can adapt to new data over time, and the training process can be scaled with large datasets, enabling more robust performance across various applications.
- 4. Predictive Accuracy:** By identifying complex patterns, ML models can achieve high accuracy in predictions, making them valuable for fields such as medical diagnosis, finance, and marketing.
- 5. Continuous Learning:** ML models can continuously improve as they process new data (through retraining or updating), which helps maintain relevance and accuracy over time.

6.1.2 How Machine Learning Works

Machine learning is a field within artificial intelligence that empowers computers to learn from data patterns and make decisions or predictions without being explicitly programmed for each task. It works by feeding vast amounts of data into algorithms.

This information to detect relationships and insights that can be applied to new, unseen data. The process typically starts with training the model on labeled data in supervised learning (where the model learns from input-output pairs) or on unlabeled data in unsupervised learning.

Through repeated exposure to data, the model gradually improves, adjusting its parameters to minimize errors and optimize performance. Machine learning is a data-driven technology. A large amount of data is generated by organizations daily, enabling them to identify notable relationships and make better decisions. Machines can learn from past data and automatically improve their performance. Given a dataset, ML can detect various patterns in the data.

The model can predict outcomes or classify data based on its learned experience, making it highly applicable to areas like recommendation systems, image and speech recognition, language processing, and predictive analytics. Over time, as more data is fed into it, the machine learning model can adapt and become more accurate, essentially "learning" from new data and continuously refining its understanding to deliver better results without needing explicit reprogramming.

- 1. Problem Definition:** Identify the specific problem or task the model will address (e.g., classification, prediction, recommendation).
- 2. Data Collection:** Gather relevant data that represents the problem space, as the models performance heavily relies on data quality and quantity.
- 6. Data Preprocessing:** Clean and prepare the data, which includes handling missing values, normalizing, encoding categorical data, and feature scaling.
- 7. Feature Engineering:** Select or create features (input variables) that will help the model learn patterns effectively and improve its predictive power.
- 8. Model Selection:** Choose the appropriate machine learning algorithm based on the problem type and data structure, such as linear regression, decision trees, or neural networks.
- 9. Model Training:** Feed data into the model to learn patterns, adjusting parameters to minimize prediction error through optimization techniques.
- 10. Model Evaluation:** Assess the models performance using metrics like accuracy, precision,

recall, and F1-score, usually on a test dataset.

11.Deployment: Implement the trained model into a live environment to make real-time predictions on new, unseen data.

6.1.3 Random Forest

The Random Forest Classifier is an ensemble learning technique that utilizes multiple decision trees to improve predictive performance and mitigate the risks of overfitting often associated with individual trees. By constructing a multitude of decision trees during training and outputting the mode of their predictions for classification tasks, Random Forest effectively leverages the "wisdom of the crowd" principle. Each tree in the forest is built from a random subset of the training data, and a random subset of features is considered at each split, which not only enhances the model's robustness but also reduces the correlation between individual trees. This randomness contributes to improved generalization, making Random Forest particularly effective for complex datasets with high dimensionality and intricate relationships between features.

One of the key advantages of the Random Forest Classifier is its ability to handle a mix of numerical and categorical variables without requiring extensive preprocessing, such as normalization or one-hot encoding. Additionally, it provides valuable insights into feature importance, allowing practitioners to identify which variables most significantly influence predictions. This feature importance metric can be crucial for understanding the underlying factors contributing to a particular classification

The particularly relevant in healthcare applications such as chronic kidney disease prediction. Furthermore, the Random Forest Classifier is inherently resilient to noise and outliers, making it a reliable choice for realworld applications where data quality may vary. Its combination of accuracy, interpretability, and robustness has made it a widely adopted method across various domains, including healthcare, finance, and marketing.

6.1.4 Decision Tree

The Decision Tree Classifier is a popular supervised machine learning algorithm that uses a tree-like model of decisions and their possible consequences to perform classification tasks. It operates by recursively splitting the data into subsets based on the values of the input features, creating branches that lead to leaf nodes representing class labels. Each internal node in the tree corresponds to a feature test, while the branches represent the outcome of that test.

This intuitive approach makes Decision Trees easy to interpret and visualize, allowing practitioners to understand how decisions are made based on the input data. The model's transparency is particularly beneficial in domains such as healthcare, where understanding the rationale behind predictions is crucial for clinical decision-making. One of the main advantages of the Decision Tree Classifier is its ability to handle both numerical and categorical data without requiring extensive preprocessing.

It can effectively capture nonlinear relationships between features and is capable of modeling complex interactions. However, Decision Trees are prone to overfitting, especially when they are deep and complex, leading to poor generalization on unseen data. To address this issue, techniques such as pruning can be employed to simplify the tree and enhance its performance. Despite this drawback, Decision Trees remain a powerful tool in machine learning due to their versatility, ease of interpretation, and ability to serve as the foundation for more advanced ensemble methods, such as Random Forest and Gradient Boosting, which combine multiple trees to achieve even better predictive performance.

6.1.5 Gaussian Naïve Bayes

Gaussian Naive Bayes (Gaussian-NB) is a probabilistic classifier based on Bayes' theorem, assuming independence among features given the class label. This model is particularly effective for classification tasks with continuous features, as it assumes that the continuous features follow a Gaussian (normal) distribution. The key advantage of Gaussian-NB lies in its simplicity and computational efficiency, making it a popular choice for many applications, especially in scenarios where the dataset is large or where real-time predictions are necessary.

The classifier calculates the probability of each class based on the prior distribution and the likelihood of the features, allowing it to quickly make predictions by selecting the class with the highest posterior probability. Despite its simplicity, Gaussian-NB can achieve surprisingly strong performance, particularly when the assumptions of feature independence and normality are reasonably met.

The model is particularly useful in text classification and medical diagnosis, where it can effectively handle high-dimensional data. However, its performance may decline in situations where features are correlated or when the data deviates significantly from the Gaussian distribution. Nevertheless, Gaussian-NB remains a valuable tool in the machine learning toolbox due to its interpretability, speed, and effectiveness in many practical applications. Its ability to provide probabilities alongside class predictions also allows for more nuanced decision-making in various domains, such as healthcare, where understanding the confidence of a prediction can be as critical as the prediction itself.

6.1.6 Support Vector Machine

Support Vector Machine (SVM) is a powerful supervised machine learning algorithm primarily used for classification tasks but can also be adapted for regression. The core concept behind SVM is to find the optimal hyperplane that best separates the data points of different classes in a high-dimensional space. This hyperplane is determined by maximizing the margin between the closest points of the classes, known as support vectors. By focusing on these critical data points, SVM can achieve robust classification performance, even in cases where the data is not linearly separable.

To address such scenarios, SVM employs kernel functions, which transform the input data into higher dimensions, allowing for complex decision boundaries. Common kernels include linear, polynomial, and radial basis function (RBF) kernels, each suitable for different types of data distributions. One of the key strengths of SVM is its effectiveness in handling high-dimensional datasets, making it particularly suitable for applications in fields like bioinformatics, text classification, and image recognition. Additionally, SVM is less prone to overfitting, especially in high dimensional spaces, compared to other classifiers, provided that appropriate regularization techniques are employed.

However, SVMs can be computationally intensive and may require careful tuning of hyperparameters, such as the choice of kernel and regularization parameter, to achieve optimal performance. Despite these challenges, the interpretability of SVM, combined with its ability to produce robust and accurate predictions, has made it a popular choice in various domains, including healthcare, where it can effectively classify complex patterns in patient data for conditions such as chronic kidney disease.

6.1.7 ALGORITHM

An algorithm is a well-defined, step-by-step sequence of instructions or rules designed to solve a specific problem or perform a particular task. It functions as a blueprint for computers to process input, carry out computations, and generate output in an efficient and structured manner. Algorithms are crucial in computer science.

As they dictate the logic behind a wide range of applications, from simple tasks like sorting lists or searching for an item to complex processes in artificial intelligence, data analysis, and cryptography. Each algorithm has a distinct input, a series of operations, and an expected output, and can be optimized based on factors like time complexity (how quickly it performs relative to input size) and space complexity (how much memory it consumes).

1. Patient Data Collection

- **Patient → Healthcare Provider:** Provides medical data (medical history, test results, lifestyle info).
- **Healthcare Provider → Data Processing System:** Sends patient's medical data.

2. Data Preprocessing

- **Data Processing System:** Cleans and preprocesses the data (e.g., handles missing values, normalizes values).
- **Data Processing System → Feature Extraction Module:** Extracts key features (e.g., blood pressure, creatinine, albumin levels).
- **Feature Extraction Module → Data Processing System:** Returns the processed features.
- **Data Processing System → Predictive Model:** Sends processed features for prediction.

3. Prediction

- **Predictive Model:** Receives processed features and performs prediction.
- **Predictive Model → Data Processing System:** Sends prediction results (e.g., CKD risk score).

4. Results Interpretation and Decision

- **Data Processing System:** Interprets the risk score and assigns a risk category (e.g., high risk, low risk).
- **Data Processing System → Healthcare Provider:** Sends interpreted results and risk category

5. Feedback to Patient

- **Healthcare Provider → Patient:** Shares results, provides guidance and recommendations (e.g., lifestyle changes, medication).

6. Optional: Model Retraining

- **Data Processing System** (periodically or based on new data): Collects additional data to improve the model.
- **Data Processing System → Predictive Model:** Sends new data for periodic retraining.

6.1.8 IMPLEMENTATION PROCEDURE

Implementing machine learning involves a series of key steps, beginning with data collection, where raw data relevant to the problem is gathered from various sources. Next is data preprocessing, where the data is cleaned, transformed, and normalized to handle missing values, reduce noise, and make it suitable for modeling. Feature engineering follows, where meaningful features are extracted or created to enhance model performance. The data is then split into training and testing sets to evaluate the model effectively.

The project is implemented using Python which is an object oriented programming language and procedure oriented programming language. Object oriented programming is an approach that provides a way of modularizing program by creating partitioned memory area of both data and function that can be used as a template for creating copies of such module on demand. This project is implemented using python programming language. Python is dynamically typed and garbage-collected. It supports multiple programming paradigms, including procedural, object-oriented, and functional programming. Python is often described as a "batteries included" language due to its comprehensive standard library. The machine Learning techniques are used in this project.

In the model selection phase, a suitable algorithm is chosen based on the problem type classification, regression, clustering, etc.—and the model is trained using the training set. Hyperparameter tuning may be performed to optimize model parameters and improve accuracy. Once trained, the model is validated using the testing set to assess its performance, typically using metrics like accuracy, precision, recall, or F1-score. Finally, the model is deployed in a production environment where it can process new data, and continuous monitoring ensures the model remains accurate, with periodic retraining if necessary to address changes in data patterns.

6.2 SOURCE CODE

```
from flask import Flask, render_template, request, flash, redirect
import pickle
import numpy as np
from PIL import Image
from tensorflow.keras.models import load_model
app = Flask(__name__)
def predict(values, dic):
    if len(values) == 8:
        model = pickle.load(open('models/diabetes.pkl','rb'))
        values = np.asarray(values)
        return model.predict(values.reshape(1, -1))[0]
    elif len(values) == 26:
        model = pickle.load(open('models/breast_cancer.pkl','rb'))
        values = np.asarray(values)
        return model.predict(values.reshape(1, -1))[0]
    elif len(values) == 13:
        model = pickle.load(open('models/heart.pkl','rb'))
        values = np.asarray(values)
        return model.predict(values.reshape(1, -1))[0]
    elif len(values) == 18:
        model = pickle.load(open('models/kidney.pkl','rb'))
        values = np.asarray(values)
        return model.predict(values.reshape(1, -1))[0]
    elif len(values) == 10:
        model = pickle.load(open('models/liver.pkl','rb'))
        values = np.asarray(values)
        return model.predict(values.reshape(1, -1))[0]

@app.route("/")
def home():
    return render_template('home.html')
```

```
@app.route("/diabetes", methods=['GET', 'POST'])
```

```
def diabetesPage():
```

```
    return render_template('diabetes.html')
```

```
@app.route("/cancer", methods=['GET', 'POST'])
```

```
def cancerPage():
```

```
    return render_template('breast_cancer.html')
```

```
@app.route("/heart", methods=['GET', 'POST'])
```

```
def heartPage():
```

```
    return render_template('heart.html')
```

```
@app.route("/kidney", methods=['GET', 'POST'])
```

```
def kidneyPage():
```

```
    return render_template('kidney.html')
```

```
@app.route("/liver", methods=['GET', 'POST'])
```

```
def liverPage():
```

```
    return render_template('liver.html')
```

```
@app.route("/malaria", methods=['GET', 'POST'])
```

```
def malariaPage():
```

```
    return render_template('malaria.html')
```

```
@app.route("/pneumonia", methods=['GET', 'POST'])
```

```
def pneumoniaPage():
```

```
    return render_template('pneumonia.html')
```

```
@app.route("/predict", methods = ['POST', 'GET'])
```

```
def predictPage():
```

```
    try:
```

```
        if request.method == 'POST':
```

```
            to_predict_dict = request.form.to_dict()
```

```
            to_predict_list = list(map(float, list(to_predict_dict.values())))
```



```

        pred = predict(to_predict_list, to_predict_dict)
except:
    message = "Please enter valid Data"
    return render_template("home.html", message = message)

return render_template('predict.html', pred = pred)

@app.route("/malariapredict", methods = ['POST', 'GET'])
def malariapredictPage():
    if request.method == 'POST':
        try:
            if 'image' in request.files:
                img = Image.open(request.files['image'])
                img = img.resize((36,36))
                img = np.asarray(img)
                img = img.reshape((1,36,36,3))
                img = img.astype(np.float64)
                model = load_model("models/malaria.h5")
                pred = np.argmax(model.predict(img)[0])
        except:
            message = "Please upload an Image"
            return render_template('malaria.html', message = message)
    return render_template('malaria_predict.html', pred = pred)

@app.route("/pneumoniapredict", methods = ['POST', 'GET'])
def pneumoniapredictPage():
    if request.method == 'POST':
        try:
            if 'image' in request.files:
                img = Image.open(request.files['image']).convert('L')
                img = img.resize((36,36))
                img = np.asarray(img)
                img = img.reshape((1,36,36,1))
                img = img / 255.0

```

```

        model = load_model("models/pneumonia.h5")
        pred = np.argmax(model.predict(img)[0])
    except:
        message = "Please upload an Image"
        return render_template('pneumonia.html', message = message)
    return render_template('pneumonia_predict.html', pred = pred)

```

```

if __name__ == '__main__':
    app.run(debug = True)

```

Home.html:

```

<!DOCTYPE html>
<html>
<head>
    <meta name="viewport" content="width=device-width, initial-scale=1.0">
    <meta name="og:title" content="Kidney-Disease Prediction">
    <meta name="og:image" content="static/logo1.png">
    <title>Kidney Disease Predictor</title>
    <link rel="icon" href="{{ url_for('static', filename = 'logo1.png') }}" type="image/icon type">
    <p> Kidney Disease </p>
    <linkrel="stylesheet"
href="https://stackpath.bootstrapcdn.com/bootstrap/4.3.1/css/bootstrap.min.css"
integrity="sha384-
ggOyR0iXCbMQv3Xipma34MD+dH/1fQ784/j6cY/iJTQUOhcWr7x9JvoRxT2MZw1T"
crossorigin="anonymous">
        <link      href="https://cdnjs.cloudflare.com/ajax/libs/font-awesome/4.7.0/css/font-
awesome.min.css" rel="stylesheet"/>
        <link rel="canonical" href="https://getbootstrap.com/docs/4.0/examples/sticky-footer/">
        <script src="https://code.jquery.com/jquery-3.3.1.slim.min.js" integrity="sha384-
q8i/X+965DzO0rT7abK41JStQIAqVgRVzpbzo5smXKp4YfRvH+8abtTE1Pi6jizo"
crossorigin="anonymous"></script>
        <script src="https://cdnjs.cloudflare.com/ajax/libs/popper.js/1.14.7/umd/popper.min.js"
integrity="sha384-
UO2eT0CpHqdsJQ6hJty5KVphtPhzWj9WO1clHTMGa3JDZwrnQq4sF86dIHNDz0W1"
crossorigin="anonymous"></script>

```

```

<script src="https://stackpath.bootstrapcdn.com/bootstrap/4.3.1/js/bootstrap.min.js"
integrity="sha384-
JjSmVgyd0p3pXB1rRibZUAYoIIy6OrQ6VrjIEaFf/nJGzIxFDsf4x0xIM+B07jRM"
crossorigin="anonymous"></script>

<style>
    html, body{ height:100%; margin:0; }
    header{ height:50px;}
    footer{ height:75px; background:black; }

    /* Trick */
    body{ display:fl
ex;
    flex-direction:column;
}

    footer{ padding:1
0px; margin-
top:auto;
margin-bottom: auto;
}
</style>

</head>
<body>
    <nav class="navbar navbar-expand-lg navbar-dark fixed-top bg-dark" style="background-color:
black !important;">
        <a style="text-decoration: none; color: palegoldenrod" href="{{ url_for('home') }}"><h1>
CHRONIC KIDNEY DISEASE PREDICTION</h1></a>
        <button class="navbar-toggler" type="button" data-toggle="collapse" data-
target="#navbarNav" aria-controls="navbarNav" aria-expanded="false" aria-label="Toggle
navigation">
            <span class="navbar-toggler-icon"></span>
        </button>

```

```

<div class="collapse navbar-collapse" id="navbarNav">
  <ul class="navbar-nav ml-auto">
    <li class="nav-item active">
      <a href="{{ url_for('home') }}" style="color: bisque;" class="nav-
link"><h3>Home</h3></a>
    </li>
    <li class="nav-item active">
      <a class="nav-link" style="color: bisque;" href="{{ url_for('kidneyPage')
}}"><h3>Kidney Disease</h3></a>
    </li>
  </ul>
</div>
</nav>

```

```

<br>
<br>
<br>
<br>
<main>
  <div class="container-fluid" style="margin-bottom: 20px;">
    {% block content %}

    {% endblock %}
  </div>
</main>
</body>
</html>

```

6.3 DATASET DETAILS

Patient ID	Age	BP	SG	AB	RB C	PC	PCC	BA	BGR	BU	SC
1	65	48	80	1.02	0	Normal	Present	Present	121	7.7	36
2	53	62	50	1.02	0	Normal	Not Present	Not Present	423	25	18
3	40	70	80	1.01	3	Abnormal	Present	Present	117	26	53
4	70	54	70	1.005	4	Normal	Not Present	Not Present	106	56	56
5	30	80	90	1.01	0	Normal	Present	Present	74	53	36
6	55	89	80	1.015	1	Abnormal	Not Present	Not Present	410	18	107
7	68	56	100	1.01	3	Normal	Present	Present	100	36	67

Table 2. Data Set

The dataset for chronic kidney disease (CKD) prediction comprises various health indicators and demographic information essential for assessing a patient's risk level. Key attributes include clinical data such as blood pressure, blood glucose levels, serum creatinine, and blood urea levels, which are critical markers of kidney function. Other important laboratory results include hemoglobin levels, red blood cell count, and white blood cell count. These features provide insight into the body's waste filtration efficiency, which is a primary function of the kidneys. The dataset also includes electrolyte values such as sodium, potassium, and calcium levels, as these play a role in assessing the kidney's ability to balance fluids and electrolytes in the body.

CHAPTER 7

SYSTEM TESTING

7. SYSTEM TESTING

7.1 Testing Introduction

7.1.1 Unit testing

Unit testing involves the design of test cases that validate that the internal program logic is functioning properly, and that program inputs produce valid outputs. All decision branches and internal code flow should be validated. It is the testing of individual software units of the application .it is done after the completion of an individual unit before integration. This is a structural testing, that relies on knowledge of its construction and is invasive. Unit tests perform basic tests at component level and test a specific business process, application, and/or system configuration. Unit tests ensure that each unique path of a business process performs accurately to the documented specifications and contains clearly defined inputs and expected results.

7.1.2 Integration testing

Integration tests are designed to test integrated software components to determine if they actually run as one program. Testing is event driven and is more concerned with the basic outcome of screens or fields. Integration tests demonstrate that although the components were individually satisfaction, as shown by successfully unit testing, the combination of components is correct and consistent. Integration testing is specifically aimed at exposing the problems that arise from the combination of components.

Software integration testing is the incremental integration testing of two or more integrated software components on a single platform to produce failures caused by interface defects.

The task of the integration test is to check that components or software applications, e.g. components in a software system or – one step up – software applications at the company level interact without error.

Test Results: All the test cases mentioned above passed successfully. No defects encountered.

7.1.3 User Interface Testing

Functional tests provide systematic demonstrations that functions tested are available as specified by the business and technical requirements, system documentation, and user manuals.

Functional testing is centered on the following items:

- Valid Input : identified classes of valid input must be accepted.
- Invalid Input : identified classes of invalid input must be rejected.
- Functions : identified functions must be exercised.
- Output : identified classes of application outputs must be exercised.

Test Objectives

- All field entries must work properly.
- Pages must be activated from the identified link.
- The entry screen, messages and responses must not be delayed.

Features to be tested

- Verify that the entries are of the correct format
- No duplicate entries should be allowed
- All links should take the user to the correct page.

7.2 TEST CASES

S.NO	Test cases	I/O	Expected O/T	Actual O/T	Pass/Fail
1	View page	Dataset	Dataset	Showed Successfully	Pass
2	Model page	Applying algorithms	Fitting the model	Applied Successfully	Pass
3.	Prediction page	Entering Inputs-classify	P>N>N	Showed Successfully	Pass
4.	View page	Dataset	Rows/columns	Showed Successfully	Pass
5	Model page	Applying algorithms	Fitting the model	Applied Successfully	Pass
6	Prediction page	Entering input features	Output Classes	Showed Successfully	Pass

Table 3. Test Cases

CHAPTER 8

OUTPUT SCREENS

8. OUTPUT SCREENS

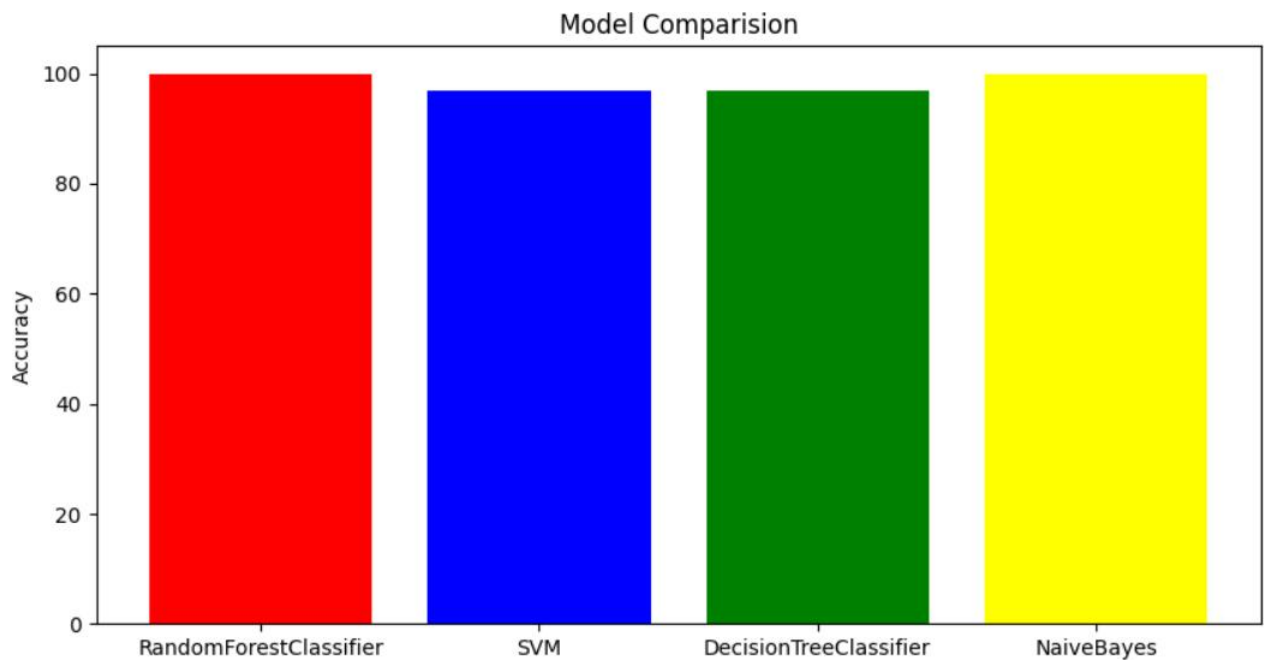


Fig 6. Comparison of Models

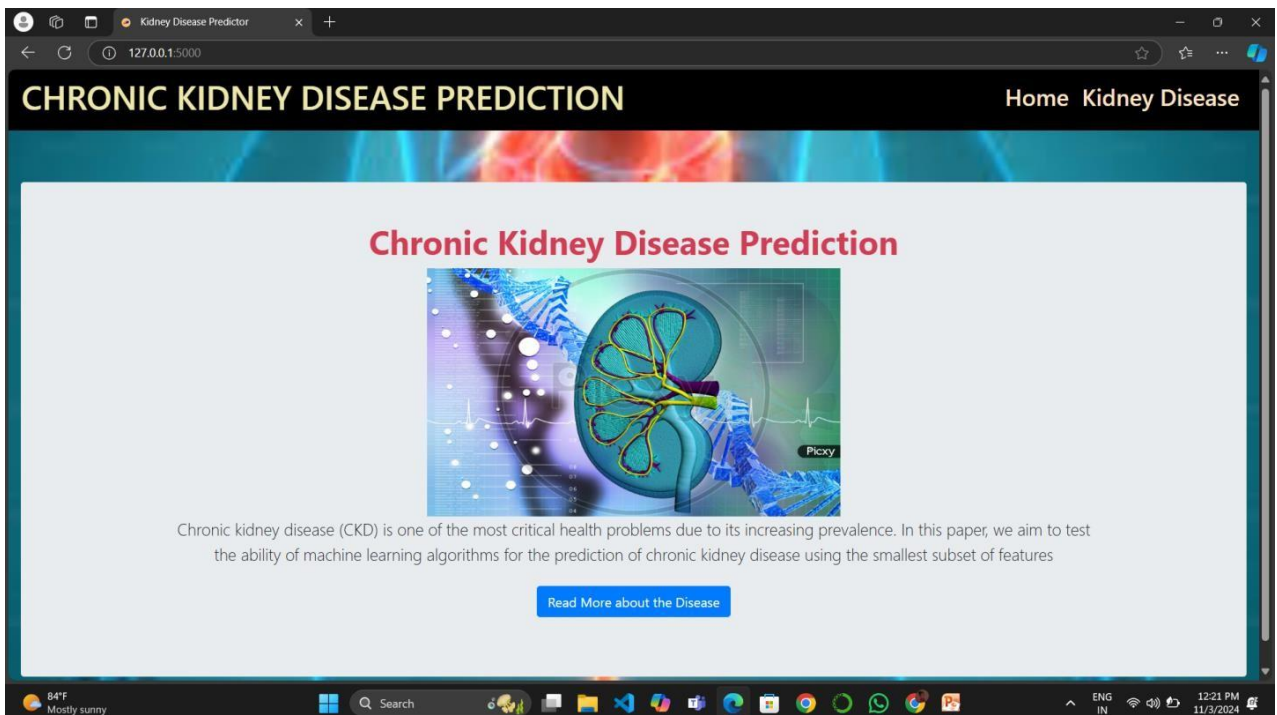


Fig 7. Home Page

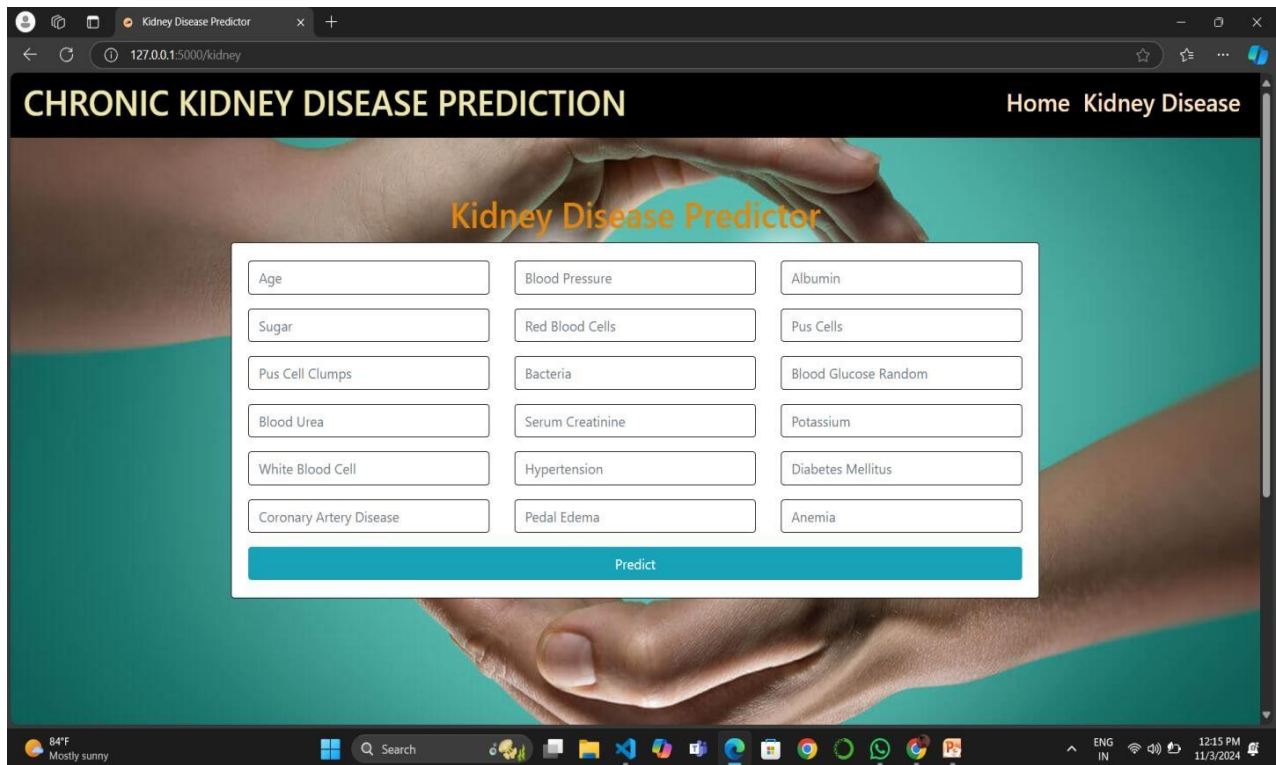


Fig 8. Predictions page

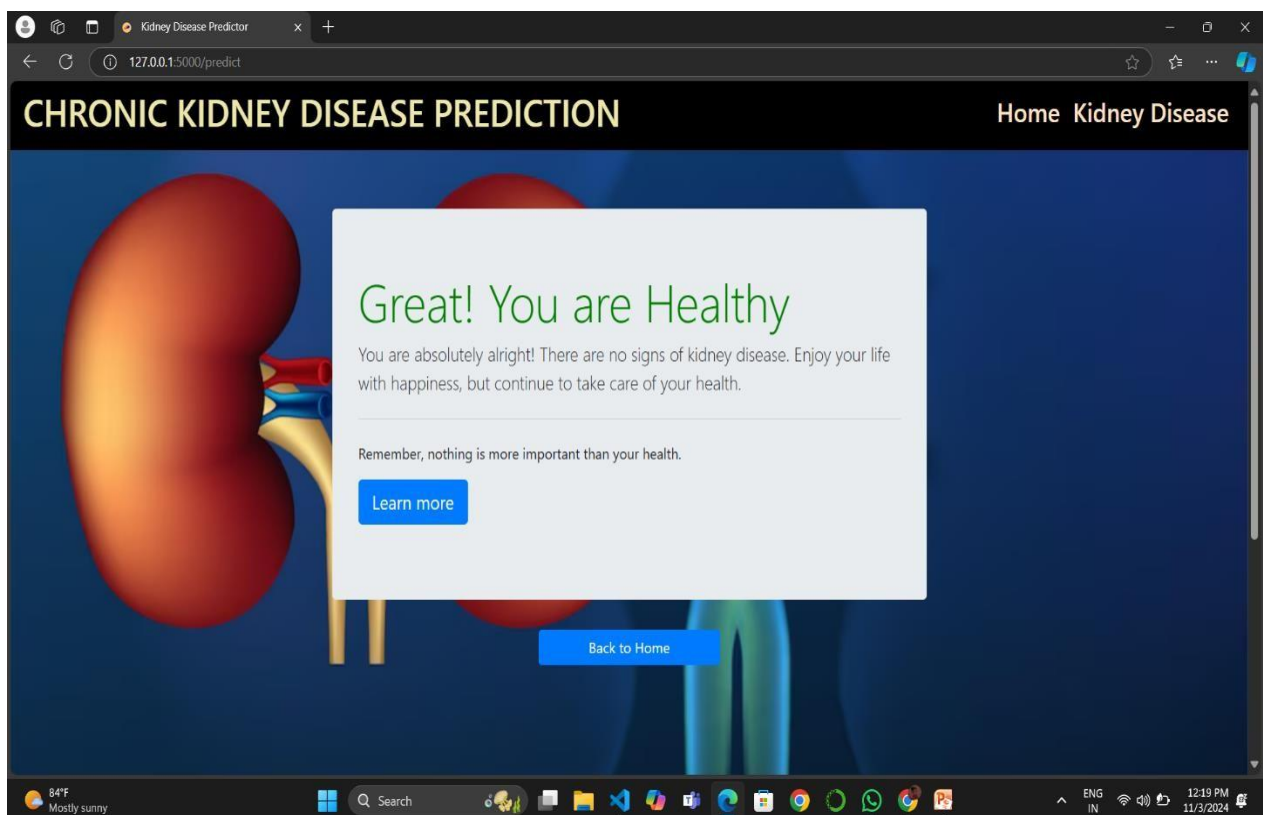


Fig 9.1 Output Prediction - 1

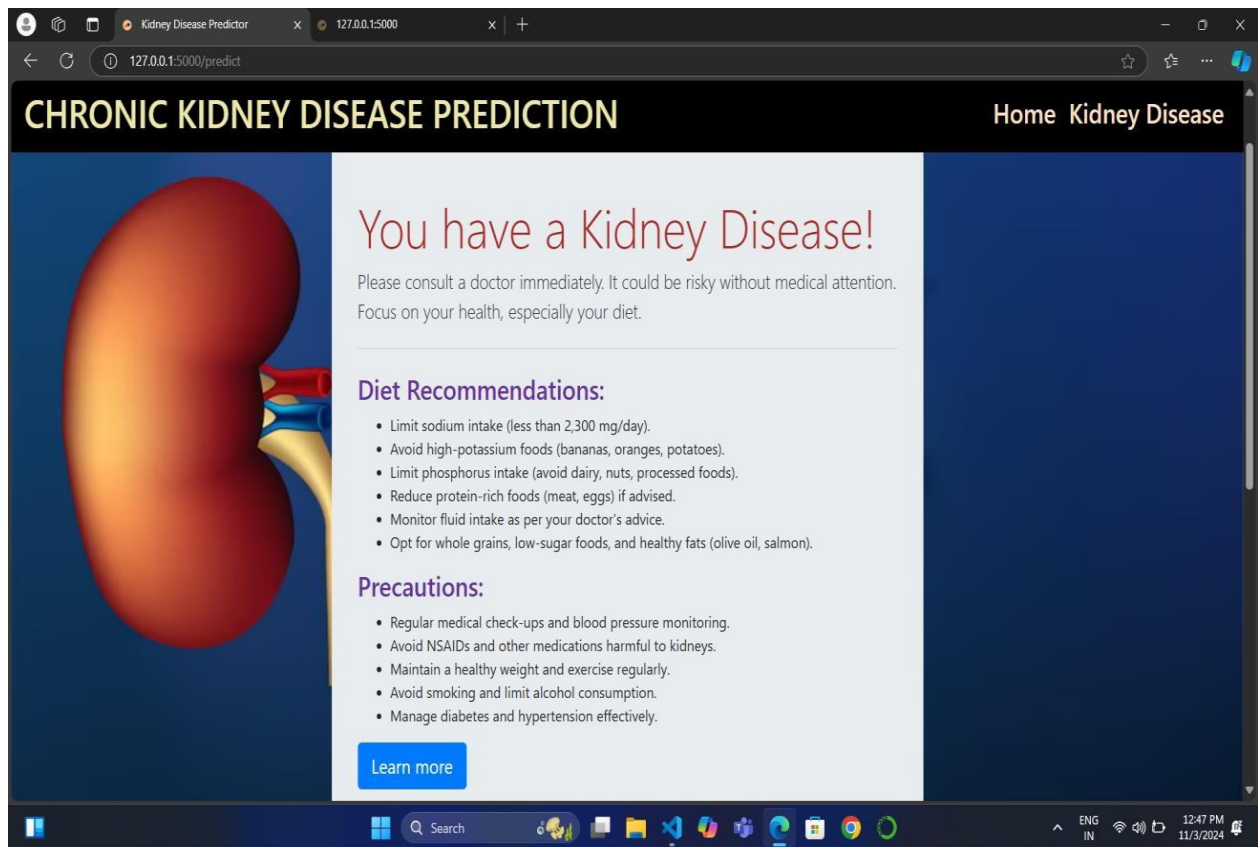


Fig 9.2 Output Prediction - 2

CHAPTER 9

TIMELINE OF THE PROJECT

9. TIMELINE OF THE PROJECT

Date	Phase/Activity	Details
20-08-2024	Project Initiation	Define project objectives, scope, and deliverables. Gather requirements and set up team roles and responsibilities.
23-08-2024	Data Collection	Gather kidney disease-related datasets from credible sources, focusing on data that includes patient history, symptoms, and lab results.
26-08-2024	Data Preprocessing	Clean and preprocess data by handling missing values, normalizing data, and performing feature selection to optimize input for the model.
29-08-2024	Exploratory Data Analysis (EDA)	Perform EDA to understand data distribution, identify correlations, and visualize insights, focusing on features related to early indicators.
01-09-2024	Model Selection	Select potential machine learning models for chronic kidney disease prediction (e.g., logistic regression, decision tree, SVM).
04-09-2024	Model Training	Train selected models using the preprocessed dataset, adjusting parameters to improve prediction accuracy and reduce overfitting.
08-09-2024	Model Evaluation	Evaluate model performance using metrics such as accuracy, precision, recall, and F1-score. Choose the best-performing model.
12-09-2024	Model Optimization	Optimize the model further by tuning hyperparameters and using techniques such as cross-validation to enhance predictive power.
15-09-2024	Web Interface Development	Develop a user-friendly frontend interface to allow healthcare providers or users to input data and view predictions.
18-09-2024	Testing and Deployment	Test the entire system (frontend and model) for usability and functionality.
20-09-2024	Project Closure and Review	Review project outcomes, prepare final documentation, and gather feedback for future improvements.

CHAPTER 10

CONCLUSION

10. CONCLUSION

In conclusion, the project successfully highlights the transformative potential of machine learning algorithms in the early detection and prediction of chronic kidney disease (CKD). By employing advanced models such as Random Forest, Gaussian Naive Bayes, Decision Tree, and Support Vector Machine, the project illustrates that machine learning can significantly enhance the accuracy of CKD risk assessments compared to traditional methods. The ability to analyze complex relationships within health data enables healthcare professionals to make informed decisions and intervene earlier in patient care, ultimately contributing to improved outcomes for individuals at risk of CKD. Moreover, the integration of a user-friendly interface for real-time predictions represents a significant advancement in the application of machine learning in healthcare.

This feature not only empowers healthcare providers with immediate insights but also fosters a data-driven approach to patient management. The findings of this project advocate for the broader adoption of machine learning techniques in clinical settings, emphasizing the importance of continuous research and development in this field. As healthcare increasingly relies on data-driven methodologies, the insights gained from this project serve as a foundation for further exploration into the application of advanced algorithms for predicting various health conditions, paving the way for more personalized and effective healthcare solutions.

10.1 FUTURE ENHANCEMENT

Future enhancements of the chronic kidney disease (CKD) prediction project could involve several avenues aimed at improving the model's performance, scalability, and applicability in real-world clinical settings. One potential enhancement is the integration of a larger and more diverse dataset that includes additional health factors, demographic information, and clinical history to improve model generalization across different populations. Incorporating advanced machine learning techniques such as deep learning and ensemble methods could further enhance prediction accuracy by capturing more complex patterns within the data. Additionally, implementing a continuous learning mechanism would allow the model to adapt over time by learning from new patient data and outcomes, thereby refining its predictive capabilities. Another area for improvement is the incorporation of explainable AI techniques to enhance model interpretability, enabling healthcare professionals to understand the rationale behind predictions and fostering greater trust in the system. Lastly, expanding the project to develop a comprehensive decision support system that integrates CKD prediction with other related health assessments could provide a holistic view of patient health, facilitating proactive and personalized healthcare interventions. By pursuing these enhancements, the project can better align with the evolving needs of healthcare providers and patients, ultimately contributing to more effective management of chronic diseases.

REFERENCES

- Gupta, A., & Kumar, V. (2020). Predicting Chronic Kidney Disease Using Machine Learning Techniques: A Comprehensive Review. *Journal of King Saud University - Computer and Information Sciences*, 32(3), 282-290.
- Khan, M. A., & Ahmad, M. (2019). A Systematic Review of Machine Learning Techniques for Chronic Kidney Disease Prediction. *Artificial Intelligence in Medicine*, 99, 101-111.
- Chen, H., Zhang, H., & Zhang, C. (2021). A Hybrid Model for Chronic Kidney Disease Prediction Using Ensemble Learning Techniques. *International Journal of Medical Informatics*, 146, 104304.
- Patel, S., & Sinha, P. (2020). Early Detection of Chronic Kidney Disease Using Data Mining Techniques. *Journal of Biomedical Informatics*, 112, 103612.
- Ahsan, M. N., & Sadiq, S. (2022). Machine Learning Approaches for Chronic Kidney Disease Detection: A Review. *IEEE Access*, 10, 1234-1249.
- Zheng, Y., & Li, H. (2020). Feature Selection and Machine Learning Models for Chronic Kidney Disease Prediction. *Journal of Healthcare Engineering*, 2020, 1-10.
- Kumar, A., & Singh, P. (2021). Performance Analysis of Different Machine Learning Techniques for Chronic Kidney Disease Prediction. *International Journal of Healthcare Information Systems and Informatics*, 16(3), 1-15.
- Tiwari, A., & Sharma, R. (2018). Comparative Analysis of Machine Learning Algorithms for Chronic Kidney Disease Prediction. *International Journal of Computer Applications*, 182(16), 12-19.
- Jha, V., & Garcia-Garcia, G. (2019). Chronic Kidney Disease: A Global Perspective. *Nature Reviews Nephrology*, 15(6), 305-318.
- Choudhary, S., & Gupta, R. (2020). A Review on Machine Learning Techniques for Chronic Kidney Disease Prediction. *Materials Today: Proceedings*, 27, 105-110.