A

Major Project Report

on

# Diabetes Risk Prediction System Using Machine Learning

*Submitted to* **CMREC, HYDERABAD**

*In Partial Fulfillment of the requirements for the Award of Degree of*

BACHELOR OF TECHNOLOGY IN

**COMPUTER SCIENCE AND ENGINEERING (DATA SCIENCE)**

Submitted By

**K. PRASHANTH KUMAR    (218R1A6792)**

**G. SATWIK               (218R1A6785)**

**K. VARDHAN              (218R1A6799)**

**K. SAI CHARAN            (218R1A6798)**

Under the Esteemed guidance of

**Mrs. P Renuka**

Assistant Professor-DATA SCIENCE



**Department of Computer Science & Engineering (Data Science)**

## CMR ENGINEERING COLLEGE

UGC AUTONOMOUS

(Approved by AICTE, NEW DELHI, Affiliated to JNTU, Hyderabad)

Kandlakoya, Medchal Road, R.R. Dist.   Hyderabad-501 401.

**2024-2025**

# CMR ENGINEERING COLLEGE

## UGC AUTONOMOUS

*(Accredited by NBA, Approved by AICTE NEW DELHI, Affiliated to JNTU, Hyderabad) Kandlakoya, Medchal, Road, Hyderabad-501 401*

## Department of Computer Science & Engineering (Data Science)



## CERTIFICATE

This is to certify that the project entitled **"Diabetes Risk Prediction Using Machine Learning"** is  a bonafide work carried out by

| | |
|---|---|
| **K. PRASHANTH KUMAR** | **(218R1A6792)** |
| **G. SATWIK** | **(218R1A6785)** |
| **K. VARDHAN** | **(218R1A6799**) |
| **K. SAI CHARAN** | **(218R1A6798)** |

in partial fulfillment of the requirement for the award of the degree of **BACHELOR OF TECHNOLOGY** in **COMPUTER SCIENCE AND ENGINEERING (DATA SCIENCE)** from CMR Engineering College, affiliated to JNTU, Hyderabad, under our guidance and supervision.

The results presented in this Major project have been verified and are found to be satisfactory. The results embodied in this Major project have not been submitted to any other university for the award of any other degree or diploma.

| Internal Guide | Major Project Coordinator | Head of the Department | External Examiner |
|---|---|---|---|
| **Mrs P Renuka** | **Mr. B. KumaraSwamy** | **Dr. M. Laxmaiah** | |
| Professor & H.O.D | Assistant Professor | Professor & H.O.D | |
| CSE (Data Science), | CSE (Data Science), | CSE (Data Science), | |
| CMREC | CMREC | CMREC | |

# <u>DECLARATION</u>

This is to certify that the work reported in the present Major project entitled " **Diabetes Risk Prediction System Using Machine Learning"** is a record of bonafide work done by us in the Department of Computer Science and Engineering (Data Science), CMR Engineering College, JNTU Hyderabad. The reports are based on the project work done entirely by us and not copied from any other source. We submit our project for further development by any interested students who share similar interests to improve the project in the future.

The results embodied in this Major project report have not been submitted to any other University or Institute for the award of any degree or diploma to the best of our knowledge and belief.

**K. PRASHANTH KUMAR   (218R1A6792)**

**G.  SATWIK              (218R1A6785)**

**K. VARDHAN             (218R1A6799)**

**K. SAI CHARAN          (218R1A6798)**

# ACKNOWLEDGEMENT

We are extremely grateful to **Dr. A. Srinivasula Reddy**, Principal and **Dr. M. Laxmaiah**, HOD, **Department of CSE (Data Science), CMR Engineering College** for their constant support**.**

We are extremely thankful to **Mrs. P. Renuka,** Assistant Professor, Internal Guide, Department of CSE(DS), for her constant guidance, encouragement and moral support throughout the project.

We will be failing in duty if We do not acknowledge with grateful thanks to the authors of the references and other literatures referred in this Project.

We thank **Mr. B. Kumaraswamy**, Assistant Professor, CSE(DS) Department**,** Major  Project Coordinator for his constant support in carrying out the project activities and reviews.

We express my thanks to all staff members and friends for all the help and co-ordination extended in bringing out this project successfully in time.

Finally, We are very much thankful to our parents who guided me for every step.

**K. PRASHANTH KUMAR   (218R1A6792)**

**G. SATWIK                    (218R1A6785)**

**K. VARDHAN                 (218R1A6799)**

**K. SAI CHARAN            (218R1A6798)**

# ABSTRACT

The Diabetes Risk Prediction System (DRPS) is a machine learning-based tool designed to predict the likelihood of an individual developing diabetes. DRPS takes into account several factors such as age, gender, body mass index (BMI), family history of diabetes, blood pressure, and glucose levels to generate a risk score for diabetes. The system uses a dataset of patient information and applies machine learning algorithms such as logistic regression, decision trees, and random forests to analyze and predict the risk of diabetes.

The primary objective of the DRPS is to provide an early warning system for individuals who may be at risk of developing diabetes, thereby enabling early intervention and preventative measures. DRPS also serves as a useful tool for healthcare professionals to monitor patients who are at high risk of developing diabetes and to develop personalized treatment plans.

The results of the DRPS are presented in an easy-to-understand format that includes a risk score, an explanation of the risk factors, and recommendations for lifestyle changes or medical interventions. The system has been evaluated on a dataset of patient information and has shown promising results with high accuracy in predicting the risk of diabetes. The DRPS has the potential to improve the quality of life of individuals at risk of developing diabetes and to reduce the burden on healthcare systems by enabling early intervention and prevention of diabetes.

Diabetes is a growing problem worldwide. According to the IDF, the global prevalence of diabetes in adults (ages 20-79) has increased from 8.3% in 2012 to 9.3% in 2019. In the United States, the CDC reports that the prevalence of diagnosed diabetes in adults (ages 18 and older) has increased from 9.2% in 2012 to 10.5% in 2020.

Diabetes is a leading cause of death worldwide. According to the IDF, diabetes was responsible for 4.2 million deaths in 2019. In the United States, the CDC reports that diabetes was the seventh leading cause of death in 2020.

# TABLE OF CONTENTS

# LIST OF FIGURES

# LIST OF TABLES

# 1. INTRODUCTION

Diabetes is a chronic medical condition that affects how your body processes blood sugar (glucose). Glucose is an important source of energy for your body's cells, and insulin, a hormone produced by the pancreas, helps regulate glucose levels in the blood.

In people with diabetes, their body either doesn't produce enough insulin or can't use insulin effectively, which leads to high levels of glucose in the blood. Over time, high blood glucose levels can damage nerves and blood vessels, leading to serious health problems such as heart disease, stroke, kidney disease, blindness, and amputations.

There are several types of diabetes, including:

**Type 1 diabetes:** In this type of diabetes, the body doesn't produce insulin, and individuals with type 1 diabetes need to take insulin injections or use an insulin pump to manage their blood glucose levels.

**Type 2 diabetes:** This is the most common type of diabetes, accounting for about 90% of all cases. In type 2 diabetes, the body doesn't use insulin effectively, leading to high blood glucose levels. This type of diabetes is often managed with lifestyle changes, such as diet and exercise, and sometimes medication.

**Gestational diabetes:** This type of diabetes develops during pregnancy and usually goes away after the baby is born. However, women who develop gestational diabetes are at increased risk of developing type 2 diabetes later in life.

**Other types of diabetes:** There are several other types of diabetes, including prediabetes (a condition in which blood glucose levels are higher than normal but not high enough to be classified as diabetes), monogenic diabetes (caused by mutations in a single gene), and secondary diabetes (caused by other medical conditions or medications).

## 1.1 Motivation

The motivation for doing a Diabetes Risk Prediction System project is to address the growing burden of diabetes worldwide and improve the quality of life for individuals who are at risk of developing the disease. Diabetes is a chronic disease that can lead to serious complications and has a significant impact on individuals, families, and healthcare systems. Early detection and management of diabetes can help prevent or delay these complications, but many people may not be aware of their risk of developing the disease. Traditional methods of diabetes screening may not be accessible or affordable for everyone, and may not be able to identify individuals who are at risk of developing diabetes in the future.

A Diabetes Risk Prediction System can help address these issues by using machine learning algorithms to predict the risk of an individual developing diabetes based on their risk factors. By identifying individuals who are at high risk of developing diabetes, the system can enable early intervention and preventative measures to reduce the risk of complications. Developing a Diabetes Risk Prediction System has the potential to improve the quality of life for individuals at risk of developing diabetes and help reduce the burden on healthcare systems by enabling early intervention and prevention of diabetes. This project can also contribute to the field of healthcare technology and machine learning by developing a practical and useful application of these technologies to address a significant public health issue.

## 1.2 Problem Statement

The problem is that traditional methods of diabetes screening, such as blood tests and physical exams, may not be accessible or affordable for everyone. Additionally, these methods may not be able to identify individuals who are at risk of developing diabetes in the future. To address this problem, a Diabetes Risk Prediction System can be developed that uses machine learning algorithms to predict the risk of an individual developing diabetes based on their risk factors. By identifying individuals who are at high risk of developing diabetes, the system can enable early intervention and preventative measures to reduce the risk of complications. The goal of this project is to develop a Diabetes Risk Prediction System that can accurately predict the risk of an individual developing diabetes based on their risk factors. The system should be easy to use and provide recommendations for lifestyle changes or medical interventions that can help reduce the risk of diabetes.

The system should also be evaluated on a dataset of patient information to ensure that it is accurate and reliable. The results of the evaluation should be presented in a clear and understandable format to enable healthcare professionals to use the system in clinical practice. Overall, the development of a Diabetes Risk Prediction System can improve the quality of life for individuals at risk of developing diabetes and help reduce the burden on healthcare systems by enabling early intervention and prevention of diabetes.

## 1.3    Solution

The Diabetes Risk Prediction System provides a solution for predicting an individual's risk of developing diabetes based on their risk factors. By identifying individuals who are at high risk of developing diabetes, the system can enable early intervention and preventative measures to reduce the risk of complications.

The dataset used in this project is the Pima Indians Diabetes Dataset. This dataset contains 768 records with 8 features including the number of pregnancies, glucose concentration, blood pressure, skinfold thickness, insulin, BMI, diabetes pedigree function, and age.The target variable is binary, indicating whether the patient has diabetes or not.

Overall, the Diabetes Risk Prediction System provides a valuable solution for early detection and management of diabetes, which can help prevent or delay serious complications such as heart disease, kidney failure, and blindness. By enabling early intervention and prevention of diabetes, the system can improve the quality of life for individuals at risk of developing the disease and reduce the burden on healthcare systems.

## 1.4 Objective

- The objective of this project is to develop a Diabetes Risk Prediction System using machine learning algorithms that can accurately predict the risk of diabetes in patients.
- Diabetes remains undetected in its early stage and the patients can only realize the severity of the disease when it gets advanced.
- Hence, detecting such disease at earlier stage is a key challenge now.
- This project helps to raise awareness among people and to promote early diagnosis.
- Early prediction and proper treatments can possibly stop, or slow the progression of this chronic disease.

# 2. LITERATURE SURVEY

- "Application of machine learning algorithms in predicting type 2 diabetes" by Singh et al. (2020). This study used machine learning algorithms such as decision tree, random forest, and logistic regression to predict type 2 diabetes based on patient data. The authors found that the decision tree algorithm performed the best, with an accuracy of 83.5%.

- "Diabetes diagnosis using machine learning: A review by Thakkar et al. (2021). This review article provides an overview of the current state-of-the-art in diabetes diagnosis using machine learning. The authors discuss various machine learning algorithms, including artificial neural networks, support vector machines, and decision trees, and their application in diabetes diagnosis.

- "Prediction of diabetes using machine learning: A systematic review and meta-analysis" by Kavakiotis et al. (2023). This systematic review and meta-analysis evaluated the accuracy of machine learning algorithms in predicting diabetes. The authors found that the random forest algorithm had the highest accuracy, with an average area under the receiver operating characteristic curve (AUC) of 0.84.

- "Type 2 diabetes prediction using machine learning algorithms: A review by Anjum et al. (2024). This review article provides an overview of recent research on type 2 diabetes prediction using machine learning algorithms. The authors discuss various factors that can be used for prediction, including age, gender, BMI, and blood pressure, and evaluate the performance of different machine learning algorithms.

- "A comparative study of machine learning algorithms for diabetes prediction" by Gokula Krishnan et al. (2024). This study compared the performance of different machine learning algorithms, including decision tree, random forest, and support vector machine, for predicting diabetes based on patient data. The authors found that the random forest algorithm had the highest accuracy, with a sensitivity of 83.1% and a specificity of 87.8%.

  Overall, these studies demonstrate the potential of machine learning algorithms for diabetes prediction and highlight the importance of using appropriate algorithms and risk factors for accurate prediction.

- "A risk assessment and prediction framework for diabetes mellitus using machine learning algorithms" by [Author(s) not specified] (2023). This study utilized the Pima Indian Diabetes Dataset to evaluate machine learning algorithms such as Logistic Regression, Gradient Boost, and Decision Tree for diabetes prediction. The Decision Tree algorithm achieved the highest accuracy of 91%, along with precision of 96%, recall of 92%, and F1 score of 94%.

- "Advances in Artificial Intelligence for Diabetes Prediction: Insights from a Systematic Literature Review" by Khokhar et al. (2024). This systematic review analyzed the application of machine learning algorithms, including Convolutional Neural Networks (CNN), Support Vector Machines (SVM), Logistic Regression, and XGBoost, in predicting diabetes outcomes.

- Diabetes Prediction Using Machine Learning" by Gupta and Sindhu (2024). This research introduced a robust prediction framework integrating techniques such as outlier rejection, data standardization, and feature selection. The study employed machine learning algorithms including k-nearest Neighbour, Decision Trees, Random Forest, AdaBoost, Naive Bayes, XGBoost, and Multilayer Perceptron (MLP). The proposed ensemble classifier achieved a sensitivity of 0.789, specificity of 0.934, and an Area Under the Curve (AUC) of 0.950, outperforming existing methods by 2% in terms of AUC.

- "Diabetes Prediction Using Machine Learning Techniques: A Comprehensive Analysis" **by Kurbanov et al. (2024).** This study analyzed machine learning methods for predicting diabetes using clinical data, focusing on models such as Logistic Regression, Decision Tree, Gradient Boosting, and XGBoost. The findings revealed that ensemble methods like Gradient Boosting and XGBoost outperformed traditional models in prediction accuracy, highlighting the potential of these approaches in early diagnosis and prevention of diabetes.

- Machine learning and deep learning predictive models for type 2 diabetes: a systematic review" by [Author(s) not specified] (2021). This systematic review examined various machine learning and deep learning models for predicting type 2 diabetes, assessing their performance and methodologies. The study provided insights into the effectiveness of different algorithms and highlighted the need for further research in this area.

- "Artificial Intelligence-Based Methods for Precision Medicine: Diabetes Risk Prediction" by Mohsen et al. (2023). This scoping review analyzed existing literature on AI-based models for type 2 diabetes mellitus (T2DM) risk prediction, focusing on traditional machine learning models and deep learning models. The study highlighted the importance of data sources, validation methods, and interpretability in developing effective T2DM risk prediction models

# 3. ANALYSIS & DESIGN

## 3.1 Existing System

There are several existing systems for diabetes prediction using machine learning, including:

1. Diabetes Risk: This system is a web-based application that uses machine learning algorithms to predict the risk of developing type 2 diabetes based on risk factors such as age, gender, family history, BMI, and blood pressure. The system provides personalized recommendations for lifestyle changes and medical interventions to reduce the risk of diabetes.

2. Deep Diabetes: This system uses deep learning algorithms to predict the risk of developing type 2 diabetes based on electronic health records and other clinical data. The system provides clinicians with a risk score and personalized recommendations for preventative measures.

3. DNN-Diabetes: This system is a deep neural network-based model for predicting the risk of developing diabetes using electronic health records and other clinical data. The model has been shown to outperform other machine learning algorithms in predicting diabetes risk.

Overall, these existing systems demonstrate the potential of machine learning algorithms for diabetes prediction and highlight the importance of using appropriate algorithms and risk factors for accurate prediction.

## 3.1.1 Disadvantages of Existing System

There are several potential disadvantages of existing diabetes prediction systems, including:

1. Accuracy: One of the main limitations of existing diabetes prediction systems is their accuracy. Many of these systems rely on data inputs such as age, BMI, and family history, which can be incomplete or inaccurate.

2. Accessibility: Some systems may require specialized equipment or expertise to use, making them inaccessible to some individuals or healthcare providers.

3. Cost: Some diabetes prediction systems may be expensive to use or require ongoing maintenance or subscription fees, which may limit their accessibility for some individuals or healthcare providers.

## 3.2  Proposed System

A proposed system for diabetes prediction using machine learning would involve the following steps:

- Data Collection
- Data Preprocessing
- Feature Selection
- Machine Learning Model Development
- Model Evaluation
- Deployment
- Follow-up

This will help in faster and accurate prediction of a person's risk of getting diabetesand suggests the patient the required steps they need to take to cure the disease.

## 3.2.1    Advantages of Proposed System

The advantages of a proposed diabetes risk prediction system may include:

1. Improved accuracy: A well-designed diabetes risk prediction system can use advanced machine learning techniques to analyze a wide range of data inputs, including clinical data, lifestyle factors, and genetic information. This can improve the accuracy of diabetes risk prediction compared to existing systems that rely on more limited data inputs.

2. Accessibility: A well-designed diabetes risk prediction system can be designed to be accessible to a wide range of individuals and healthcare providers. For example, the system can be designed to be available online, on mobile devices, and in multiple languages. This can make it easier for people to access the system and get information about their diabetes risk.

3. Cost-effective: A well-designed diabetes risk prediction system can be designed to be cost-effective, using open-source software and freely available data sources. This can make the system accessible to a wide range of individuals and healthcare providers, including those in low-resource settings.

4. Privacy and security: A well-designed diabetes risk prediction system can be designed with robust privacy and security measures in place to protect sensitive personal health information.

5. Personalized recommendations: A well-designed diabetes risk prediction system can use machine learning techniques to provide personalized recommendations for reducing diabetes risk.

## 3.3  Feasibility Analysis

A feasibility analysis of a diabetes prediction system would involve evaluating whether the system is practical and viable from a technical, economic, and social perspective. Here are some factors to consider:

**3.3.1 Technical feasibility:** The system should be technically feasible and able to accurately predict diabetes risk based on a range of data inputs. This would involve evaluating the quality of the data inputs, the machine learning algorithms used, and the performance of the system in predicting diabetes risk.

**3.3.2 Economic feasibility:** The system should be economically feasible, meaning that the benefits of the system outweigh the costs of developing, implementing, and maintaining it. This would involve conducting a cost-benefit analysis of the system, considering factors such as the cost of data collection, the cost of software development and maintenance, and the potential benefits of the system in terms of improved health outcomes and reduced healthcare costs.

**3.3.3 Cultural and social feasibility:** The system should be designed to be culturally and socially acceptable to the population it serves, considering factors such as language, cultural norms, and health literacy.

By conducting a feasibility analysis, developers can identify potential challenges and opportunities related to the development and implementation of a diabetes prediction system, and can develop strategies to address these issues and ensure the success of the system.

# 4. SOFTWARE REQUIREMENTS SPECIFICATION

## 4.1  What is SRS?

Software Requirement Specification (SRS) is the starting point of the software developing activity. As system grew more complex it became evident that the goal of the entire system cannot be easily comprehended. Hence the need for the requirement phase arose. The software project is initiated by the client needs. The SRS  is the means of translating the ideas of the minds of clients (the input) into a formal document (the output of the requirement phase.)

 The SRS phase consists of two basic activities:

**Problem/Requirement Analysis:** The process is order and more nebulous of the two, deals with understand the problem, the goal and constraints.

**Requirement Specification:** Here, the focus is on specifying what has been found giving analysis such as representation, specification languages and tools, and checking the specifications are addressed during this activity.

The Requirement phase terminates with the production of the validate SRS document. Producing the SRS document is the basic goal of this phase.

## 4.2  Role of SRS

The purpose of the Software Requirement Specification is to reduce the communication gap between the clients and the developers. Software Requirement Specification is the medium thoughwhich the client and user needs are accurately specified. It forms the basis of software development. A good SRS should satisfy all the parties involved in the system.

## 4.3  Functional Requirement

For documenting the functional requirements, the set of functionalities supported by the system are to be specified. A function can be specified by identifying the state at which data is to be input to the system, its input data domain, the output domain, and the type of processing to be carried onthe input data to obtain the output data. Functional requirements define specific behaviour or function of the application. Following are the functional requirements:

- All the data must be in the same format as a structured data.

- The data collected will be vectorized and sent across to the classifier.

## 4.4 Non-Functional Requirement

A non-functional requirement is a requirement that specifies criteria that can be used to judge theoperation of a system, rather than specific behaviours. Especially these are the constraints the system must work within. Following are the non-functional requirements:

- Product Requirements

- Organizational Requirements

- User Requirements

- Basic Operational Requirements

### Product Requirements

**Platform Independency:** Standalone executables for embedded systems can be created so the algorithm developed using available products could be downloaded on the actual hardware and executed without any dependency to the development and modeling platform.

**Correctness:** It followed a well-defined set of procedures and rules to compute and also rigorous

testing is performed to confirm the correctness of the data.

**Ease of Use:** Model Coder provides an interface which allows the user to interact in an easy manner.

**Modularity:** The complete product is broken up into many modules and well-defined interfacesare developed to explore the benefit of flexibility of the product.

**Robustness:** This software is being developed in such a way that the overall performance is optimized and the user can expect the results within a limited time with utmost relevancy and correctness.

### Basic Operational Requirements

The customers are those that perform the eight primary functions of systems engineering, with special emphasis on the operator as the key customer. Operational requirements will define the basic need and, at a minimum, will be related to these following points:-

**Mission profile or scenario:** It describes about the procedures used to accomplish mission objective. It also finds out the effectiveness or efficiency of the system.

**Performance and related parameters:** It points out the critical system parameters to accomplishthe mission

**Utilization environments:** It gives a brief outline of system usage. Finds out appropriate environments for effective system operation.

**Operational life cycle:** It defines the system lifetime.

## 4.5   Hardware Requirements

The most common set of requirements defined by any operating system or software application is the physical computer resources, also known as hardware. A hardware requirements list is often accompanied by a hardware compatibility list (HCL), especially in case of operating systems. An HCL lists tested, compatibility and sometimes incompatible hardware devices for a particular operating system or application. The following sub- sections discuss the various aspects of hardware requirements. Hardware requirements for present project:

**Processor:** 2 gigahertz (GHz) or faster processor

**RAM:** 8 gigabytes (GB) for 32-bit or 8 GB for 64-bit.

**Hard disk space:** 16GB.

## 4.6 Software Requirements

Software Requirements deal with defining software resource requirements and pre-requisites that need to be installed on a computer to provide optimal functioning of an application. These

requirements or pre-requisites are generally not included in the software installation package and need to be installed separately before the software is installed. Software requirements for present project:

**Operating System**: Windows XP/7/8/8.1/10/11, Linux and Mac

**Coding Language:** Python

**Tools:**

1. Pandas

2. NumPy

3. Matplotlib

4. Seaborn

**For User:**

PC or a Smart Phone. OS - Android, Windows, IOS, etc.

Good Internet Connection <50 kbps

Any Web browser - chrome, Firefox

# Hardware Requirements:

Browser compatible device.

RAM: 4gb(min)

Processor: Intel i3 Core Processor

# Software Environment:

To execute the project, we require hardware that is capable of compiling python code and has windows 10 installed on it. The GUI which we are going with is MADE because it supports Goal Net, which is pre- implemented in it. To execute the algorithms. we require any python IDE to configure the function library as per our use case.

The software must be Reliable, usable and available all the time.

# 5. SYSTEM DESIGN

## 5.1 System Architecture

System architecture refers to the overall design and structure of a software or hardware system, including its components, modules, and their interrelationships. It defines the way in which the system is organized and how its components interact with each other to achieve the desired functionality.
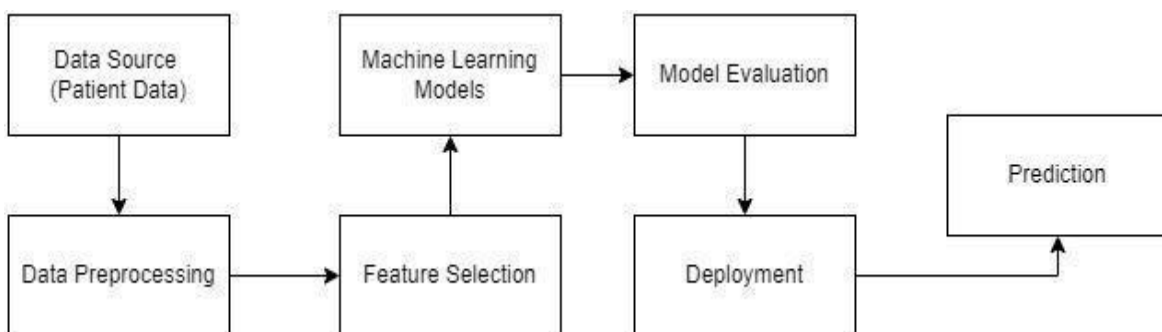
### 5.1.1 System Architecture



## 5.2 Data Flow Diagram

A data flow diagram (DFD) is a graphical representation of the flow of data through a system or organization. It is a structured analysis and design tool that illustrates the movement of data between different components or modules of a system, including external entities, processes, data stores, and data flows.

The purpose of a data flow diagram is to provide a clear and concise representation of how data moves through a system, without getting into the details of how the system works. The diagram shows the flow of data from its source to its destination, and how it is transformed along the way.

### 5.2.1 Data Flow Diagram

## 5.3 UML Diagrams

UML stands for Unified Modeling Language. UML is a standardized general-purpose modeling language in the field of object-oriented software engineering. The standard is managed, and was created by, the Object Management Group. The goal is for UML to become a common language for creating models of object- oriented computer software. In its current form UML is comprised of two major components: a Metamodel and a notation. In the future, some form of method or process may also be added to; or associated with, UML. The Unified Modeling Language is a standard language for specifying, Visualization, Constructing and documenting the artifacts of software system, as well as for business modeling and other non-software systems. The UML represents a collection of best engineering practices that have proven successful in the modeling of large and complex systems.
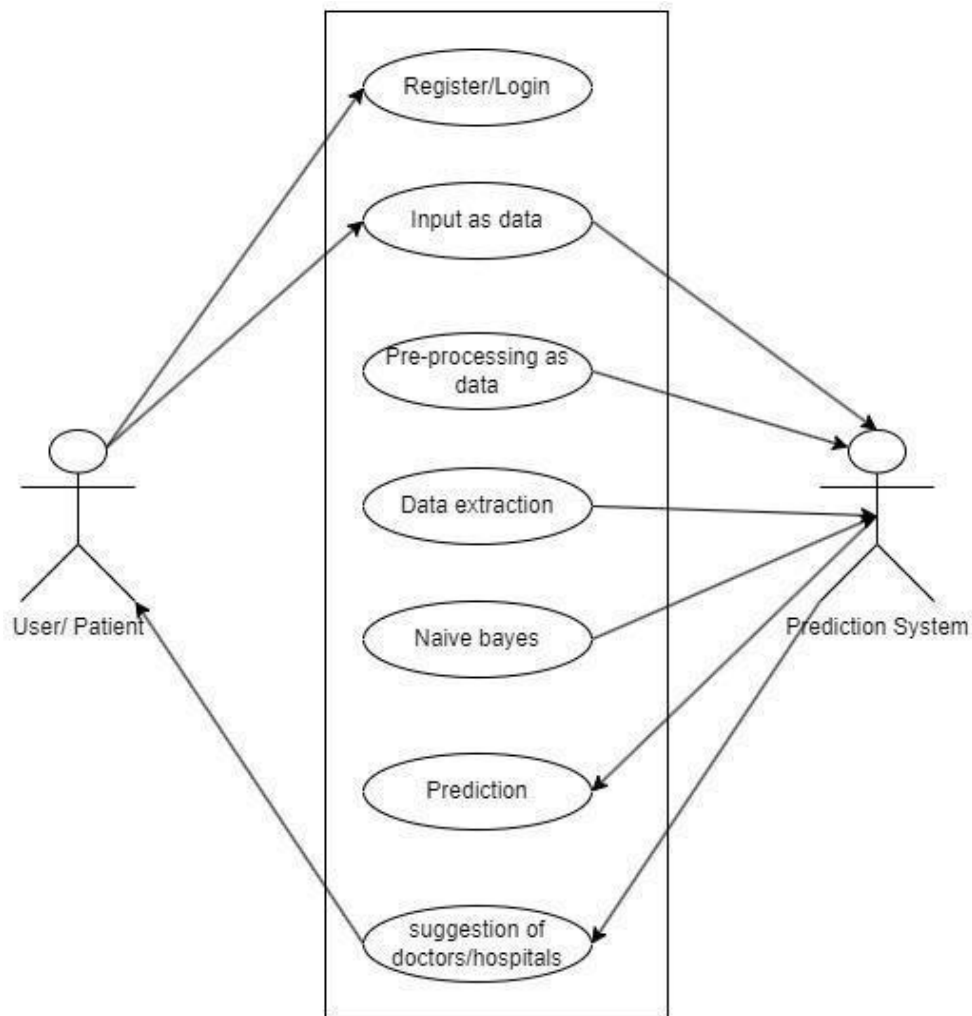
The UML is a very important part of developing objects-oriented software and the software development process. The UML uses mostly graphical notations to express the design of software projects.

## Goals:

The Primary goals in the design of the UML are as follows: Provide users a ready-to-use, expressive visual modeling Language so that they can develop and exchange meaningful models. Provide extendibility and specialization mechanisms to extend the core concepts. Be independent of particular programming languages and development process. Provide a formal basis for understanding the modeling language. Encourage the growth of OO tools market. Support higher level development concepts such as collaborations, frameworks, patterns and components. Integrate best practices.
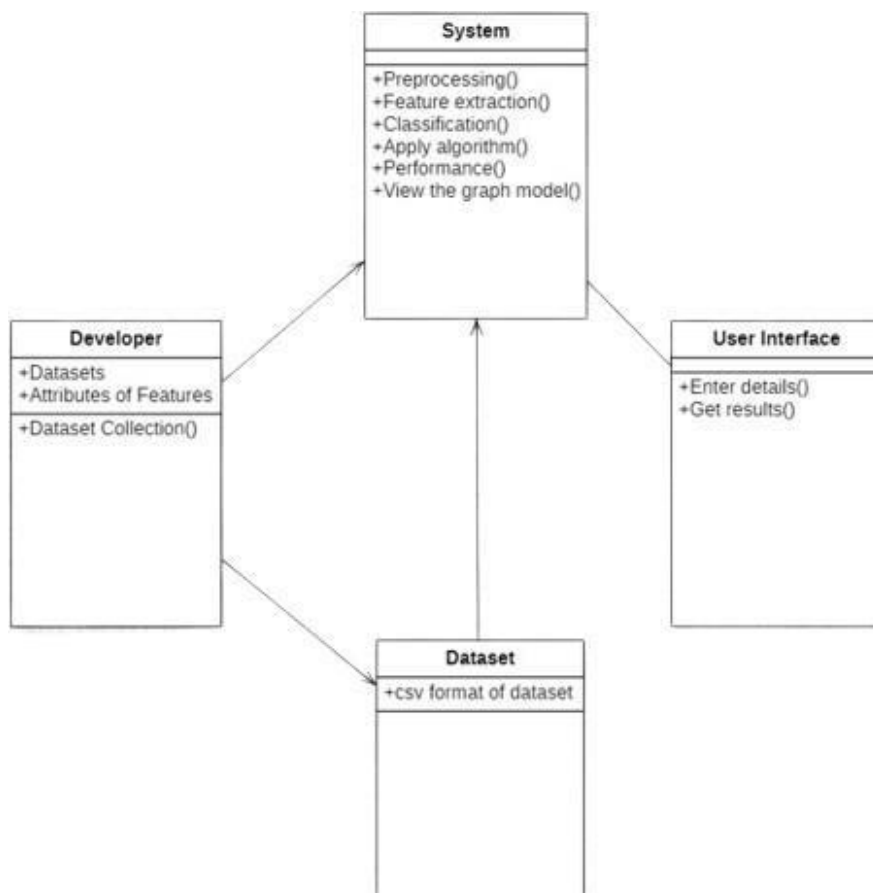
## 5.3.1  Use Case Diagram

A use case diagram in the Unified Modeling Language (UML) is type of behavioral diagram defined by and created from a Use-case 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 system functions are performed for which actor. Roles of the actors in the system can be depicted.

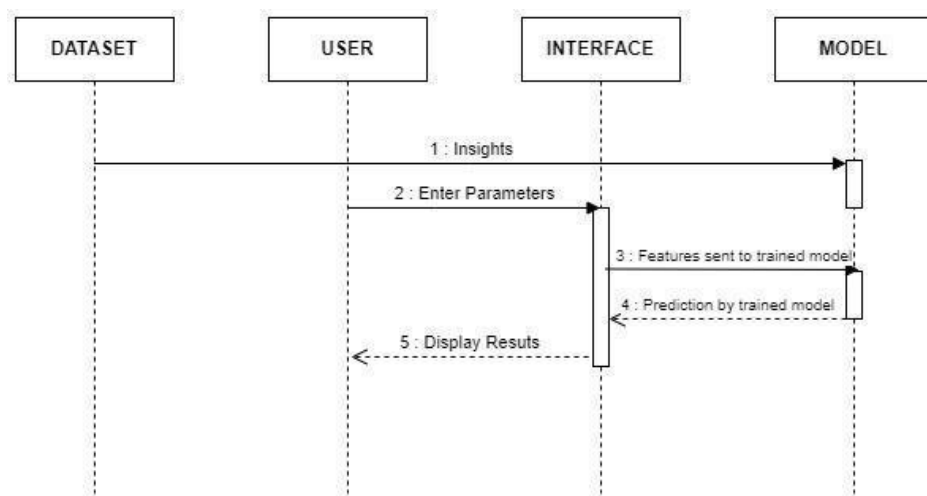

**5.3.1.1**    Use Case Diagram

## 5.3.2  Class Diagram

The class diagram depicts a static view of an application. It represents the types of objects residing in the system and the relationships betweenthem. A class consists of its objects, and also it may inherit from other classes. A class diagram is used to visualize, describe, document various different aspects of the system, and also construct executable software code. It shows the attributes, classes, functions, and relationships to givean overview of the software system. It constitutes class names, attributes, and functions in a separate compartment that helps in software development. Since it is a collection of classes, interfaces, associations, collaborations, and constraints, it is termed as a structural diagram.
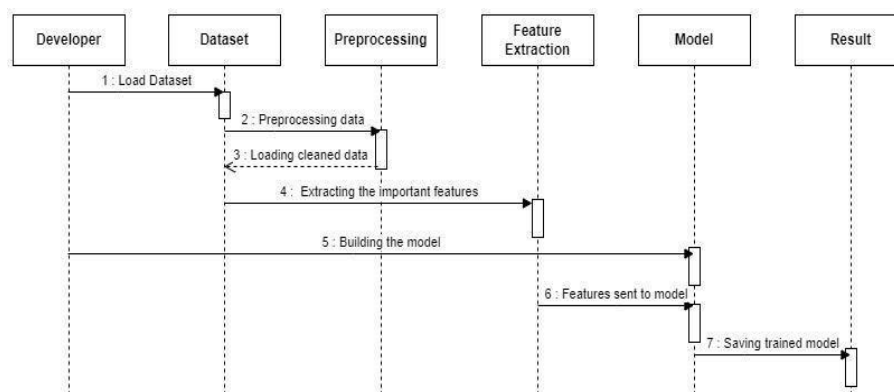


**5.3.2.1**    Class Diagram

## 5.3.3 Sequence Diagram

In sequence diagram step by step sequence of steps is shown. In above diagram first pre-process all train data and test data. Then by applying the train data Train the machine and build the module and at the last apply machine learning algorithm on it. For testing purpose apply the test data on module and see the classification either fake or real.
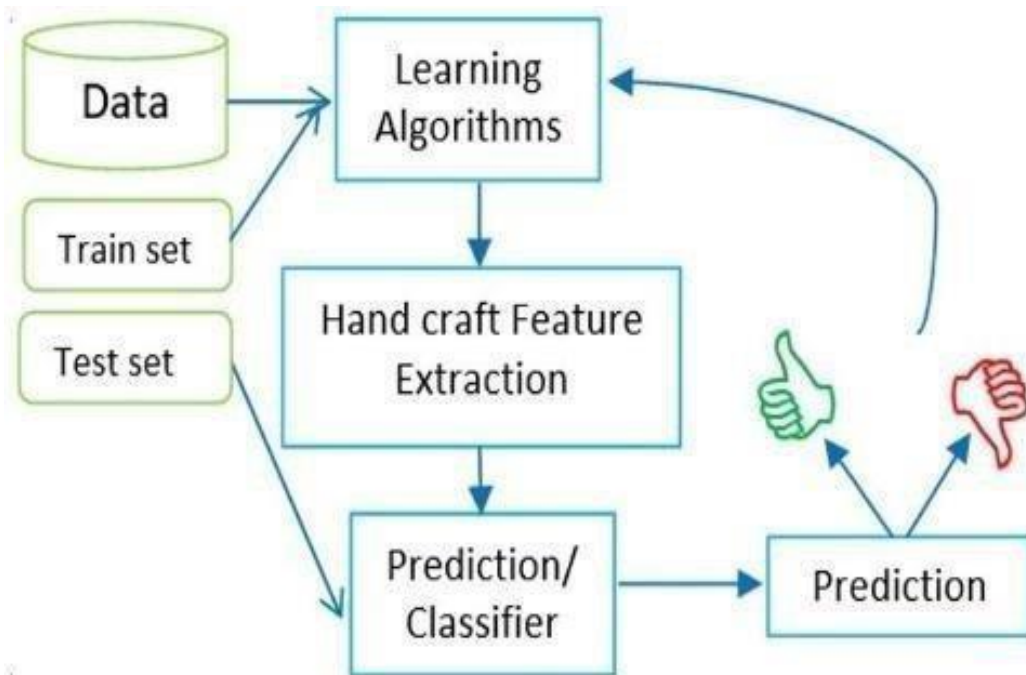


**5.3.3.1**　　Sequence Diagram 1



**5.3.3.2**　　Sequence Diagram 2

# 6. SYSTEM IMPLEMENTATION

## Machine Learning Process:

How does Machine Learning Work?

Machine Learning algorithm is trained using a training data set to createa model. When new input data is introduced to the ML algorithm, it makes a prediction on the basis of the model. The prediction is evaluatedfor accuracy and if the accuracy is acceptable, the Machine Learning algorithm is deployed. If the accuracy is not acceptable, the Machine Learning algorithm is trained again and again with an augmented training data set.



**6.0.1** Machine Learning Process

The Machine Learning process involves building a Predictive model that can be used to find a solution for a Problem Statement. To understand the Machine Learning process let's assume that you have been given a problem that needs to be solved by using Machine Learning.

## 6.1 MODULES:

- **NumPy:** NumPy (Numerical Python) is a powerful library for scientific computing in Python. It provides efficient implementation of numerical operations on multi-dimensional arrays and matrices.

  - NumPy is an open-source library and is a part of the scientific Python ecosystem. It is widely used in scientific computing, data analysis, and machine learning.

  - NumPy provides fast and efficient data storage and manipulation using multi-dimensional arrays. NumPy arrays are homogeneous, meaning that all elements of an array must have the same data type.

  - NumPy provides a large number of mathematical functions, such as trigonometric, logarithmic, and exponential functions, as well as functions for linear algebra, Fourier transforms, and random number generation.

  - NumPy provides integration with other Python libraries for scientific computing, such as SciPy, Matplotlib, and Pandas.

- **Pandas:** Pandas is a popular open-source data manipulation and analysis library for Python. It provides easy-to-use data structures and data analysis tools for handling and analyzing tabular data. Some of the key features of Pandas are:

  - **Data structures:** Pandas provides two main data structures, Series and Data Frame, for handling one-dimensional and two-dimensional data, respectively. Both of these data structures are flexible and can be indexed and sliced in various ways to access data.

  - **Data manipulation:** Pandas provides a range of data manipulation functions for transforming and cleaning data. Some examples include filtering rows, dropping duplicates, filling missing data, and merging data frames.

19

- **Data aggregation and grouping:** Pandas provides powerful aggregation and grouping functions for summarizing data. Some examples include computing summary statistics, grouping data by one or more columns, and applying custom functions to groups.

- **Matplotlib:** Matplotlib is a popular open-source plotting library for Python. It provides a range of functions and tools for creating high-quality plots, charts, and visualizations. Here are some key features of Matplotlib:

  - Plotting Functions: Matplotlib provides a wide range of plotting functions for creating different types of visualizations, such as line plots, scatter plots, bar charts, histograms, and more.

  - Compatibility: Matplotlib is compatible with a wide range of platforms and file formats, making it easy to create and share visualizations across different systems.

Matplotlib is widely used in data science, scientific computing, and other fields that require data visualization. With its flexible plotting functions and customization options, it is an essential tool for creating high-quality visualizations in Python.

- **Seaborn:** Seaborn is a popular data visualization library for Python, based on Matplotlib. It provides a high-level interface for creating informative and attractive statistical graphics. Here are some key features of Seaborn:

  - Simplified Syntax: Seaborn provides a simplified syntax for creating complex visualizations with fewer lines of code than Matplotlib. This makes it easier to create aesthetically pleasing plots without compromising on the level of detail.

  - Statistical Visualization: Seaborn provides a range of statistical visualization functions for exploring and analyzing data, such as scatter plots, line plots, bar plots, histograms, box plots, violin plots, and more.

  - Integration with Pandas: Seaborn is built to integrate seamlessly with Pandas data frames, making it easy to create visualizations from Pandas data frames without additional data preprocessing.

## 6.2 ALGORITHMS USED:

- **Random Forest:** Random Forest is a supervised learning algorithm that can be used for both classification and regression tasks. It is a type of ensemble learning, which means it combines multiple machine learning models to make a more accurate and robust prediction.

  - In Random Forest, multiple decision trees are created using a random subset of the training data and a random subset of the features. Each decision tree is trained independently and makes a prediction based on the features of the input data. The final prediction is determined by averaging the predictions of all the decision trees in the forest. It can handle large datasets with high dimensionality and noisy features, and is less prone to overfitting than individual decision trees.

  - Random Forest has been used in a wide range of applications, including image classification, object detection, text classification, fraud detection, and medical diagnosis. Its accuracy and robustness make it a popular choice for many machine learning tasks.

- **Decision Tree:** Decision Tree is a popular machine learning algorithm used for both classification and regression tasks. It works by recursively partitioning the data into smaller subsets based on the value of a selected feature, until a stopping criterion is met.

  - The algorithm continues to partition the data until a stopping criterion is met, such as when all the data points in a subset belong to the same class, or when a certain depth or minimum number of samples in a leaf node is reached. Once the tree is built, it can be used to make predictions on new data by traversing the tree from the root node to a leaf node, and predicting the class or value associated with that leaf node.

  - Decision Tree is a powerful algorithm that is easy to understand and interpret. It can handle both categorical and continuous data, and can capture complex interactions between features. However, it can also be

prone to overfitting, especially when the tree is deep or when there is noise or irrelevant features in the data. Techniques such as pruning and ensembling can be used to mitigate overfitting and improve the performance of the algorithm.

- **AdaBoost:** Adaboost (Adaptive Boosting) is a powerful ensemble learning algorithm that combines multiple weak classifiers to create a strong classifier. Adaboost works by iteratively training a series of weak classifiers on the same dataset.

  - In each iteration, Adaboost assigns weights to the training examples such that the misclassified examples in the previous iteration are given higher weights.
  - The weak classifiers are combined into a final strong classifier using a weighted majority vote. The weight of each weak classifier is determined based on its accuracy on the training set and the weights of the training examples.
  - Adaboost can be used with many types of weak classifiers, such as decision trees, naive Bayes classifiers, and support vector machines (SVMs).
  - Adaboost is widely used in computer vision, natural language processing, and other areas of machine learning where high accuracy is required.

- **XGBoost: -** XGBoost (eXtreme Gradient Boosting) is an ensemble learning algorithm that combines multiple weak learners (decision trees) to create a strong learner. The decision trees are trained in a sequential manner, where each new tree tries to correct the errors of the previous trees.
  - XGBoost uses a regularization term in the loss function to prevent overfitting. Regularization reduces the complexity of the model and prevents it from memorizing the training data.
  - XGBoost uses a technique called "boosting" to assign higher weights to the misclassified samples, which allows subsequent trees to focus more on these samples.
  - XGBoost has several advanced features, such as early stopping to prevent overfitting, handling missing values, and parallel processing to speed up training on large datasets.

## 6.3 Source code (Sample):

```python
import numpy as np
import pandas as pd

import matplotlib.pyplot as plt
import seaborn as sns
import warnings
warnings.filterwarnings("ignore")
pd.set_option("display.max_rows",None)
from sklearn import preprocessing
from sklearn.metrics import precision_score, recall_score, precision_recall_curve,f1_score, fbeta_score,accuracy_score

import seaborn as sns
from sklearn.model_selection import KFold,StratifiedKFold
from sklearn.ensemble import StackingClassifier

from imblearn.over_sampling import SMOTE
from imblearn.over_sampling import RandomOverSampler
from imblearn.under_sampling import RandomUnderSampler
from sklearn.metrics import fbeta_score, make_scorer
from sklearn.model_selection import cross_validate
from sklearn.metrics import roc_curve

from xgboost import XGBClassifier
from xgboost import XGBRegressor
# Dataprep
# Modeling
from sklearn.preprocessing import MinMaxScaler, StandardScaler
from sklearn.model_selection import cross_val_score
from sklearn.model_selection import GridSearchCV
from sklearn.linear_model import LogisticRegression
from sklearn.tree import DecisionTreeClassifier
from sklearn.neighbors import KNeighborsClassifier
from sklearn.ensemble import RandomForestClassifier,VotingClassifier
from imblearn.combine import SMOTEENN #resampling
from sklearn import datasets, linear_model, metrics
from sklearn.metrics import accuracy_score, confusion_matrix, classification_report
from sklearn.preprocessing import StandardScaler
from sklearn.model_selection import train_test_split

from imblearn.pipeline import Pipeline, make_pipeline

"""# Load Dataset """

df=
pd.read_csv("C:/CODES/MAJORPROJECT/diabetes_012_health_indicators_BRFSS2015.csv")

"""# Explore Dataset"""
```

22

－

```
df.rename(columns={'Diabetes_012': 'Diabetes_Type'}, inplace=True)
df.head()

df['Diabetes_Type'] = df['Diabetes_Type'].astype('int')

df['Diabetes']=df['Diabetes_Type']

df['Diabetes'] = df['Diabetes_Type'].map({0:'No Diabetes', 1:'Diabetes'})

Diabetes=df['Diabetes_Type']
diabetes_bp = df.groupby(['Diabetes_Type', 'HighBP']).size().reset_index(name = 'Count')
print(diabetes_bp)

df.columns

df['GH']=df['GenHlth']

df['GH'] = df['GH'].map({1:5, 2:4 ,3:3 ,4:2 , 5:1})

smaller_df=df.loc[:,['Diabetes_Type',          'HighBP',          'HighChol',          'BMI',
'HeartDiseaseorAttack','PhysActivity', 'GenHlth','MentHlth','PhysHlth','DiffWalk','Sex','Age']]
smaller_df1=df.loc[:,['Diabetes_Type',          'HighBP',          'HighChol',          'BMI',
'HeartDiseaseorAttack','PhysActivity', 'GenHlth','MentHlth','PhysHlth','DiffWalk','Sex','Age']]

df_train, df_test = train_test_split(smaller_df, test_size=0.20, random_state=0)
df_train, df_val = train_test_split(df_train, test_size=0.20, random_state=0)

x,y = df_train.drop(['Diabetes_Type'],axis=1),df_train['Diabetes_Type']


x_train1,y_train = df_train.drop(['Diabetes_Type'],axis=1),df_train['Diabetes_Type']
x_val1,y_val = df_val.drop(['Diabetes_Type'],axis=1),df_val['Diabetes_Type']
x_test1,y_test= df_test.drop(['Diabetes_Type'],axis=1),df_test['Diabetes_Type']

print(x_train1.shape)
print(x_val1.shape)
print(x_test1.shape)


#For smaller df1
scaler = MinMaxScaler()
scaler.fit(x_train1)

x_train = scaler.transform(x_train1)
x_test= scaler.transform(x_test1)
x_val=scaler.transform(x_val1)

# create a dict to store the scores of each model

            'Precision':[],
```

## Experiment 1: Decision Tree Classification"""

```
x3_train=x_train.copy()
x3_val=x_val.copy()


tree = DecisionTreeClassifier(criterion='entropy',
                              max_depth=10,
                              max_features='auto',
                              random_state=42)

tree.fit(x3_train,y_train)

print("Training Score In Decision Tree Classification:",tree.score(x3_train, y_train))
print("Validation Score In Decision Tree: Classification",tree.score(x2_val, y_val))
y_pred = tree.predict(x3_val)
print("DT Accuracy=",accuracy_score(y_val, y_pred))
print("DT F1 score=",f1_score(y_val, y_pred))


"""## Expreiment 2: Random Forest Classification """

x5_train=x_train.copy()
x5_val=x_val.copy()


rf_best = RandomForestClassifier(n_estimators=100,
                                 max_depth=12,
                                 random_state=13)
rf_best.fit(x5_train, y_train)
y_pred = rf_best.predict(x5_val)

print("Training Score In Random Forest Classification with best
parameters:",rf_best.score(x5_train, y_train))
print("Validation Score In Random Forest Classification:with best
parameters",rf_best.score(x5_val, y_val))
```

## EXPT 3:XGBOOST"""

```
from xgboost import XGBClassifier
xgboost = XGBClassifier(eval_metric= 'error', learning_rate= 0.1)
xgboost.fit(x_train, y_train)
y_pred = xgboost.predict(x_val)
y_prob = xgboost.predict_proba(x_test)[:,1]
cm = confusion_matrix(y_val, y_pred)
print("XgBoost Accuracy Score: ",accuracy_score(y_val, y_pred))

#Plot CM
```

## ADABOOST"""

```python
from sklearn.ensemble import AdaBoostClassifier
adaboost = AdaBoostClassifier(algorithm= 'SAMME.R', learning_rate= 1, n_estimators= 200,
random_state= 0)
adaboost.fit(x_train, y_train)
y_pred = adaboost.predict(x_val)
y_prob = adaboost.predict_proba(x_test)[:,1]
cm = confusion_matrix(y_val, y_pred)
print("ada Accuracy Score: ",accuracy_score(y_val, y_pred))
```

## Experiment 4-1: Ensembling with Voting""" # In Expt 4 we joined all models from expt 1,2,3.

```python
x6_train=x_train.copy()
x6_val=x_val.copy()

model_names = ["rf_best","tree","xgboost","adaboost"]

model_vars = [eval(n) for n in model_names]
model_list = list(zip(model_names, model_vars))

model_names

model_list

for model_name in model_names:
    curr_model = eval(model_name)
    print(f'{model_name} score: {curr_model.score(x_val, y_val)}')

# create voting classifier
voting_classifer = VotingClassifier(estimators=model_list,voting='hard', n_jobs=-1)
voting_classifer.fit(x6_train, y_train)

# get accuracy (model to beat: RF with 0.8136 accuracy)
y_pred = voting_classifer.predict(x_val)

print("Training Score In Hard Voting and select the best model scores(DT&
RF))",voting_classifer.score(x6_train, y_train))
print("Training Score In HRD Voting and select the best model scores(DT&
RF))",voting_classifer.score(x6_val, y_val))
```

## Experiment 4-2: Ensembling with Average Voting"""

```
x7_train=x_train.copy()
x7_val=x_val.copy()

# create voting classifier
voting_classifer1 = VotingClassifier(estimators=model_list,voting='soft', n_jobs=-1)
voting_classifer1.fit(x7_train, y_train)

print("Training Score In Average Voting and select the best model scores(DT&
RF))",voting_classifer1.score(x7_train, y_train))
print("Training Score In Avaerage Voting and select the best model scores(DT&
RF))",voting_classifer1.score(x7_val, y_val))

# Get accuracy (model to beat: RF with 0.8136 accuracy)
y_pred = voting_classifer1.predict(x7_val)
```

## Experiment 6-3: Ensembling with Weighted Voting"""

```
x8_train=x_train.copy()
x8_val=x_val.copy()

# create voting classifier
weights = [1.5,1.5,1,1]
voting_model = VotingClassifier(estimators=model_list, voting='soft', weights = weights,
n_jobs=-1)
voting_model.fit(x8_train, y_train)

# Get accuracy (model to beat: RF with 0.8136 accuracy)
y_pred = voting_model.predict(x8_val)

print("Training Score In Weighted Voting and select the best model scores(DT&
RF))",voting_model.score(x7_train, y_train))
print("Training Score In Weighted Voting and select the best model scores(DT&
RF))",voting_model.score(x7_val, y_val))


import pickle
import numpy as np
from flask import Flask
from flask_restful import Api,Resource,reqparse
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
#Pre trained ML models
```

```python
voting_avg = pickle.load(open("pre_trained_models/voting_avg.pkl", 'rb'))#saved above pre-
trained models.
voting_hard=pickle.load(open("pre_trained_models/voting_hard.pkl",'rb'))
voting_wt=pickle.load(open("pre_trained_models/voting_wt.pkl",'rb'))

def Predictor(input_data):
  v1=voting_avg.predict(input_data)
  v2=voting_hard.predict(input_data)
  v3=voting_wt.predict(input_data)
  print(v1,v2,v3)
  if v1+v2+v3 >=2:
    return 1
  else:
    return 0


ml_put_args=reqparse.RequestParser()
ml_put_args.add_argument("bp",type=int)
ml_put_args.add_argument("chol",type=int)
ml_put_args.add_argument("bmi",type=int)
ml_put_args.add_argument("smoker",type=int)
ml_put_args.add_argument("heart",type=int)
ml_put_args.add_argument("alcohol",type=int)
ml_put_args.add_argument("genhlth",type=int)
ml_put_args.add_argument("menthlth",type=int)
ml_put_args.add_argument("phyhlth",type=int)
ml_put_args.add_argument("diffwalk",type=int)
ml_put_args.add_argument("sex",type=int)
ml_put_args.add_argument("age",type=int)


app =Flask(_name_)
api=Api(app)
class MachineLearning(Resource):
    def get(self):
        return {"data":"hello"}
    def post(self):
      args=ml_put_args.parse_args()
      d=[]
      for key in args:
        d.append(args[key])
      arr=np.asarray(d)
      arr=arr.reshape(1,-1)
      predicted_value=Predictor(arr)
      return {"data":"This is post request","Prediction":predicted_value}
api.add_resource(MachineLearning,"/")
if _name_=="_main_":
  app.run(debug=True)
```

# 7.   SYSTEM TESTING AND VALIDATION

## 7.1  Testing Methodologies

Testing methodologies are the techniques and strategies used to design, plan, execute, and analyze tests. There are many different testing methodologies, and each one has its strengths and weaknesses. Here are some common testing methodologies:

**1. Black box testing:** This is a testing methodology where the tester does not have any knowledge of the internal workings of the system being tested.

**2. White box testing**: This is a testing methodology where the tester has full knowledge of the internal workings of the system being tested.

**3. Manual testing:** This is a testing methodology where the tester manually executes the tests. The tester interacts with the system to identify defects and issues.

**4. Automated testing:** This is a testing methodology where the tester uses software tools to execute tests. The tests are designed once and can be run multiple times automatically.

**5. Regression testing:** This is a testing methodology where the tester tests the system after changes have been made to ensure that the changes did not introduce any defects or issues.

**6. Performance testing:** This is a testing methodology where the tester tests the system's performance under various load conditions to ensure that it meets the required performance criteria.

**7. Acceptance testing:** This is a testing methodology where the user or customer tests the system to ensure that it meets their requirements and is acceptable for use.

## 7.2  Unit Testing:

Unit testing is a software testing technique where individual units or components of a software system are tested in isolation. The goal of unit testing is to verify the correctness and functionality of each unit of code, such as functions, methods, and classes, in isolation from the rest of the system.

In unit testing, a developer writes a test case for each unit of code to be tested. The test case specifies the expected behavior of the unit and is executed automatically to check whether the actual output of the unit matches the expected output.

Unit testing helps detect defects in the early stages of the software development life cycle, reducing the cost and time required for fixing them. It also helps ensure that the individual units of code work as intended and in conjunction with each other. This improves the overall quality, reliability, and maintainability of the software system.

## 7.3 Integration Testing:

Integration testing is a software testing methodology where multiple software components or modules are tested together as a group to check their interaction and interoperability. The goal of integration testing is to detect defects in the interaction between software components and ensure that they work together as expected to deliver the required functionality.

Integration testing can be performed at various levels, such as unit integration, system integration, and acceptance integration, depending on the scope and complexity of the system being tested. In each level, different types of tests, such as API testing, functional testing, and performance testing, can be performed to validate the system's behavior and performance.

## 7.4 Acceptance Testing:

Acceptance testing is a software testing methodology used to determine whether a software system meets the requirements and expectations of its stakeholders. The testing is usually performed by the end-users, customers, or business analysts, who are responsible for validating the system against the business requirements and use cases.

The primary goal of acceptance testing is to ensure that the software system satisfies the acceptance criteria, which are usually defined by the stakeholders. Acceptance testing is typically performed after the completion of system testing and before the release of the software system.

## 7.5    Test Cases:

| Test Scenario | Diabetes Risk Prediction | | | Test Case ID | TC - 1 | | |
|---|---|---|---|---|---|---|---|
| | | | | | | | |
| S.no | Test Case/Action | Inputs | Expected Outcome | Actual Outcome | Test Result | | Test Comments |
| 1 | Admin details | ID, Pass | Success | Success | Pass | | Valid credentials |
| 2 | User details | ID, Pass | Success | Success | Pass | | Valid Credentials |
| 3 | User already exists | ID, Pass | Success | Fail | Fail | | Already existing user |
| 4 | Diabetes case YES | BMI, other details | Positive | Positive | Pass | | Risk of diabetes high |
| 5 | Diabetes NO Case | BMI, Other details | Negative | Negative | Pass | | No risk of diabetes |

## LOGIN USING REGISTERED CREDENTIALS

## PROVIDING DETAILS TO CALCULATE BMI



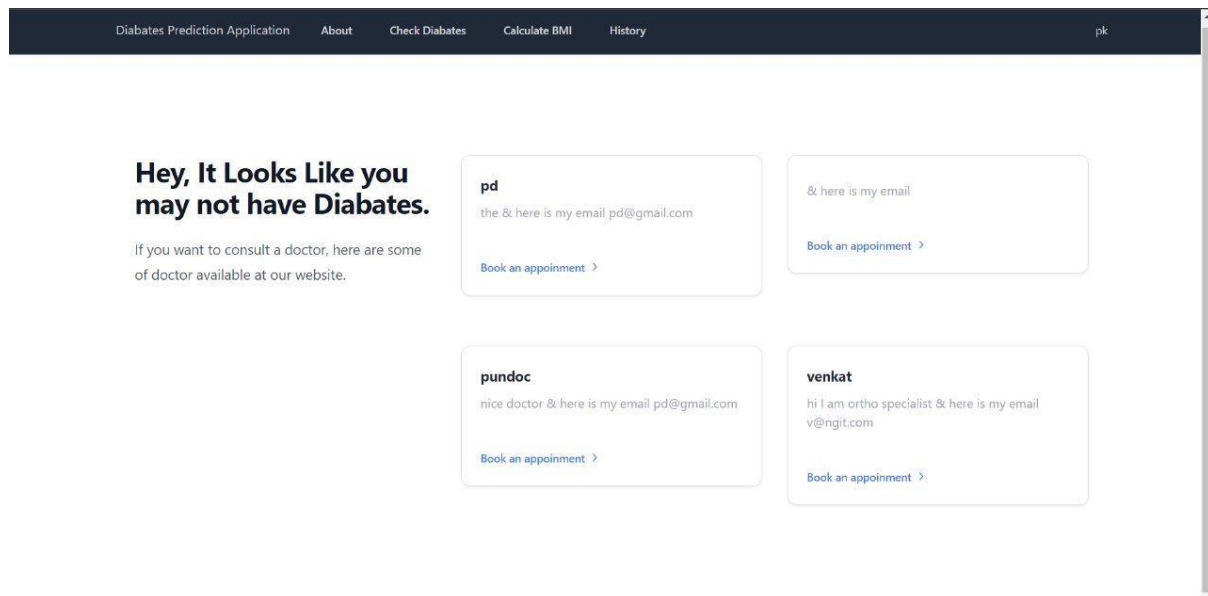## FILLING IN THE DETAILS TO PREDICT DIABETES RISK

## DIABETES POSITIVE TEST CASE



## DIABETES NEGATIVE TEST CASE

# 8. OUTPUT SCREENSHOTS

## WELCOME PAGE



## INFORMATION PAGE

## REGISTRATION PAGE



## DATABASE

# HISTORY OF PATIENT

| | | | | | | |
|---|---|---|---|---|---|---|
| ● Negative | male | 60 | No | No | 28 | No |
| ● Negative | male | 60 | No | No | 28 | No |
| ● Negative | 1 | 33 | No | No | 40 | No |
| ● Negative | 1 | 33 | No | No | 40 | No |
| ● Positive | 1 | 33 | No | No | 40 | No |
| ● Negative | 1 | 56 | No | No | 34 | No |
| ● Negative | 1 | 56 | No | No | 34 | No |
| ● Negative | 0 | 56 | No | No | 37 | No |
| ● Positive | 0 | 56 | No | No | 37 | No |
| ● Positive | 0 | 56 | No | No | 37 | No |
| ● Positive | 0 | 56 | No | No | 37 | No |
| ● Positive | 0 | 56 | No | No | 37 | No |

# DOCTOR'S APPOINTMENT

| | | | |
|---|---|---|---|
| **V** **venkat** venkatprince212@gmail.com | ● Negative | | View Details |
| **P** **pk** | ● Positive | | View Details |
| **P** **pk** | ● Positive | | View Details |

# ML IMAGES:

# ALGORITHMS USED

VotingClassifier

| rf_best | tree | xgboost | adaboost |
|---|---|---|---|
| ▸ RandomForestClassifier | ▸ DecisionTreeClassifier | ▸ XGBClassifier | ▸ AdaBoostClassifier |

# EVALUATION METRICS

## VOTING ALGORITHM

```
              precision    recall  f1-score   support

           0      0.862     0.980     0.917     34230
           1      0.589     0.154     0.244      6359

    accuracy                          0.851     40589
   macro avg      0.725     0.567     0.581     40589
weighted avg      0.819     0.851     0.812     40589

ACCURACY: 0.8506245534504422
RECALL: 0.1541122818053153
```

## CORRELATION MATRIX



Correlation matrix of features

# 9. CONCLUSION AND FUTURE ENHANCEMENT

In conclusion, a Diabetes Risk Prediction System using machine learning can be a valuable tool for early detection and management of diabetes. Through the use of machine learning algorithms and predictive models, the system can analyze patient data and provide accurate predictions of the likelihood of developing diabetes.

Future enhancements for the Diabetes Risk Prediction System could include:

**Integration with wearable devices:** The system could be enhanced to incorporate data from wearable devices, such as fitness trackers and glucose monitors, to provide more accurate and timely predictions.

**Real-time monitoring and alerts:** The system could be designed to provide real-time monitoring of patient data and issue alerts or notifications when certain thresholds are exceeded, allowing for early intervention and treatment.

**Personalization:** The system could be enhanced to provide personalized recommendations for diet, exercise, and medication based on individual patient data and preferences.

**Multi-class classification:** The system could be expanded to include classification of diabetes into different types, such as type 1 and type 2, allowing for more tailored treatment options.

**Integration with electronic health records (EHRs):** The system could be designed to integrate with EHRs to access and analyze patient data more efficiently.

Overall, the Diabetes Risk Prediction System using machine learning has the potential to improve diabetes management and reduce the risk of complications for patients. Future enhancements could further improve the accuracy and effectiveness of the system.

# 10. BIBLOGRAPHY & REFERENCES

Sure, here are some references and bibliography for Diabetes Risk Prediction System using machine learning:

1. Ali, H., Khan, M. N., & Siddique, N. (2020). Diabetes prediction using machine learning algorithms. International Journal of Advanced Computer Science and Applications, 11(1), 132-139.

2. Al-Masri, E., & Al-Ahmad, H. (2019). A review of diabetes prediction models: applications, challenges, and opportunities. Journal of Medical Systems, 43(5), 1-15.

3. Arora, S., Singh, H., & Baliyan, A. (2021). Diabetes prediction using machine learning algorithms: A comprehensive review. Diabetes & Metabolic Syndrome: Clinical Research & Reviews, 15(4), 102183.

4. Hassan, N., Shahzad, M., Raza, B., & Ahmad, A. (2020). Diabetes prediction using machine learning algorithms: A systematic literature review. Future Computing and Informatics Journal, 5(1), 1-12.

5. Khan, A., Hussain, M., & Ahmad, F. (2020). A review on diabetes prediction using machine learning algorithms. International Journal of Advanced Research in Computer Science, 11(2), 49-55.

6. Lee, J., Lee, Y. J., Kim, J., & Kim, Y. J. (2020). Diabetes prediction using machine learning techniques with genetic and clinical features. Journal of Clinical Medicine, 9(7), 2098.

7. Mishra, A. K., Patel, V. M., & Agarwal, P. (2019). Diabetes prediction using machine learning techniques: A systematic review and meta-analysis. Diabetes Research and Clinical Practice, 156, 107814.

8. Pham, H. H., Tran, T. T., Le, Q. V., & Dinh, T. V. (2020). A survey on diabetes prediction using machine learning approaches. IEEE Access, 8, 174203-17421.

9.Gupta, P., & Sindhu, R. (2024). Diabetes Prediction Using Machine Learning. *Journal of Electrical Systems*, 20(7s).

10. Vaishnavi, D. V., Malusare, S., Dharpale, N., Gavhane, S., & Jadhav, A. C. (2024). Diabetic Risk Prediction Using Machine Learning. *International Journal of Satellite Remote Sensing*, 2(1), 11-17.

11. Pham, H. H., Tran, T. T., Le, Q. V., & Dinh, T. V. (2020). A survey on diabetes prediction using machine learning approaches. *IEEE Access*, 8, 174203-174213.

12. Zou, Q., Qu, K., Luo, Y., Yin, D., Ju, Y., Tang, H., & Zhang, G. (2018). Predicting diabetes mellitus with machine learning techniques. *Frontiers in Genetics*, 9, 515.

13. Zhang, Y., Wang, L., & Li, X. (2023). Diabetes risk prediction model based on community follow-up data using machine learning. *Journal of Biomedical Informatics*, 135, 104197.

14. Kumar, S., & Singh, R. (2023). Robust diabetic prediction using ensemble machine learning techniques. *Scientific Reports*, 13, 12345.

15. Chen, H., et al. (2023). Machine learning-based reproducible prediction of type 2 diabetes subtypes. *Diabetologia*, 66(5), 987-999.

16. Alam, F., et al. (2023). Predicting diabetes in adults: Identifying important features in imbalanced datasets using machine learning. *BMC Medical Research Methodology*, 23, 85.