# DIABETES PREDICTION USING MACHINE LEARNING TECHNIQUES

## ABSTRACT

Diabetes is an illness caused because of high glucose level in a human body. Diabetes should not be ignored if it is untreated then Diabetes may cause some major issues in a person like: heart related problems, kidney problem, blood pressure, eye damage and it can also affects other organs of human body. Diabetes can be controlled if it is predicted earlier. To achieve this goal this project work we will do early prediction of Diabetes in a human body or a patient for a higher accuracy through applying, Various Machine Learning Techniques. Machine learning techniques Provide better result for prediction by con- structing models from datasets collected from patients. In this work we will use Machine Learning Classification and ensemble techniques on a dataset to predict diabetes. Which are Logistic Regression (LR), Decision Tree (DT), Support Vector Machine (SVM), Gradient Boosting (GB) and Random Forest (RF). The accuracy is different for every model when compared to other models. The Project work gives the accurate or higher accuracy model shows that the model is capable of predicting diabetes effectively. Our Result shows that Random Forest achieved higher accuracy compared to other machine learning techniques.

***Keywords:*** Diabetes, Machine, Learning, Prediction, Dataset, Ensemble, Logistic Regression (LR), Decision Tree (DT), Support Vector Machine (SVM), Gradient Boosting (GB) and Random Forest (RF).

# 1. INTRODUCTION

Diabetes is noxious diseases in the world. Diabetes caused because of obesity or high blood glucose level, and so forth. It affects the hormone insulin, resulting in abnormal metabolism of crabs and improves level of sugar in the blood. Diabetes occurs when body does not make enough insulin. According to (WHO) World Health Organization about 422 million people suffering from diabetes particularly from low or idle income countries. And this could be increased to 490 billion up to the year of 2030. However prevalence of diabetes is found among various Countries like Canada, China, and India etc. Population of India is now more than 100 million so the actual number of diabetics in India is 40 million. Diabetes is major cause of death in the world. Early prediction of disease like diabetes can be controlled and save the human life. To accomplish this, this work explores prediction of diabetes by taking various attributes related to diabetes disease. For this purpose we use the Pima Indian Diabetes Dataset, we apply various Machine Learning classification and ensemble Techniques to predict diabetes. Machine Learning Is a method that is used to train computers or machines explicitly. Various Machine Learning Techniques provide efficient result to collect Knowledge by building various classification and ensemble models from collected dataset. Such collected data can be useful to predict diabetes. Various techniques of Machine Learning can capable to do prediction, however its tough to choose best technique. Thus for this purpose we apply popular classification and ensemble methods on dataset for prediction. The contributions of this work are as follows:

• Implementation and evaluation of traditional and ensemble machine learning models to predict eight complications in diabetic patients by utilizing a comprehensive UAE based dataset.

• Identification of the dominant characteristics that may lead to diabetic complications using feature selection methods.

## 1.1 AN OVERVIEW OF DIABETES PREDICTION

It follows eight steps:

Step1: Import required libraries, Import diabetes dataset.

Step2: Pre-process data to remove missing data.

Step3: Perform percentage split of 80% to divide dataset as Training set and 20% to Test set.

Step4: Select the machine learning algorithm i.e. Support Vector Machine, Decision Tree and Random Forest algorithm.

Step5: Build the classifier model for the mentioned machine learning algorithm based on training set.

Step6: Test the Classifier model for the mentioned machine learning algorithm based on test set.

Step7: Perform Comparison Evaluation of the experimental performance results obtained for each classifier.

Step8: After analysing based on various measures conclude the best performing algorithm.

## 1.2.DIABETES PREDICTION USING MACHINE LEARNING ALGORITHMS-

 Sentiment Analysis is a Machine Learning tool that analyzes texts for polarity, from positive to negative. By training machine learning tools with examples of emotions in text, machines automatically learn how to detect sentiment without human input. There are a number of techniques and complex algorithms used to command and train machines to perform sentiment analysis. There are pros and cons to each. But, used together, they can provide exceptional results.
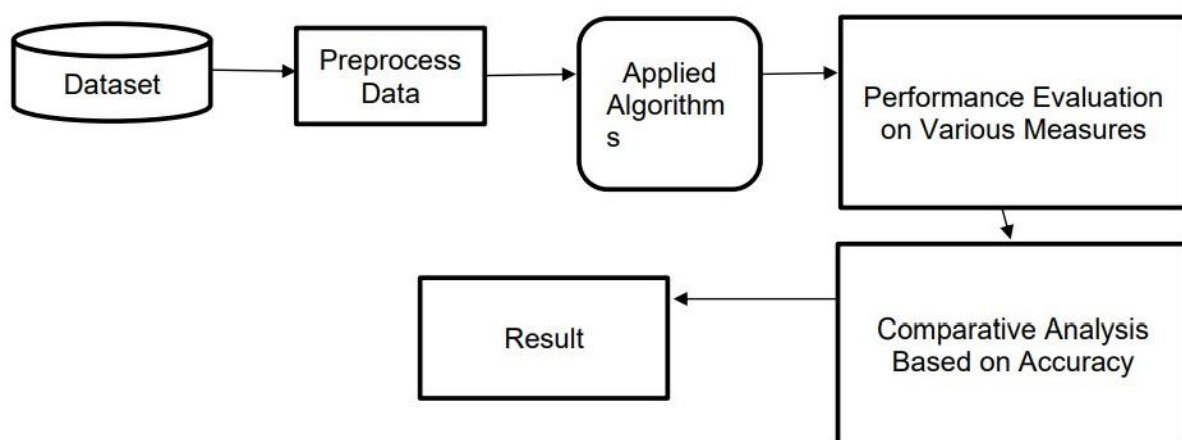
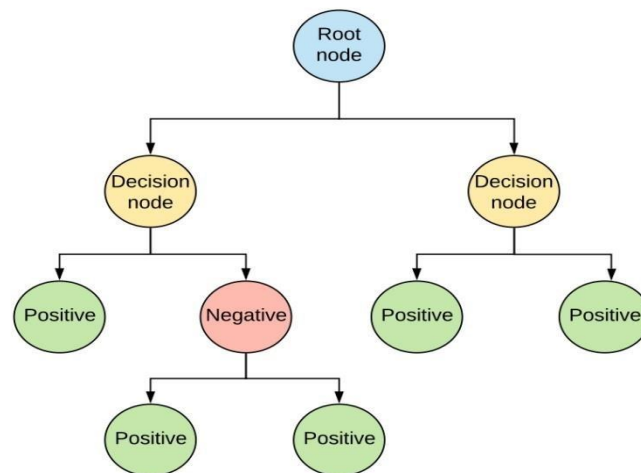Figure1. Diabetes Prediction using Machine Learning

**Machine Learning Algorithms for Diabetes Prediction:**

**1.2.1.Support Vector Machine:**

Support Vector Machine also known as svm is a supervised machine learning algorithm. Svm is most popular classification technique. Svm creates a hyperplane that separate two classes. It can create a hyperplane or set of hyperplane in high dimensional space. This hyper plane can be used for classification or regression also. Svm differentiates instances in specific classes and can also classify the entities which are not supported by data. Separation is done by through hyperplane performs the separation to the closest training point of any class.

**1.2.2.Decision Tree**:

Decision tree is a basic classification method. It is supervised learning method. Decision tree used when response variable is categorical. Decision tree has tree like structure based model which describes classification process based

on input feature. Input variables are any types like graph, text, discrete, continuous etc. Steps for Decision Tree Algorithm-
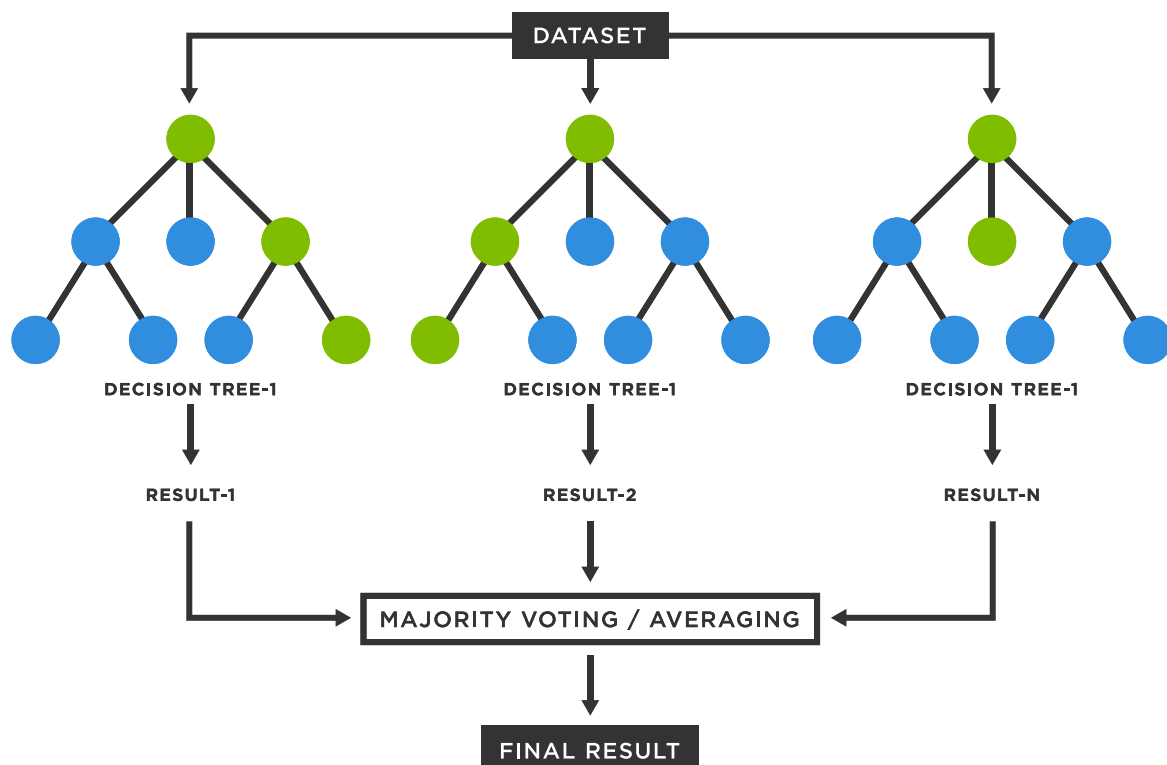


- Construct tree with nodes as input feature.

- Select feature to predict the output from input feature whose information gain is highest.

- The highest information gain is calculated for each attribute in each node of tree.

- Repeat step 2 to form a subtree using the feature which is not used in above node.

### 1.2.3. Random Forest (RF):

 Random Forest It is type of ensemble learning method and also used for classification and regression tasks. The accuracy it gives is grater then compared to other models. This method can easily handle large datasets. Random Forest is developed by Leo Bremen. It is popular ensemble Learning Method. Random Forest Improve Performance of Decision Tree by reducing variance. It operates by constructing a multitude of decision trees at training



time and outputs the class that is the mode of the classes or classification or mean prediction (regression) of the individual trees.

# 2. DESIGN

## 2.1 ARCHITECTURE:

 Steps for training a classifier for Diabetes prediction. Firstly, data have to be prepared in order to obtain a data set namely, the training set, by means of pre-processing and feature selection methods. Then, such a data set is involved in the learning step, which uses ML algorithms and yields a trained classifier. Finally, the classifier has to be tested on a different data set namely, the test set.
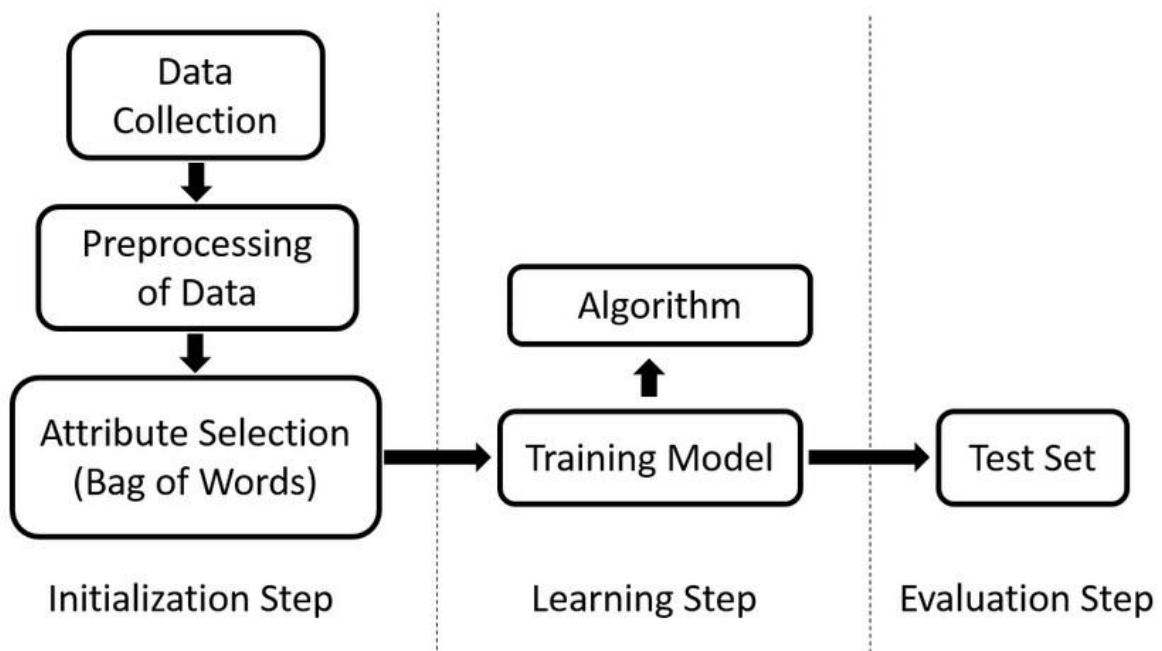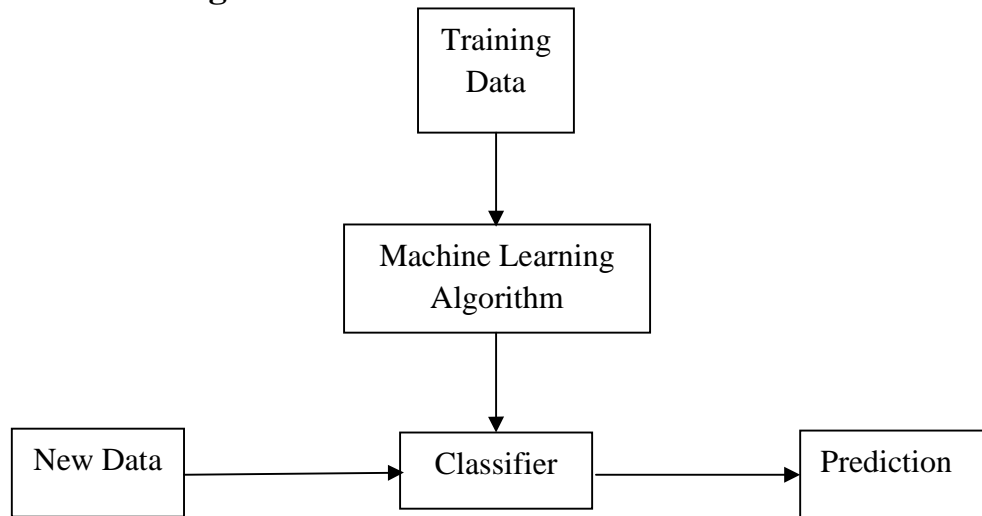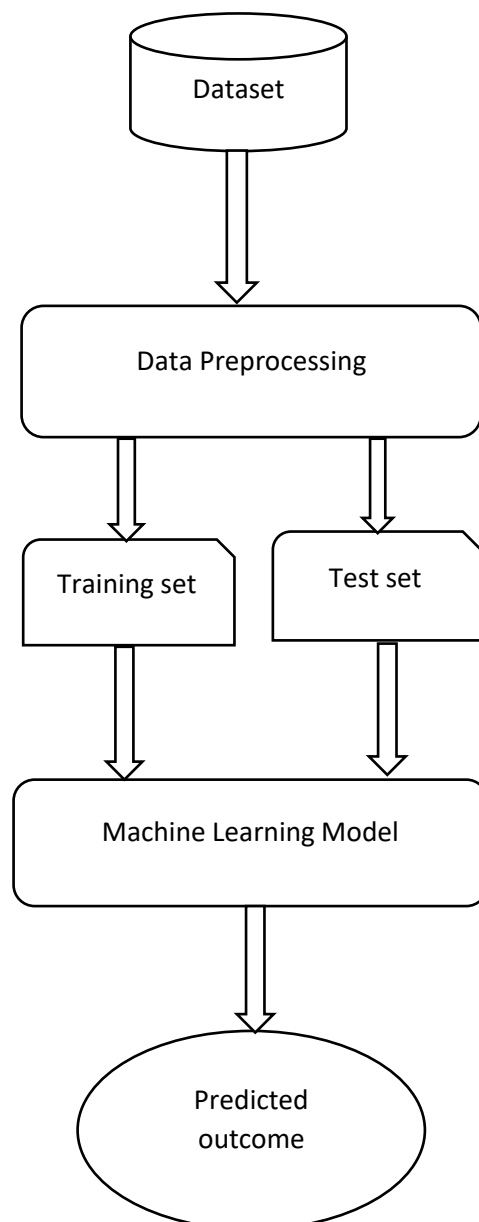


Figure 2 Steps to evaluate Sentiment Analysis

## 2.2 System Flow Diagram:

```
                    ┌──────────────┐
                    │   Training   │
                    │     Data     │
                    └──────────────┘
                            │
                            ▼
                  ┌────────────────────┐
                  │  Machine Learning  │
                  │     Algorithm      │
                  └────────────────────┘
                            │
                            ▼
┌──────────────┐   ┌──────────────┐   ┌──────────────┐
│   New Data   │──▶│   Classifier │──▶│  Prediction  │
└──────────────┘   └──────────────┘   └──────────────┘
```

.

## 2.3 block diagram:

```
                    ╭──────────────╮
                    │   Dataset    │
                    ╰──────────────╯
                            │
                            ▼
              ┌──────────────────────────┐
              │     Data Preprocessing    │
              └──────────────────────────┘
                    │                │
                    ▼                ▼
            ┌──────────────┐  ┌──────────────┐
            │ Training set │  │   Test set   │
            └──────────────┘  └──────────────┘
                    │                │
                    ▼                ▼
              ┌──────────────────────────┐
              │   Machine Learning Model  │
              └──────────────────────────┘
                            │
                            ▼
                    ╭──────────────╮
                    │   Predicted  │
                    │   outcome    │
                    ╰──────────────╯
```

# 3. IMPLEMENTATION

We experimented with different classification and ensemble algorithms to predict diabetes. In the following, we briefly discuss the phase.

**3.1. Dataset Description**- the data is gathered from UCI repository which is named as Pima Indian Diabetes Dataset. The dataset have many attributes of 768 patients.

```
diabetes_df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 768 entries, 0 to 767
Data columns (total 9 columns):
 #   Column                    Non-Null Count  Dtype
---  ------                    --------------  -----
 0   Pregnancies               768 non-null    int64
 1   Glucose                   768 non-null    int64
 2   BloodPressure             768 non-null    int64
 3   SkinThickness             768 non-null    int64
 4   Insulin                   768 non-null    int64
 5   BMI                       768 non-null    float64
 6   DiabetesPedigreeFunction  768 non-null    float64
 7   Age                       768 non-null    int64
 8   Outcome                   768 non-null    int64
dtypes: float64(2), int64(7)
memory usage: 54.1 KB
```

The 9th attribute is class variable of each data points. This class variable shows the outcome 0 and 1 for diabetics which indicates positive or negative for diabetics.

Distribution of Diabetic patient- We made a model to predict diabetes however the dataset was slightly imbalanced having around 500 classes labelled as 0 means negative means no diabetes and 268 labelled as 1 means positive means diabetic.
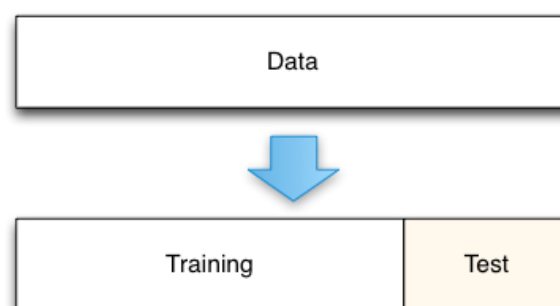
**3.2 Data Preprocessing**- Data preprocessing is most important process. Mostly healthcare related data contains missing vale and other impurities that can cause effectiveness of data. To improve quality and effectiveness obtained after mining process, Data preprocessing is done. To use Machine Learning Techniques on the dataset effectively this process is essential for accurate result and successful

prediction. For Pima Indian diabetes dataset we need to perform pre processing in two steps.

0. Missing Values removal- Remove all the instances that have zero (0) as worth. Having zero as worth is not possible. Therefore this instance is eliminated. Through eliminating irrelevant features/instances we make feature subset and this process is called features subset selection, which reduces dimensionality of data and help to work faster.

1. Splitting of data- After cleaning the data, data is normalized in training and testing the model. When data is spitted then we train algorithm on the training data set and keep test data set aside. This training process will produce the training model based on logic and algorithms and values of the feature in training data. Basically aim of normalization is to bring all the attributes under same scale.

**3.3.SPLITTING OF DATA:**The Processed data is splitted into the training set and testing set based on the size of the dataset.



**3.4.Apply Machine Learning**- When data has been ready we apply Machine Learning Technique. We use different classification and ensemble techniques, to predict diabetes. The methods applied on Pima Indians diabetes dataset. Main objective to apply Machine Learning Techniques to analyze the performance of these methods and find accuracy of them, and also been able to figure out the

responsible/important feature which play a major role in prediction. The Techniques are follows-

**3.4.1.Support Vector Machine**- Support Vector Machine also known as svm is a supervised machine learning algorithm. Svm is most popular classification technique. Svm creates a hyperplane that separate two classes. It can create a hyperplane or set of hyperplane in high dimensional space. This hyper plane can be used for classification or regression also. Svm differentiates instances in specific classes and can also classify the entities which are not supported by data. Separation is done by through hyperplane performs the separation to the closest training point of any class.

Algorithm-

- Select the hyper plane which divides the class better.

- To find the better hyper plane you have to calculate the distance between the planes and the data which is called Margin.

- If the distance between the classes is low then the chance of miss conception is high and vice versa. So we need to

- Select the class which has the high margin. Margin = distance to positive point + Distance to negative point.

**3.4.2.Decision Tree**- Decision tree is a basic classification method. It is supervised learning method. Decision tree used when response variable is categorical. Decision tree has tree like structure based model which describes classification process based on input feature. Input variables are any types like graph, text, discrete, continuous etc. Steps for Decision Tree Algorithm-

- Construct tree with nodes as input feature.

- Select feature to predict the output from input feature whose information gain is highest.

- The highest information gain is calculated for each attribute in each node of tree.

- Repeat step 2 to form a subtree using the feature which is not used in above node.

**3.4.3.Random Forest-** It is type of ensemble learning method and also used for classification and regression tasks. The accuracy it gives is grater then compared to other models. This method can easily handle large datasets. Random Forest is developed by Leo Bremen. It is popular ensemble Learning Method. Random Forest Improve Performance of Decision Tree by reducing variance. It operates by constructing a multitude of decision trees at training time and outputs the class that is the mode of the classes or classification or mean prediction (regression) of the individual trees.

- The first step is to select the R features from the total features m where R<<M.

- Among the R features, the node using the best split point.

- Split the node into sub nodes using the best split.

- Repeat a to c steps until l number of nodes has been reached.

- Built forest by repeating steps a to d for a number of times to create n number of trees.

  The random forest finds the best split using the Gin-Index Cost Function which is given by:

  The first step is to need the take a glance at choices and use the foundations of each indiscriminately created decision tree to predict the result and stores the anticipated outcome at intervals the target place. Secondly, calculate the votes for each predicted target and ultimately, admit the high voted predicted target as a result of the ultimate prediction from the random forest formula. Some of the options of Random Forest does correct predictions result for a spread of applications are offered.

# 4.RESULTS

By comparing all models on each other, RANDOM FOREST MODEL has more accuracy in predicting diabetes.

| Classifier | Accuracy |
|---|---|
| Random forest | 0.7597 |
| Decision Tree | 0.6948 |
| Support Vector Machine | 0.7727 |

Table 1: Evaluation of Models

# 5.CONCLUSION AND FUTURE SCOPE

The main aim of this project was to design and implement Diabetes Prediction Using Machine Learning Methods and Performance Analysis of that methods and it has been achieved successfully. The proposed approach uses various classification and ensemble learning method in which SVM, Knn, Random Forest, Decision Tree, Logistic Regression and Gradient Boosting classifiers are used. And 77% classification accuracy has been achieved. The Experimental results can be asst health care to take early prediction and make early decision to cure diabetes and save humans life. In future, the designed system with the used machine learning classification algorithms can be used to predict or diagnose other diseases. The work can be extended and improved for the automation of diabetes analysis including some other machine learning algorithms.