# An Expertise System for Insulin Dosage Prediction using Machine Learning Techniques

## Dr.K.Vijay Kumar, K.Mounika Sai Sadhvi, M.Bhargavi, K.Sri Varshini, K.Sirisha

Department of Computer Science and Engineering, Vignan"s Institute of Engineering for Women,

Visakhapatnam, A.P, India

**Abstract:** Diabetes Mellitus is a chronic metabolic disorder. Normally, with a proper adjusting of blood glucose levels (BGLs), diabetic patients could live a normal life without the risk of having serious complications that normally developed in the long run. However, blood glucose levels of most diabetic patients are not well controlled for many reasons. Although the traditional prevention techniques such as eating healthy food and conducting physical exercise are important for the diabetic patients to control their BGLs, however taking the proper amount of insulin dosage has the crucial rule in the treatment process. In this project we are using Gradient Boosting Classifier to predict diabetes and then using Linear Regression algorithm to predict insulin dosage in diabetis detected patients. To implement this project we are using PIMA diabetes dataset and UCI insulin dosage dataset. We are training both algorithms with above mention dataset and once after training we will upload test dataset with no class label and then Gradient Boosting will predict presence of diabetes and Linear Regression will predict insulin dosage if diabetes detected by Gradient Boosting.

**Keywords:** Diabetis Mellitus, Blood Glucose Levels (BGLs), Insulin Dosage, Gradient Boosting Classifier, Linear Regression.

## 1. INTRODUCTION:

Diabetes prevalence has been rising more rapidly in middle and low-income countries. The number of people with diabetes has risen from 108 million in 1980 to 422 million in 2014. The global prevalence of diabetes among adults over 18 years of age has risen from 4.7% in 1980 to 8.5% in 2014. A lack of insulin, or an inability to adequately respond to insulin, can each lead to the development of the symptoms of diabetes. For a diabetes patient, insulin dose is necessary to control the level of glucose. The doctor of the patient also has to know the required insulin dose from previous records of doses & from patient's current calculated blood sugar level. This has inspired us to make a research on how to predict the insulin dose level of a patient before every meal.

Lastly, we want to tell that there are many researches going on to discover newer methods to predict insulin dose. Despite having all those methods, we believe our research data will help doctors to predict almost accurate insulin dose of diabetes patients.

## 2. LITERATURE SURVEY

**Predicting Diabetes Mellitus using Data Mining Technique:** Diabetes is a chronic disease caused due to the expanded level of sugar addiction in the blood. Various automated information systems were outlined utilizing various classifiers for anticipate and diagnose the diabetes. Data mining approach helps to diagnose patient's diseases. Diabetes Mellitus is a chronic disease to affect various organs of the human body. Early prediction can save human life and can take control over the diseases. Selecting legitimate classifiers clearly expands the correctness and adeptness of the system. Due to its continuously increasing rate, more and more families are unfair by diabetes mellitus. Most diabetics know little about their risk factor they face prior to diagnosis. This paper explores the early prediction of diabetes using data mining techniques. The dataset has taken 768 instances from PIMA Indian Diabetes Dataset to determine the accuracy of the data mining techniques in prediction. Then we developed five predictive models using 9 input variables and one output variable from the Dataset information; we evaluated the five models in terms of their accuracy, precision, sensitivity, specificity and F1 Score measures. The purpose of this study is to compare the performance analysis of Naïve Bayes, Linear Regression, Artificial neural networks (ANNs), C5.0 Decision Tree and Support Vector Machine (SVM) models for predicting diabetes using common risk factors. The decision tree model (C5.0) had given the best classification accuracy, followed by the linear regression model, Naïve Bayes, ANN and the SVM gave the lowest accuracy Index Terms—Data mining, Prediction, Naïve Bayes, Logistic Regression, C5.0 Decision Tree, Artificial Neural Networks (ANN) and Support Vector Machine (SVM).

**Analysis of Various Data Mining Techniques to Predict Diabetes Mellitus:** Data mining approach helps to diagnose patient's diseases. Diabetes Mellitus is a chronic disease to affect various organs of the human body. Early prediction can save human life and can take control over the diseases. This paper explores the early prediction of diabetes using various data mining techniques. The dataset has taken 768 instances from PIMA Indian Dataset to determine the accuracy of the data mining techniques in prediction. The analysis proves that Modified J48 Classifier provide the highest accuracy than other techniques.

**Comparison Data Mining Techniques to Prediction Diabetes Mellitus:** Diabetes is one of the chronic diseases caused by excess sugar in the blood. Various methods of automated algorithms in various to anticipate and diagnose diabetes. One approach to data mining method can help diagnose the patient's disease. In the presence of predictions can save human life and begin prevention before the disease attacks the patient. Choosing a legitimate classification clearly expands the truth and accuracy of the system as levels continue to increase. Most diabetics know little about the risk factors they face before the diagnosis. This method uses developing five predictive models using 9 input variables and one output variable from the dataset information. The purpose of this study was to compare performance analysis of Naive Bayes, Decision Tree, SVM, K-NN and ANN models to predict diabetes mellitus.

Diabetes is a disease caused by hyperglycaemia (high blood glucose level). Diabetes affects an estimated 3-4% of the world's population (half of whom are undiagnosed), making it one of the major chronic illnesses prevailing today. It is caused by hyperglycaemia resulting from defects in insulin secretion, insulin action, or both. The chronic hyperglycaemia of diabetes is associated with long-term damage, dysfunction, and failure of various organs, especially the eyes, kidneys, nerves, heart, and blood vessels. This deficiency leads to destruction of the b-cells of the pancreas with consequent insulin deficiency to abnormalities that result in resistance to insulin action and reaction process. Deficiency of insulin results from inadequate insulin secretion. This Improper insulin secretion and defects in insulin action is the primary cause of hyperglycaemia. So, the significance of insulin dose I s clearly visible.

The prediction of glucose concentrations could facilitate the appropriate patient reaction in crucial situations such as hypoglycaemia. Thus, several recent studies have considered advanced data driven techniques for developing accurate predictive models of glucose metabolism. The fact that the relationship between input variables (i.e., medication, diet, physical activity, stress etc.) and glucose levels is nonlinear, dynamic, interactive and patientspecific necessitates the application of non-linear regression models such as artificial neural networks, support vector regression and Gaussian processes.

The aim of project to predict diabetes to predict. In addition to the general guidelines that the patient follows during his daily life, several diabetes management systems have been proposed to further assist the patient in the self-management of the disease. One of the essential components of a diabetes management system concerns the predictive modelling of the glucose metabolism.

## 2. PROPOSED SYSTEM

### 3.1. Methodology

In this project we are using Gradient Boosting Classifier to predict diabetes and then using Linear Regression algorithm to predict insulin dosage in diabetic detected patients. To implement this project we are using PIMA diabetes dataset and UCI insulin dosage dataset.The objective of the PIMA dataset is to diagnostically predict whether or not a patient has diabetes, based on certain diagnostic measurements included in the dataset.The UCI Machine Learning Repository is a collection of databases, domain theories, and data generators that are used by the machine learning community for the empirical analysis of machine learning algorithms. We are training both algorithms with above mention dataset and once after training we will upload test dataset with no class label and then Gradient Boosting will predict presence of diabetes and Linear Regression will predict insulin dosage if diabetes detected by Gradient Boosting.

### 3.2. Dataset

The PIMA Indian Diabetes Dataset, originally from the National Institute of Diabetes and Digestive and Kidney Diseases, contains information of 768 women from a population near Phoenix, Arizona, USA. The outcome tested was Diabetes, 258 tested positive and 500 tested negative. Therefore, there is one target (dependent) variable and the 8 attributes (TYNECKI, 2018): pregnancies, OGTT(Oral Glucose Tolerance Test), blood pressure, skin thickness, insulin, BMI(Body Mass Index), age, pedigree diabetes function. The Pima population has been under study by the National Institute of Diabetes and Digestive and Kidney Diseases at intervals of 2 years since 1965. As epidemiological evidence indicates that T2DM results from interaction of genetic and environmental factors, the Pima Indians Diabetes Dataset includes information about attributes that could and should be related to the onset of diabetes and its future complications.

**Figure 1: Diabetes Dataset**



**Figure 2: Insulin Dataset**



**Figure 3: Testvalues Dataset**
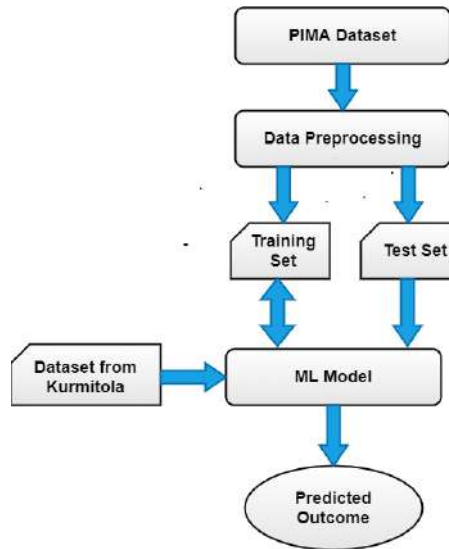
### 3.3. Architecture



Figure 4: Architecture Diagram

## 3.    ALGORITHMS AND FLOW CHART

4.

### 4.1. Algorithm1 (Gradient Boosting)

Gradient boosting is a machine learning technique used in regression and classification tasks, among others. It gives a prediction model in the form of an ensembleof weak prediction models, which are typically decision trees. When a decision tree is the weak learner, the resulting algorithm is called gradient-boosted trees. XGBoost is an optimized distributed gradient boosting library designed to be highly efficient, flexible and portable. It implements machine learning algorithms under the Gradient Boosting framework.

Gradient boosting involves three elements:

**Loss function** :A loss function to be optimized. The loss function used depends on the type of problem being solved.

**Weak learner** :A weak learner to make predictions. Decision trees are used as the weak learner in gradient boosting.

**Addictive model** :An additive model to add weak learners to minimize the loss function. Trees are added one at a time, and existing trees in the model are not changed.

A gradient descent procedure is used to minimize the loss when adding trees.

**Step 1:** Creating classification dataset with make_classification

**Step 2**: Building Gradient Boosting Classifier

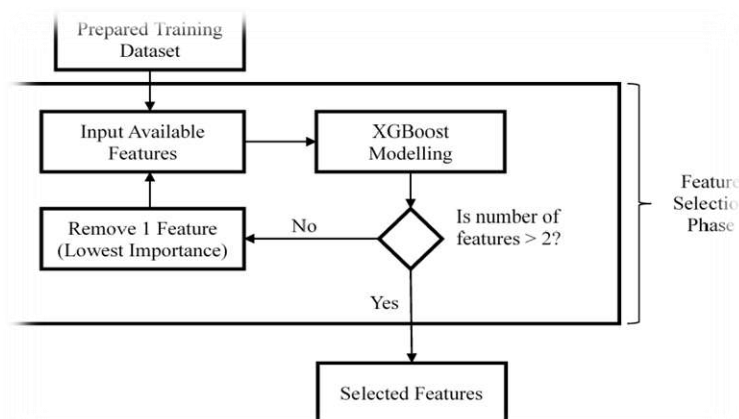 **Step 3:** Performing prediction with a classification model



**Figure 5: flowchart of GradientBoosting**

### 4.2. Algorithm2 (Linear Regression)

Linear regression analysis is used to predict the value of a variable based on the value of another variable. The variable you want to predict is called the dependent variable. The variable you are using to predict the other variable's value is called the independent variable. It is a regression model that uses a straight line to describe the relationship between variables. It finds the line of best fit through your data by searching for the value of the regression coefficient(s) that minimizes the total error of the model.

**Step 1**: Import the packages and classes that you need.

**Step 2**: Provide data to work with, and eventually do appropriate transformations.

**Step 3**: Create a regression model and fit it with existing data.

**Step 4:** Check the results of model fitting to know whether the model is satisfactory.

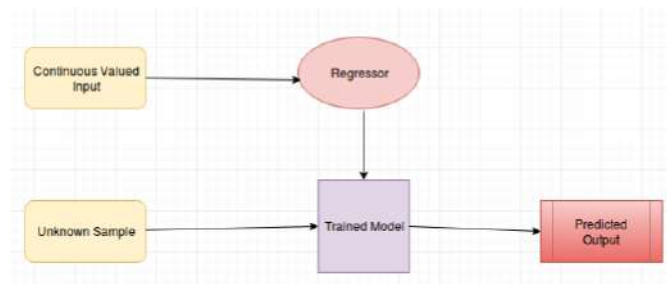**Step 5**: Apply the model for predictions



**Figure 6: Linear Regression**

## 5. RESULTS AND DISCUSSION

The outcomes of the proposed system is the level of severity identified from the diabetic patients after performing various processing techniques using machine learning algorithms.
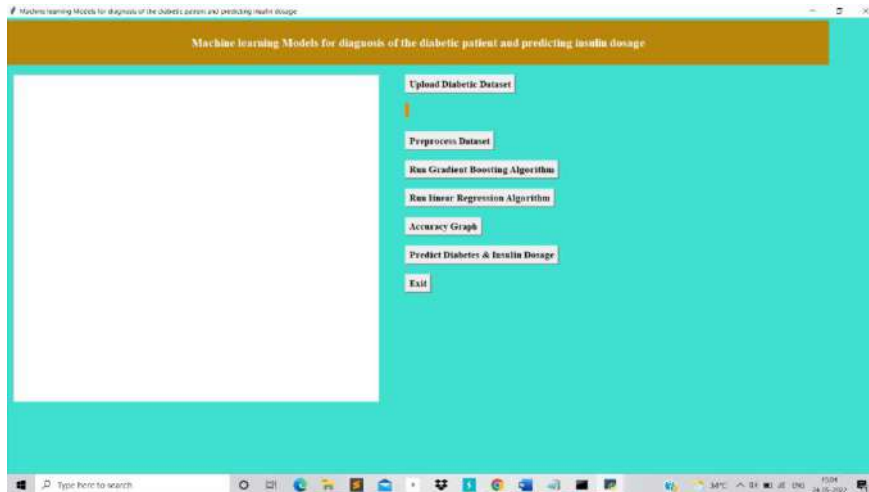
### 5.1. Processing Screens
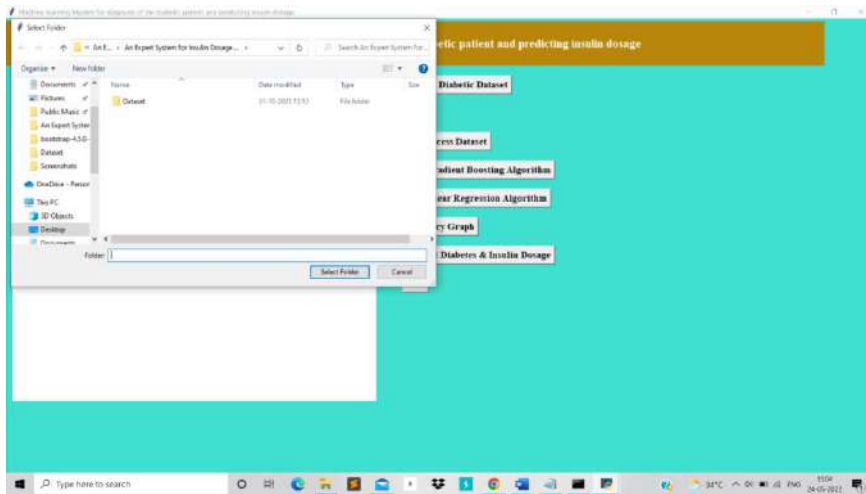


**Figure 7:Upload Diabetic Dataset**

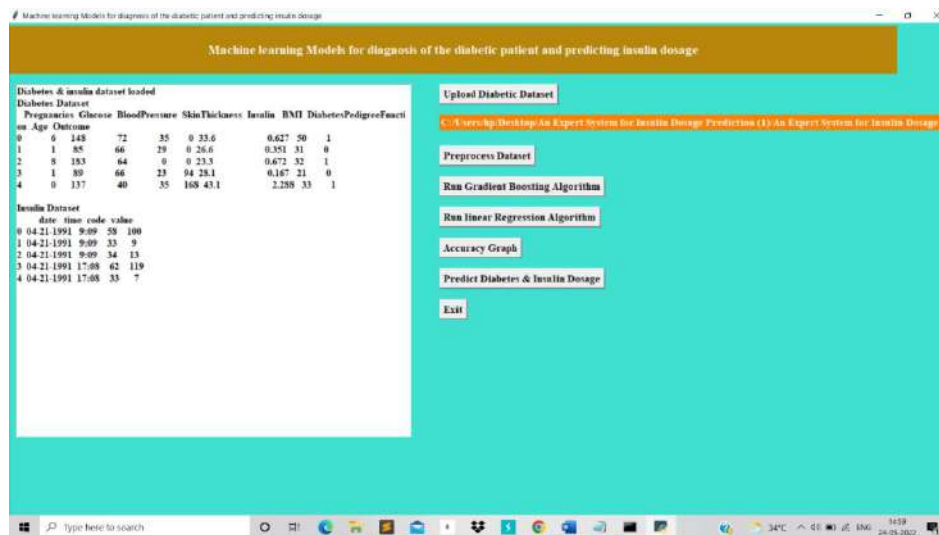**Figure 8: uploading entire Dataset folder to load both diabetes and insulin datasets**



**Figure 9: Both diabetis and insulin dataset loaded**



**Figure 10: red dots indicate presence of diabetes and blue represents no diabetes detected**

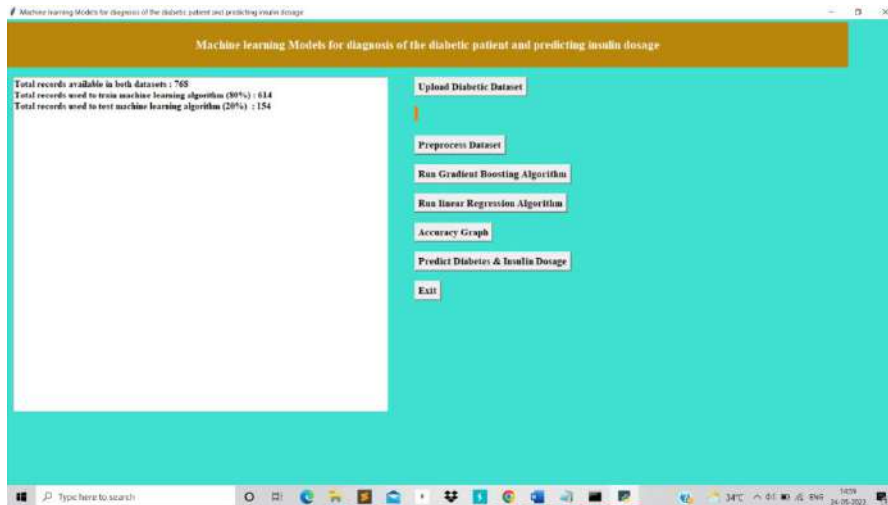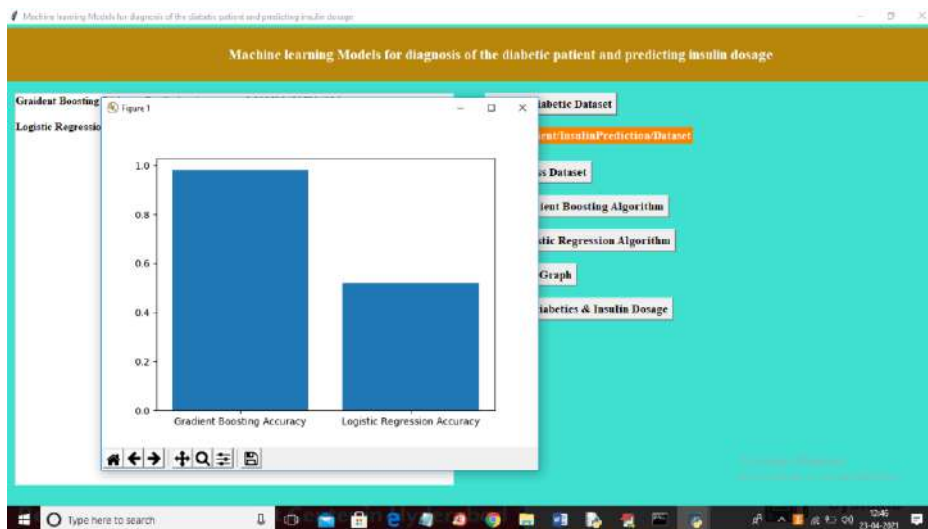**Figure 11: 'PreprocessDataset' to remove missing values and to split dataset into train and test**
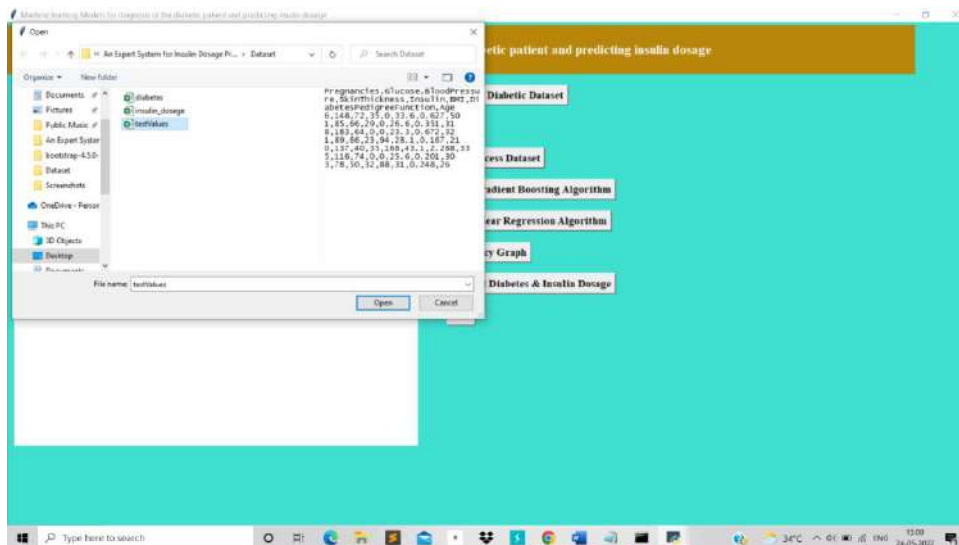


**Figure 12: representing accuracy of algorithms**



**Figure 13: Uploading 'testValues.csv' file**
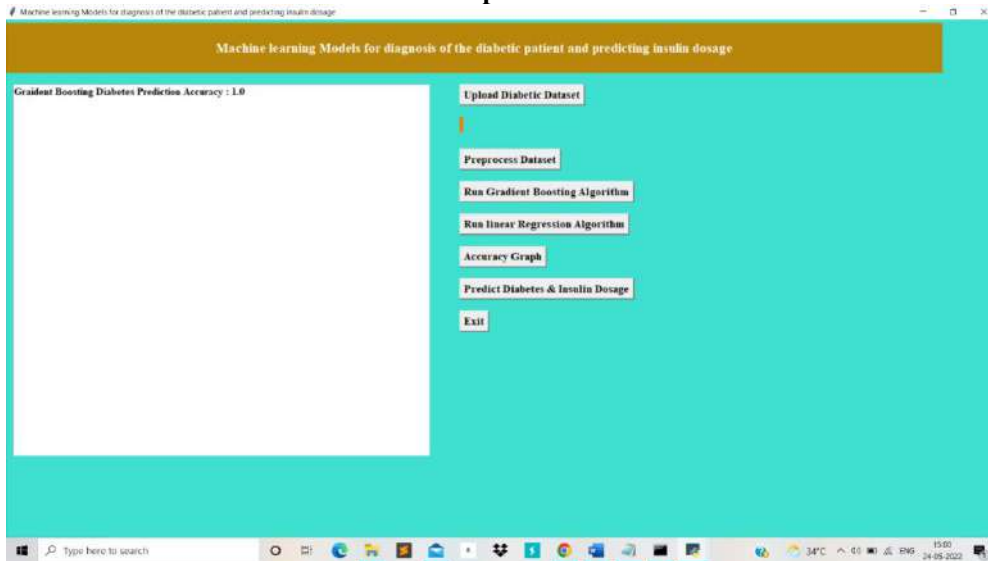
**5.2. Output Screens**



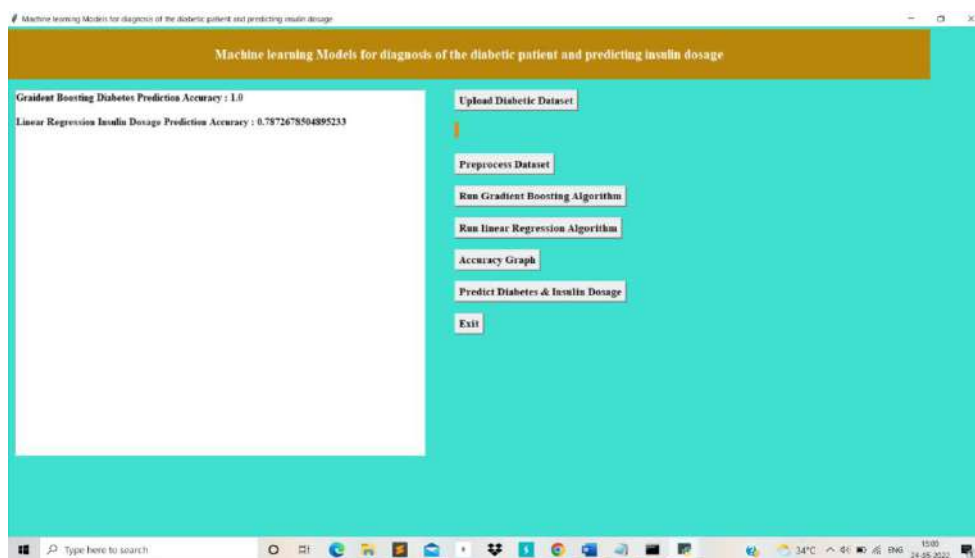**Figure 14: Diabetis is predicted with 100% accuracy**



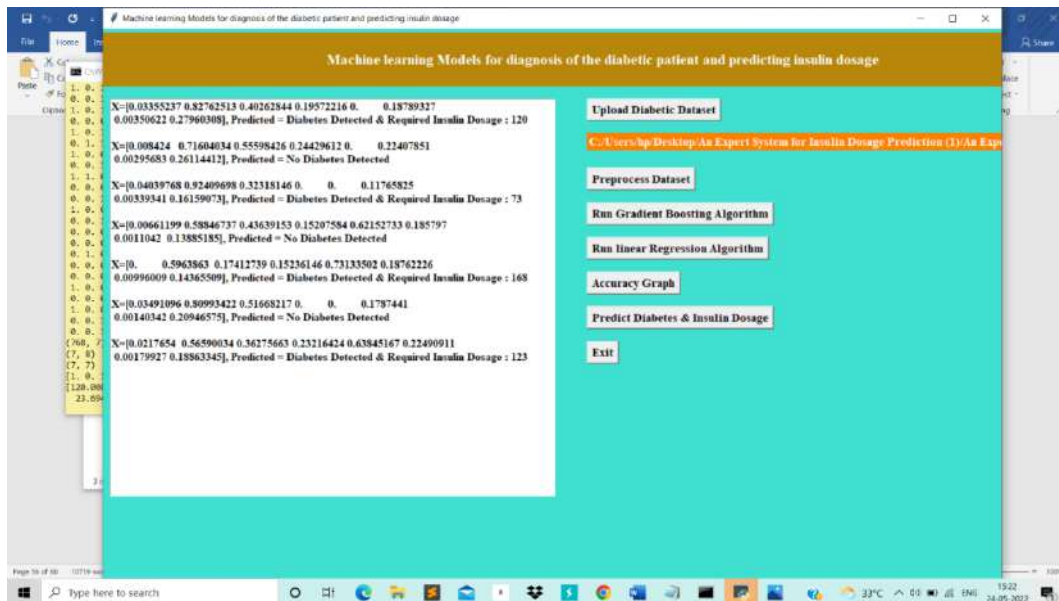**Figure 15: Insulin dosage is predicted with 78% accuracy**

**Figure 16: predicted result as 'No Diabetes Detected' or 'Diabetes Detected' and if diabetes detected then insulin dosage predicted**

## 6. CONCLUSION

This paper was aimed at modelling neural network for the prediction of amount of insulin dosage suitable for diabetic patients. A model based on Gradient boosting trained with BP was used. The model uses four input information about each patient its length, weight blood sugar, and gender. Many experiments were conducted on 180 patient's data. we are using Gradient boosting to predict diabetes and then applying Linear Regression Algorithm to predict insulin dosage if diabetes detected by Gradient boosting algorithm. The Gradient boosting model converged fast and gave results with high performance.

## REFERENCES

[1] American Diabetes Association. Diagnosis and Classification of Diabetes Mellitus. Diabetes Care, Vol. 31, No. 1, 2008, 55-60, 1935- 5548.

[2] I. EleniGeorga, C. Vasilios Protopappas and I. DimitriosFotiadi,. "Glucose Prediction in Type 1 and Type 2 Diabetic Patients Using Data Driven Techniques," Knowledge-Oriented Applications in Data Mining, InTech pp 277-296, 2011.

[3] V. Tresp,, T. Briegel, and J. Moody, "Neural-network models for the blood glucose metabolism of a diabetic," IEEE Transactions on Neural Networks, Vol. 10, No. 5, 1999, 1204-1213, 1045-9227.

[4] C. M. Bishop, . Pattern Recognition and Machine Learning, Springer, 2006, New York.

  [5] S. Haykin, Neural networks and learning machines. Pearson 2008.

[6] W.D. Patterson, Artificial Neural Networks- Theory and Applications, Prentice Hall , Singapore. 1996.

[7] M. Pradhan and R. Sahu, "Predict the onset of diabetes disease using Artificial Neural Network (ANN)," International Journal of Computer Science & Emerging Technologies, 303 Volume 2, Issue 2, April 2011.

[8] W, Sandham, D,Nikoletou, D.Hamilton, K, Paterson, A. Japp and C. MacGregor, "BLOOD Glucose Prediction for Diabetes THERAPY USING A RECURRENT Artificial Neural Networks ", EUSIPCO, Rhodes, 1998, PP. 673-676

[9] M, Divya, R. Chhabra, S. Kaur and S. , "Diabetes Detection Using Artificial Neural Networks & Back-Propagation Algorithm", International Journal of Scientific & Technology Research". V 2, ISSUE 1, JANUARY 2013.

[10] G, Robertson, E. Lehmann, W. Sandham, and D. Hamilton Blood, "Glucose Prediction Using Artificial Neural Networks Trained with the AIDA Diabetes Simulator: A Proof-of-Concept Pilot Study". Journal of Electrical and Computer Engineering , V 2011 (2011), Article ID 681786, 11 pages.