

International Advanced Research Journal in Science, Engineering and Technology ISO 3297:2007 Certified ∺ Impact Factor 7.105 ∺ Vol. 9, Issue 6, June 2022

DOI: 10.17148/IARJSET.2022.96105

# **Diabetes Prediction Using Machine Learning**

# Manjunath Ganganagoudar<sup>1</sup>, H K Madhu<sup>2</sup>

Student, Dept. of MCA, Bangalore Institute of Technology, Bengaluru, India<sup>1</sup>

Associate Professor, Dept. of MCA, Bangalore Institute of Technology, Bengaluru, India<sup>2</sup>

**Abstract:** Diabetes is an illness caused by high altitude of the glucose in the humans body. Diabetes it should not be disregarded and left untreated. Diabetes can cause serious problems such as Blood Heart, kidney problems, circulatory strain, eye harm, and it will also affect different organs in the humans body. Early prediction helps managesdiabetes. To accomplish this objective, Projectwill apply an assortment of Machine Learning methods to predict and improve the accuracy of earlydiabetes in the human body and patients. Machine learning strategies give improved results to forecast by building model from datasets gathered from patients. In the task, applymachine learning order and group strategies to outdataset to predict diabetes. These are the Decision Tree tool algorithm, Support Vector Machine algorithm, XgBoost Classifier algorithm, and Random Forest algorithm. The accuracy of each model is different from the other models. The work of Project Providesan accurate or more accurate model andshows that the model can effectively predict diabetes. Our results show that Random Forest algorithm achieves great accuracy compared with other machine learning.

## 1. INTRODUCTION

Diabetes is a typical constant illness that will be a serious danger for human wellbeing. Diabetes is a harmful disease in the world. Diabetes due to obesity and hyperglycaemia. It affects the hormone insulin causes crabs to have abnormal metabolisms .Diabetes will accouris when our body doesn't convey sufficient insulin. As per the World Health Organization, around 422 million individuals have diabetes.

However, diabetes is widespread in various countries such as Canada, China and India. The number of the diabetics patients in India is 40 million, as in India population at present nearly of 100 million. Diabetes is the main source of death around the world.

Diabetes will be detected when the blood sugar levels are higher than the normal. This is brought about by high insulin discharge or natural impacts. Diabetes can make an assortment of mischief our body and can harm tissues, kidneys, eyes and veins. Diabetes will be divided into 2 categories: type 1 diabetes categories and type 2 diabetes categories.

Patients with the type 1 are usually under an age of 30. The clinical side effects are increased thirst and continuous pee. This kind of diabetes requires treatment and can not be eliminated with the medication.

Type 2 diabetes is more normal in moderately aged and more seasoned individuals and may indicate high blood pressure, obesity, and other illnesses. Due to our standard of living, diabetes is increasing in people's daily lives. Therefore, it is worth studying how to analyse diabetes. You will be diagnosed early so that you can control the diabetes. Machine learning can settle on primer choices about diabetes by consulting with a doctor based on physical examination data. As of late, numerous algorithms have been used to anticipate diabetes, Machine learning techniques some of these Random Forest, Decision Trees, SVMs. These machine learning methods can be utilized to predict diabetes by building predictive models derived from medical datasets. By extracting this knowledge, we can predict diabetics. Use the best prediction techniques based on the attributes of a given dataset to get the perfect accuracy for predicting diabetes.

## 2. LITERATURESURVEY

K. Vijiya Kumaretal. [1] The Random Forest Algorithm proposed for diabetes expectation utilizes AI procedures to foster a framework that can all the more precisely perform early forecast of diabetes in patients. An proposed model gave an best outcomes for the diabetes forecast, and outcomes demonstrated an way that the expectation framework could foresee diabetes sickness actually, proficiently, and above all, right away.

Nonso Nuna Moko et al. [2] recommended anticipating the improvement of diabetes: the gathering way to deal with administered learning they used. Five large used classifiers are used in troupe, and a meta-classifier is used to add up to the outcomes. Results are presented and differentiated with relative assessments utilizing the identical dataset in the composition. It was demonstrated the way that the beginning of diabetes can be anticipated all the more precisely utilizing the proposed strategy.



International Advanced Research Journal in Science, Engineering and Technology

ISO 3297:2007 Certified 💥 Impact Factor 7.105 💥 Vol. 9, Issue 6, June 2022

#### DOI: 10.17148/IARJSET.2022.96105

N. Joshi et al. [3] Diabetes forecast utilizing the introduced AI procedures intends to foresee diabetes utilizing three different directed AI techniques, including SVM, Decision Tree, Random Forest, and XgBoost. This task proposes powerful procedures for early discovery of diabetes.

Muhammad Azeem Sarwaretal. [4] A proposed investigation of foreseeing diabetes utilizing AI calculations in medical care applied six different AI calculations to make sense of and look at the exhibition and exactness of the applied calculations. increment. Looking at the changed AI strategies utilized in this study uncovers the best calculation for anticipating diabetes.

Yasodhaet al.[5] utilizes the characterization on different kinds of datasets that can be achieved to choose if an individual is diabetic or not. The diabetic patient's informational index is laid out by social event information from medical clinic distribution center which contains 200 examples with nine credits. These occasions of this dataset are alluding to two gatherings for example blood tests and pee tests. In this review the execution should be possible by utilizing WEKA to order the information and the information is evaluated through 10-overlay cross approval approach, as it performs very well on little datasets, and the results are looked at. The innocent Bayes, J48, REP Tree and Random Tree are utilized. It was reasoned that J48 works best appearance an exactness of 60.2% among others.

Aiswaryaet al. [6] expects to find answers for distinguish the diabetes by the exploring and analyzing the examples start in the information through order examination by utilizing Decision Tree and Naïve Bayes calculations. The examination desires to propose a quicker and more effective technique for recognizing the illness that will help in very much coordinated fix of the patients.

Utilizing PIMA dataset and cross approval approach the review inferred that J48 calculation gives a precision pace of 74.8% while the innocent Bayes gives an exactness of 79.5% by utilizing 70:30 split.

Gupta et al. [7] means to find and compute the exactness, awareness and explicitness level of various arrangement strategies and furthermore attempted to think about and dissect the aftereffects of a few order techniques in WEKA, the review looks at the presentation of same classifiers when executed on certain devices which incorporates RapidMiner and similar boundaries (for example precision, awareness and explicitness). They applied JRIP, Graft and Bayes Net calculations. The outcome shows that Graft shows most elevated precision i.e 81.3%, responsiveness is 59.7% and particularity is 81.4%. It was additionally reasoned that WEKA works best than MATLAB and RapidMiner.

Deeraj Shetty et al. [8] proposed diabetes infection expectation utilizing information mining collect Intelligent Diabetes Disease Prediction System that gives examination of diabetes ailment using diabetes patients data set. In this framework, they propose the utilization of calculations like Bayesian and KNN (K-Nearest Neighbor) to apply on diabetes patients data set and break down them by taking different properties of diabetes for expectation of diabetes disease

.Jian-xunChen, Shih-LiSu and Che-Ha Chang[9]discussed helping in exploring the patient's records of their forecast and their condition. This paper shows the likelihood to give customize diabetes mellitus care arranging productively.

Lakshmi K.S and G.Santhosh Kumar [10]as indicated by them Hospital data sets act as well off data hotspot for the productive drug determination. IN this they utilized NLP devices alongside joined with information digging calculations for the extraction of rules.

#### **3. METHODOLOGY**

The purpose of this work is to investigate models for more accurate prediction of diabetes. Several characterization and outfit calculations were tested to anticipate diabetes. The phases are briefly described below.

**1. Dataset Description** – An data is taken incarnation a UCI vault called the Pima India Diabetes Dataset. In this dataset it had the information about the 780 patients information.

#### Table 1: Data record portrayal

The data table which had been contain the information Patients in the 10 column. The class variable of every information point is the 10th property. This class variable shows the outcome for diabetics (0 or 1), demonstrating whether it is positive or negative for diabetes.

**Diabetes Distribution**(**DD**) – We have created a model for predict diabetes, but about 500 classes are labelled 0, negative means no diabetes, 268 is labelled 1 and that they are diabetic. To affirm, the dataset was slightly imbalanced. **2. Data pre-processing** 



#### International Advanced Research Journal in Science, Engineering and Technology

ISO 3297:2007 Certified 💥 Impact Factor 7.105 💥 Vol. 9, Issue 6, June 2022

#### DOI: 10.17148/IARJSET.2022.96105

Data pre-processing is the main interaction. Most wellbeing related information contains missing qualities and different foreign substances that can influence the legitimacy of the information. Information pre-processing is performed to work on the quality and viability accomplished after the mining system. For effective application of machine learning methods to datasets, this cycle is fundamental for precise results and effective predictions. For the Indian Pima Diabetes Dataset, pre-treatment is required in two ways.

**1. Delete Missing Value**-Delete all examples with zero (0) as the worth. It cannot have zero as a value. Therefore, this instance will be deleted. Create a feature subset by eliminating irrelevant features / instances. You can work more quickly and reduce the dimensions of the data by using this procedure, known as feature subset selection.

**2. Data partitioning**- Subsequent to tidying up the information, the information is standardized during preparing and testing of the model. At the point when the information is let out, utilize the preparation dataset to prepare the calculation and put the test dataset away.Based on the functional logic, algorithms, and values of the training data, this training process generates a training model. Putting all of the attributes on the same scale is the fundamental aim of normalisation.

#### **3.Apply Machine Learning**

When the information was prepared, apply machine learning methods. Utilize distinctive classification and gathering methods foresee diabetes. Strategies connected to Pima Indian diabetes datasets. The most objective is for apply machine learning procedures to dissect the execution of the strategies, discover exactness, be able to find mindful / imperative highlights that play a key part in expectation. ... The procedure is as takes after: The technique is as follow:

**1.Support Vector Machine algorithm** – is a managed Machine learning calculation. SVM is most well-known characterization method. Svm makes a hyperplane that isolates the two classes. You can make hyperplanes or sets of hyperplane was high-layered spaces. This hyperplanes likewise utilized characterization, relapse. Svm can likewise recognize occasions into explicit classes and classify elements that are not upheld by data. Division is finished by having the hyperplane perform the partition to the nearest preparing point of any class. Algorithm-

- Select a hyperplane that better partitions the class.
- To find a preferable hyperplane over, we want to ascertain distance of the plane and the data, called the edge.
- If opening between class is little, there is a high likelihood of premature delivery as well as the other way around.

• Select class with the most noteworthy wiggle room. Edge = distance to positive direct + distance toward negative point.

**2.Decision Tree**- The choice tree is the essential characterization strategy. It is a managed learning technique. Choice tree utilized when the reaction variable is of absolute kind. Choice trees have a construction based model, for example, a tree that depicts an order cycle in light of info highlights. Input factors can be of any kind, like chart, text, discrete, ceaseless, etc. Choice Tree

#### Algorithm Steps-

- Fabricate a tree involving hubs as information highlights.
- Select the capability to foresee the result of the information capability with the most elevated data gain.
- The most noteworthy data gain is determined for each tree hub characteristic.
- Rehash stage 2 to shape a subtree with capabilities not utilized in the above hub

**3:Random Forest**: This is sort of troupe learning strategy that is likewise used for characterization, relapse errands. The precision is superior to different models. This technique makes it simple to work with immense datasets. Irregular Forest was created by an Leo Bremen. This a famous troupe learning technique. Work on the exhibition of choice trees by lessening irregular timberland conveyance. It works by building different choice trees during preparing and yielding a class that is the strategy of the class or the characterization o mean expectation of each tree. Algorithm-

- The underlying step is to pick the R highlights from the complete elements m where R<<M.
- Among the R incorporates, hub utilizing the best separated.
- Part an hub into subs hubs utilizing the best separated.
- Go over a to c strides until l number of centre points has reached.
- Constructed timberland reiterating stages a to d number of times to makes n numbers of trees.

# IARJSET



International Advanced Research Journal in Science, Engineering and Technology

ISO 3297:2007 Certified 💥 Impact Factor 7.105 💥 Vol. 9, Issue 6, June 2022

#### DOI: 10.17148/IARJSET.2022.96105

An arbitrary backwoods s finds best separated utilizing the Gin-Index Cost Function which is given by: An underlying step to require the take take a gander at decisions and utilize the groundworks every unpredictably gone with choice tree to anticipate the outcome and stores the normal outcome ranges the goal spot. Besides, compute the decision in Favor of each expected goal and finally, surrender the high casted a voting form expected focus due to a conclusive expectation form an irregular woods equation. A part of the decisions of Random Forest redresses expectations result for a spread of uses are promoted.

#### 4. MODEL BUILDING

This is the main stage, including building a model for foreseeing diabetes. In it, we executed the different AI calculations depicted above for diabetes expectation.

Minutes of the proposed procedure -

Stage 1: Imports an expected libraries and import the diabetes dataset.

Stage 2: Pre-process the information and eliminate the missing information.

Stage 3: Perform a 80% percent split to part the dataset into preparing datasets and 20% into test datasets.

Stage 4: Select an AI calculation. H.

Support Vector Machine method, Decision Tree method, XgBoost, Random Forest.

Stage 5: Build a classifier model for the AI calculation depicted above in view of the preparation set.

Stage 6: Test the classifier model of the AI calculation depicted above in view of the test set.

Stage 7: Perform a near assessment of the exploratory exhibition results got with every classifier.

Stage 8: After breaking down in view of different estimations, complete the best preforming calculation.

#### **5.EXPERIMENTAL RESULTS**

In this work various steps were taken. The proposed approach utilizes different grouping and gathering strategies and carried out utilizing python. These techniques are standard Machine Learning strategies used to get the best accuracy from information. In this work we see that Random Forest classifier accomplishes better contrasted with others. We have involved best Machine Learning methods for expectation and to accomplish execution exactness. Figure shows the aftereffect of these Machine Learning strategies.

#### Figure 3: Precision Result of Machine learning strategies

```
from sklearn import metrics
predictions = rfc.predict(X_test)
print("Accuracy_Score =", format(metrics.accuracy_score(y_test, predictions)))
Accuracy_Score = 0.7716535433070866
```

Characterization report and matrix of the Random Forest

<pre>from sklearn.metrics import classification_report, confusion_matrix</pre>						
<pre>print(confusion_matrix(y_test, predictions)) print(classification_report(y_test,predictions))</pre>						
[[136]	26] 6011					
[ ]2 (	00]]	precision	recall	f1-score	support	
	0	0.81	0.84	0.82	162	
	1	0.70	0.65	0.67	92	
accuracy				0.77	254	
macro	o avg	0.75	0.75	0.75	254	
weighte	d avg	0.77	0.77	0.77	254	





#### International Advanced Research Journal in Science, Engineering and Technology

ISO 3297:2007 Certified 💥 Impact Factor 7.105 💥 Vol. 9, Issue 6, June 2022

#### DOI: 10.17148/IARJSET.2022.96105

#### 6.CONCLUSION

The fundamental objectives of this task were the turn of events and execution of diabetes expectation utilizing AI methods and the exhibition examination of these procedures, which were effectively accomplished. The proposed approach utilizes different grouping and gathering learning techniques utilizing SVM,

Irregular Forest, choice trees, strategic, and XgBoost classifiers. Furthermore, 77% grouping precision was accomplished. Trial results assist medical services suppliers with settling on early expectations and early choices to treat diabetes and save lives.

#### REFERENCES

1.K.VijiyaKumar, B.Lavanya and I.Nirmala, S.Sofia Caroline, "Random Forest Algorithm for the Prediction of Diabetes".Proceeding of International Conference on Systems Computation Automation and Networking, 2019.

2. Nonso Nnamokao, Abir Hussain and David England, "Vorhersage des Beginns von Diabetes: ein Ensembleüberwachter Lernansatz". IEEE of a Congress of Evolutionary Computation (CEC), 2018.

3. Tejas N. Joshi, Prof. Pramila M. Chawan, "Diabetes Prediction Using Machine Learning Techniques". Int. Diary of Engineer-ing Research and Application, Vol. 8, Issue 1, (Part - II) January 2018, pp.-09-13

4. Nahla B., Andrew et al, "Savvy Support Vector Machine for Diagnosis of Diabetes. Data Technology in Biomedicine", IEEE Transaktionen .. 14, (July 2010), 1114-20

5. A.K., Dewangan, and P., Agrawal, Classification of Diabetes Mellitus Using Machine Learning Techniques, International Journal of Engineering and Applied Sciences, vol. 2, 2015.

6.Debadri Dutta, Debpriyo Paul, Parthajeet Ghosh, "Analyzing Feature Importances for Diabetes Prediction using Machine Learning". IEEE, pp 942-928, 2018.

7.Md. Faisal Faruque, Asaduzzaman, Iqbal H. Sarker, "Performance Analysis of Machine Learning Techniques to Predict Diabetes Mellitus". International Conference on Electrical, Computer and Communication Engineering (ECCE), 7-9 February, 2019.

8.Aljumah, A.A., Ahamad, M.G., Siddiqui, M.K., 2013. Application of data mining: Diabetes health care in young and old patients. Journal of King Saud University - Computer and Information Sciences 25, 127–136. doi:10.1016/j.jksuci.2012.10.003.

9 . Arora, R., Suman, 2012. Comparative Analysis of Classification Algorithms on Different Datasets using WEKA. International Journal of Computer Applications 54, 21–25. doi:10.5120/8626-2492.

10.Deeraj Shetty, Kishor Rit, Sohail Shaikh, Nikita Patil, "Diabetes Disease Prediction Using Data Mining "International Conference on Innovations in Information, Embedded and Communication Systems (ICHECS), 2017.