

# Diabetes Prediction in Women using Machine Learning Techniques

Spoorthy Y

Dept of Electronics and Communication  
GSSSIETW, Mysuru

Sunitha T

Dept of Electronics and Communication  
GSSSIETW, Mysuru

**Abstract**— Diabetes is an ailment caused in view of high glucose level in a human body. Diabetes sought not be overlooked in the event that it is untreated then Diabetes may cause some significant issues in an individual like: heart related issues, kidney issue, circulatory strain, eye harm and it can likewise influences different organs of human body. Diabetes can be controlled on the off chance that it is anticipated before. To accomplish this objective this venture work we will do early expectation of Diabetes in a human body or a patient for a higher precision through applying, Various Machine Learning Techniques. Machine learning procedures Provide better outcome for expectation by building models from datasets gathered from patients. In this work we will utilize Machine Learning Classification and troupe procedures on a dataset to anticipate diabetes. The precision accomplished by practical classifiers Random forest (RF), decision tree (DT), Support vector machine (SVM) and XGboost. Among them four, RF gives the best outcomes to diabetes beginning with a precision pace of 94.07% on the PIMA dataset. Henceforth, this proposed framework gives a powerful prognostic apparatus to medical services authorities. The outcomes got can be utilized to foster a novel programmed guess apparatus that can be useful in ahead of schedule location of the sickness. The result of the examination affirms that RF furnishes the best outcomes with the most encouraging removed highlights. RF accomplishes the precision of 92.34% which can be utilized for additional advancement of the programmed forecast device. The exactness of the

**Keywords**—Diabetes, Machine, Learning, Prediction, Dataset, Ensemble

## I. 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 it's tough to choose best technique. Thus for this purpose we apply popular classification and ensemble methods on dataset for prediction.

## II. LITERATURE REVIEW

Information mining procedures have overruled the current systems with better expectation, exactness, and accuracy. Additionally, Machine learning is an innovation of man-made brainpower that learns connections between hubs without monastery preparing them [2]. The significant capacity of AI methods to drive the forecast model without solid preparing identified with the hidden component. Information mining and AI strategies help to distinguish the information which stays covered up while utilizing the state of the art approach [3]. In this part, we will audit some past examinations to demonstrate the idea of information mining strategies convenience in the driving forecast model, fundamentally for diabetes.

Swapna G and others [4] made an examination that Machine learning practice has demonstrated valuable and proficient to build a forecast model for diabetes utilizing HRV signals in the DL approach. The creator was persuaded through the passings brought about by diabetes consistently on the planet which required staying away from the complexity of the sickness. The creator fostered another prescient model utilizing a convolutional neural organization (CNN), long transient memory (LSTM) and a gathering model for identifying compound ordered attributes of the information HRV information. Then, at that point SVM has been applied to those distinguished qualities for characterizing the information. The proposed framework can be helpful for medical services authorities and clinicians to break down diabetes utilizing ECG signals. Nesreen Samer El\_Jerjawi and Samy S. Abu-Naser [5] proposed a forecast model for diabetes utilizing ANN (Artificial Neural Network) that can be valuable for medical care official and professionals. The creator was spurred by the profoundly perilous confusion of the sickness. He fostered an ANN model for limiting the blunder work in the preparation. So the normal mistake work determined was 0.01% and precision achieved through ANN was 87.3%.

Sajida Perveen et al. [6] suggested AdaBoost innovation. An AdaBoost gathering model is better than the packing and J48 for the arrangement of a diabetic patient. The creator is enlivened through the flourishing effect of diabetes

everywhere on the world, in this way, the expectation and counteraction of diabetes mellitus are accomplishing importance in the medical care local area [7]. The creator introduced an expectation model with further developed execution for arrangement of Canadian populace diabetic patients across three distinct ages. There were three gathering models (sacking, AdaBoost and J48) which were applied on test information to assess the presentation and precision. Results show that AdaBoost beats others as far as precision. As indicated by creators AdaBoost can be applied to another infection like coronary illness, hypertension for better expectation.

Carrera et al. [8] proposed a computer-assisted procedure for the discovery of diabetic retinopathy, in light of the computerized signals handling of retinal pictures. The significant goal of this proposed approach is the classification of the situation of non-proliferative diabetic retinopathy at any of the retinal picture. The principle benefit of this methodology is that it is powerful in nature yet exactness and precision are expected to improve for the recorded application matter. Diabetes retinopathy is ongoing and has become the main way of life sickness. A since a long time ago run of this sickness can cause cardiovascular breakdown, kidney disappointment; inappropriate working of stomach, delayed raised glucose levels and some more. By considering this issue Huang et al. [9] proposed SVM and entropy techniques for the utilization of three distinctive datasets (diabetic retinopathy Debrecen, vertebral section, and mammographic mass) for estimating the precision. The creators tried the joined methodology of DL and SVM for the assessment of the preparation dataset layer by layer and the most basic trait was taken to develop the choice tree. The recommended strategy achieves promising grouping exactness. Thus, it has been seen that diabetic retinopathy Debrecen and mammographic mass gives more precise results and productivity.

### III. METHODOLOGY

#### A. Dataset

The dataset utilized for the examination is PIMA Indian dataset (PID) by NIDDK. The primary inspiration driving utilizing the PIMA dataset is that a large portion of the populace in this day and age follows a comparative way of life having a higher reliance on handled food varieties with a decrease in active work. PID is a drawn out associate examination since 1965 by NIDDK in light of the greatest danger of diabetes. The dataset contained certain demonstrative boundaries and estimation through which the patient can be related to any sort of constant illness or diabetes before time. The entirety of the Participants in PID is females and somewhere around 21 years of age. PID made out of a sum of 768 examples, from which 268 examples were recognized as diabetic and 500 were non-diabetics. The 8 most impacting ascribes that contributed towards the forecast of diabetes are as per the following: a few pregnancies the patient has had, BMI, insulin level, age, Blood Pressure, Skin thickness, Glucose, Diabetes Pedigree Function with mark result (0 and 1).

The Pima Indian dataset is taken from the URL [https://data.](https://data.world/data-society/pima-indians-diabetes-database)

[world/data-society/pima-indians-diabetes-database](https://data.world/data-society/pima-indians-diabetes-database) diabetes-data set and parts in a 80/20% proportion into the preparation and approval set. The approval part is 20% of the information dataset which has been chosen to coordinate the determination of hyper parameters.

In fact, the approval set performs preparing of hyper-parameters before the improvement [1]. Crossvalidation has been utilized for assessing the measurable exhibition of the learning model. It executes two sub-measures as testing and preparing. The preparation sub process is utilized to prepare a model and afterward the learning model is applied in the Testing sub process to gauge the precision. justification picking Pima Indian dataset is the high predominance of type 2 diabetes in the Pima gathering of Native Americans living in the space which is presently known as focal and southern Arizona. This gathering has made due with a less than stellar eating routine of starches for quite a long time due to the hereditary inclination [2]. As of late, the Pima bunch acquires a high sign of diabetes because of the abrupt shift from conventional yields to handled food sources.

#### B. Data Preprocessing

Data preprocessing is most significant cycle. Generally medical care related information contains missing value and different pollutions that can cause viability of information. To work on quality and viability got in the wake of mining measure, Data preprocessing is finished. To utilize Machine Learning Techniques on the dataset adequately this cycle is fundamental for precise outcome and effective forecast. For Pima Indian diabetes dataset we need to perform pre preparing in two stages.

Missing Values expulsion Remove every one of the examples that have zero (0) as worth. Having zero as worth is preposterous. Consequently this case is disposed of. Through disposing of insignificant highlights/examples we make include subset and this interaction is called highlights subset determination, which lessens diamentonality of information and help to work quicker.

Parting of information After cleaning the information, information is standardized in preparing and testing the model. At the point when information is spitted then we train calculation on the preparation informational collection and keep test informational index to the side. This preparation interaction will create the preparation model dependent on rationale and calculations and upsides of the element in preparing information. Fundamentally point of standardization is to bring every one of the properties under same scale.

#### C. Random Forest

It is kind of group learning technique and furthermore utilized for arrangement and relapse errands. The exactness it gives is grater then, at that point contrasted with different models. This strategy can without much of a stretch handle huge datasets. Arbitrary Forest is created by Leo Breiman. It is mainstream group Learning Method. Arbitrary Forest Improve Performance of Decision Tree by decreasing change. It works by building a large number of choice trees at preparing time and yields the class that is the method of the classes or grouping or mean forecast (relapse) of the individual trees.

- The initial step is to choose the "R" highlights from the absolute highlights "m" where  $R \ll M$ .

- Among the "R" includes, the hub utilizing the bestsplit point.
- Split the hub into sub hubs utilizing the best split.
- Repeat a to c strides until "l" number of hubs has been reached.
- Built backwoods by rehashing stages a to d for various occasions to make "n" number of trees.

The initial step is to require the take a look at decisions and utilize the establishments of each unpredictably made choice tree to foresee the outcome and stores the expected result at stretches the objective spot. Furthermore, figure the decisions in favor of each anticipated objective and eventually, concede the high casted a ballot anticipated objective because of a definitive forecast from the arbitrary timberland equation. A portion of the alternatives of Random Forest remedies forecasts result for a spread of utilizations are advertised.

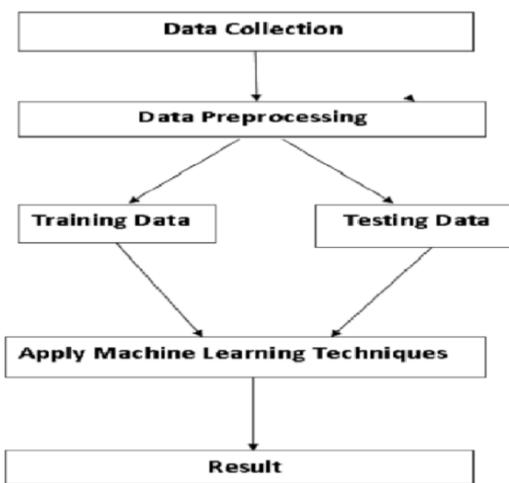


Figure 1: Overview of the Process

**D. Model Building**

This is most significant stage which incorporates model structure for forecast of diabetes. In this we have carried out different AI calculations which are talked about above for diabetes forecast.

- Step1: Import required libraries, Import diabetes dataset.
- Step2: Pre-measure information to eliminate missing information.
- Step3: Perform rate split of 80% to isolate dataset as Preparing set and 20% to Test set.
- Step4: Select the AI calculation for example Support Vector Machine, Decision Tree, Random Forest and Gradient boosting calculation.
- Step5: Build the classifier model for the referenced machine learning calculation dependent on preparing set.
- Step6: Test the Classifier model for the referenced machine learning calculation dependent on test set.
- Step7: Perform Comparison Evaluation of the trial execution results acquired for every classifier.
- Step8: After investigating dependent on different measures finish up the best performing calculation.

**IV. RESULTS & DISCUSSION**

In this examination work, results were accomplished by applying four grouping calculations (DT, SVM, XGboost, and

RF) to show boost precision in diabetes expectation. From these four classifiers, RF and Xgboost give promising precision (94.07%) which can be demonstrated as an unmistakable apparatus for the expectation of diabetes at a beginning phase. In our proposed framework we utilize the PIMA dataset and apply it's anything but a DL approach. Further, it can help the medical services expert and can be the second assessment for the improvement of choices relying upon extricated highlights. Numerous scientists have been recently chipped away at the PIMA dataset with an assorted calculation to anticipate diabetes. Accordingly a portion of the scientist's work has been addressed with their applied strategies and accomplished precision. Table 1 to table 4 shows all the promising work done on Pima dataset till time and our proposed technique accomplished the most elevated exactness for example on PIMA Indian dataset.

Table 1 Decision tree (Accuracy 89.34%)

| Actual Values<br>Predicted Values | True No | True Yes | Class Precision |
|-----------------------------------|---------|----------|-----------------|
| Predicted No                      | 137     | 4        | 97.16%          |
| Predicted yes                     | 3       | 63       | 95.45%          |
| Class recall                      | 97.86%  | 94.03%   |                 |

Table 2 Support vector machine (Accuracy 83.54.34%)

| Actual Values<br>Predicted Values | True No | True Yes | Class Precision |
|-----------------------------------|---------|----------|-----------------|
| Predicted No                      | 128     | 8        | 94.12%          |
| Predicted yes                     | 12      | 59       | 83.10%          |
| Class recall                      | 91.43%  | 88.06%   |                 |

Table 3 Random Forest (Accuracy 92.34%)

| Actual Values<br>Predicted Values | True No | True Yes | Class Precision |
|-----------------------------------|---------|----------|-----------------|
| Predicted No                      | 139     | 3        | 97.89%          |
| Predicted yes                     | 1       | 64       | 98.46%          |
| Class recall                      | 99.29%  | 95.52%   |                 |

Table 4 XGboost (Accuracy 90.34%)

| Actual Values<br>Predicted Values | True No | True Yes | Class Precision |
|-----------------------------------|---------|----------|-----------------|
| Predicted No                      | 118     | 27       | 81.38%          |
| Predicted yes                     | 22      | 40       | 64.52%          |
| Class recall                      | 84.29%  | 59.70%   |                 |

The exactness got through assorted classifiers is displayed beneath by the disarray framework which comprises of class accuracy, diabetes expectation indeed, diabetes forecast no, class review.

In this work we see that arbitrary woodland classifier accomplishes better contrasted with others. By and large we have utilized best Machine Learning methods for expectation and to accomplish elite precision. Figure 2 shows the consequence of these Machine Learning strategies.

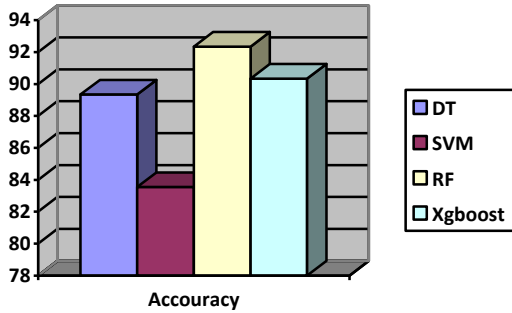


Figure 2: Comparison of RF with other ML algorithms

Here include assumed significant part in forecast is introduced for irregular backwoods calculation. The amount of the significance of each element assuming significant part for diabetes have been plotted, where X-hub addresses the significance of each component and Y-Axis the names of the highlights Figure 3.

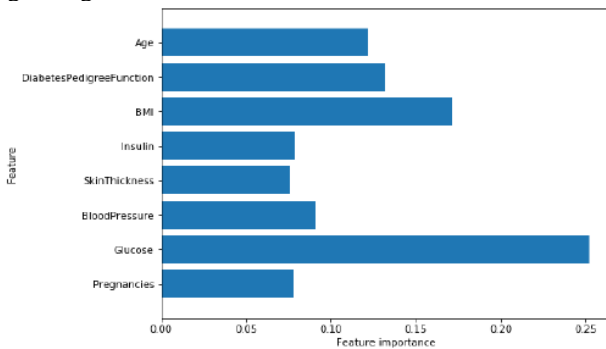


Figure 3: Feature Importance Plot for Random Forest

V. CONCLUSION

This paper intended to carry out an expectation model for the estimation of diabetes. As talked about before, an enormous piece of the human populace is in the hold of diabetes infection. On the off chance that stays untreated, it's anything but a colossal danger for the world. Thusly In our proposed research, we have tried assorted classifiers on the PIMA dataset and demonstrated that information mining and AI calculation can lessen the danger factors and work on the result as far as proficiency and precision. The result accomplished on the PIMA Indian dataset is higher than other

proposed strategies on the equivalent dataset utilizing information mining calculations as talked about in Table 1. Precision accomplished by the four classifiers (DT, SVM, XGboost, and RF) exists in the reach 90–92.34% which is impressively high than accessible techniques. Among the four proposed classifiers, RF is considered as the most effective and promising for examining diabetes with a precision pace of 92.34.07%. Later on, we expect to foster a vigorous framework as an application or a site that can utilize the proposed RF calculation to help medical care experts in the early recognition of diabetes.

ACKNOWLEDGMENT

We would like to thank our Guide Dr Jayanth J for his consistent support in completing the work

REFERENCES

- [1] "Global Report on Diabetes, 2016". Available at: [https://apps.who.int/iris/bitstream/handle/10665/204871/9789241565257\\_eng.pdf;jsessionid=2BC28035503CF9FF295E70CFB4A0E1DF?sequence=1](https://apps.who.int/iris/bitstream/handle/10665/204871/9789241565257_eng.pdf;jsessionid=2BC28035503CF9FF295E70CFB4A0E1DF?sequence=1).
- [2] "Diabetes: Asia's 'silent killer'", November 14, 2013". Available at: [www.bbc.com/news/world-asia-24740288](http://www.bbc.com/news/world-asia-24740288).
- [3] Mathers CD, Loncar D. Projections of global mortality and burden of disease from 2002 to 2030. 2015;3(11). <https://doi.org/10.1371/journal.pmed.0030442>.
- [4] Swapna G, Vinayakumar R, Soman KP. Diabetes detection using deep learning algorithms. ICT Express. 2018;4(4):243–6. <https://doi.org/10.1016/j.ict.2018.10.005>. Elsevier B.V.
- [5] El-Jerjawi NS, Abu-Naser SS. Diabetes prediction using artificial neural network. International Journal of Advanced Science and Technology. 2018;121:55–64. [doi.org/10.14257/ijast.2018.121.05](https://doi.org/10.14257/ijast.2018.121.05).
- [6] Perveen S, Shahbaz M, Keshavjee K, Guergachi A. Metabolic syndrome and development of diabetes mellitus: predictive modeling based on machine learning techniques, IEEE Access. 2019;7: 1365–75. <https://doi.org/10.1109/ACCESS.2018.2884249>.
- [7] Perveen S, et al. Performance analysis of data mining classification techniques to predict diabetes. Procedia Computer Science. 2016;82:115–21. <https://doi.org/10.1016/j.procs.2016.04.016> Elsevier Masson SAS.
- [8] Carrera EV, Carrera R. Automated detection of diabetic retinopathy using SVM, 2017. pp. 6–9.
- [9] "Machine Learning: Pima Indians Diabetes", April 14, 2018. Available at: <https://www.andreagranti.it/2018/04/14/machinelearning-pima-indians-diabetes/>.
- [10] Palaniappan S. Intelligent heart disease prediction system using data mining techniques, (march 2008). 2017. <https://doi.org/10.1109/AICCSA.2008.4493524>.