

## Mellitus Disease Prediction Using Machine Learning Algorithms

<sup>1</sup>Dr.Aziz Makandar, Dept. of Computer Science, Karnataka State Akkamahadevi Women's University, Vijayapur.

<sup>2</sup>Ms.Chaitra Sahadev Kohalli, PG Scholar, Dept. of Computer Science, Karnataka State Akkamahadevi Women's University, Vijayapur.

**Abstract** — This research investigates five K-Nearest Neighbor, Naive Bayes, Decision Tree Classifier, Random Forest, and Support Vector Machine are examples of popular AI algorithms that are used to predict diabetes disease. By integrating all of the dataset's previous gambling sections after describing and putting cross-validation into practice, we also found a consistent precision. All other classifiers provided a constant precision of above 70%, however the KNN classifier allowed us to reach the greatest and most reliable precision of 76%. We looked into why some Machine Learning classifiers do not provide consistent and high precision by imagining the preparation and evaluating exactness, as well as by inspecting model overfitting and model underfitting. The primary objective of this study is to obtain the most accurate results for the diagnosis and prognosis of precision and computing time.

**Record Terms** — Diabetes infection, Machine Learning (ML), Confusion Matrix, Scikit-Learn, Body Mass Index (BMI), Precision, Recall, F1-Score, Pandas, NumPy, and Python.

Date of Submission: 10-10-2022

Date of acceptance: 22-10-2022

### I. INTRODUCTION

We explored by examining model prediction error and model underfitting, as well as visualising the preparation and evaluating exactness, we may understand why some machine learning classifiers do not give consistent and high accuracy. The main goal of this study is to produce the best diabetes illness prediction results possible in terms of accuracy and calculation speed.

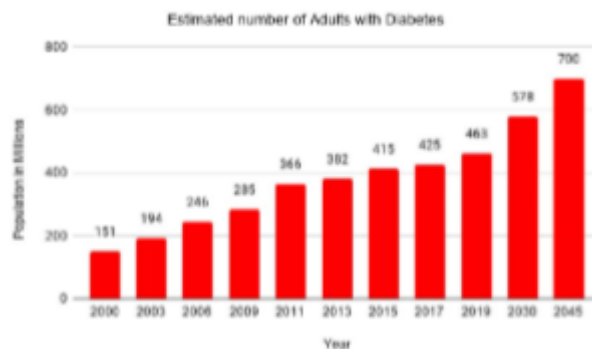


Fig. 1. Number of diabetic patients assessed as for year [1]

In the year 2019, around 463 million grown-ups matured 20 to 79 have diabetes (Global Diabetes Association IDF). 79% of the grown-up populace lived in nations with low and center pay gatherings. It is assessed that by 2045, very nearly 700 million individuals would have diabetes (IDF). Diabetes is spreading all over the planet because of regular and acquired factors. The numbers are quickly expanding because of various variables, including horrible food sources, genuine deferral, and others. Diabetes is a hormonal condition where the body's powerlessness to make insulin causes strange sugar digestion in the body, lifting blood glucose levels in the body of a particular individual. The recognizable characteristics incorporate areas of strength for a, thirst, and an unquenchable inclination to pee. Certain gamble factors, for example, age, BMI, glucose levels, pulse, etc, assume a crucial part in the contamination's responsibility. We can find in Fig. 1 that the quantity of cases is ceaselessly expanding and that the unique cases are not dialing back. It is critical to underline this since diabetes has become seemingly the most risky and fast disease that has killed numerous people everywhere. Computer

based intelligence is profoundly notable these days since it is utilized anyplace there is a ton of information, and we frantically need information from it. By and large, we can separate AI calculations into two classes, but this isn't comprehensive.

- Unsupervised Learning: In unsupervised learning, the data is not identified or prepared. We just applied the facts in actual life to come up with a few examples, if possible.

- Supervised Learning: In supervised learning, we train the model using the names associated with the data, and then group or test the new information with marks.

With the ascent of AI and its connected calculations, it has become evident that the significant difficulties and detours in acknowledgment existed beforehand may now be overwhelmed easily, while as yet giving a certain and precise end. At this point, it is perceived that AI has become altogether more effective and versatile in a joint effort with the clinical field. By zeroing in on an individual's qualities, simulated intelligence can give early recognition of a disease. Such early measures can both defer the beginning of sickness and go about as an obstacle to empowering the contamination to spread to a principal level. The objective of the work introduced in this paper is to reproduce diabetic sickness involving man-made intelligence computations for early thought of an individual.

## **RELATED WORKS**

They examined information for diabetes illness expectations based on Big Data of medical services using the WEKA device in [2]. They used a dataset from UCI that was freely available and applied various AI classifiers on it. They utilized Credulous Bayes, Backing Vector Machine, Irregular Timberland, and Basic Truck as classifiers. They started by acquiring the dataset, preprocessing it in the Weka contraption, and afterward doing the 70:30 train and test split for applying different machine computations. They didn't continue with the cross-endorsement adventure on the grounds that getting the best and most exact outcomes is additionally basic. The specialists in [3] likewise utilized a publically accessible dataset called the Pima Indians Diabetes Data set to lead their examination. The structure of how they play out the assumption begins with dataset choice and closures with data pre-taking care of. They utilized three arrangement computations, like straightforward Bayes, SVM, and Choice tree, in the wake of preprocessing the information. The most significant level of accuracy they had the option to accomplish with their examination was 76.30 percent. They haven't drilled Cross-endorsement, either, as [2]. The makers offered a cerebrum network-based diabetes sickness forecast concerning the Indians Pima Diabetes Dataset in [4]. They utilized a couple of mystery layers to identify designs in the information, and they anticipated the result utilizing those models. They call their proposed calculations ADAP, which is a custom cerebrum network with different parts and related burdens and units organized. They had the option to accomplish a half breed point of 0.76 for responsiveness and explicitness, and are presently endeavoring to repeat their outcomes later on. The creators of [5] utilized an assortment of artificial intelligence calculations, including assist vector with machining, unpredictable woodlands, key backslide, Choice tree, and others, as well as an assortment of sickness datasets, to show the utility of AI in illness expectation and examination. They additionally utilized the conventional procedure to coordinating the assessment, which included information pre-taking care of, feature extraction and determination, classifier arrangement, and testing for the end results. They utilized incorporate choice to reduce down on computational expenses. They likewise divided each dataset into a 90% preparation put and a 10% testing set together to accomplish the best outcomes. They finished cross-endorsement for every calculation, as well as the precision measure, and exhibited changed results considering different k attributes for k-overlay cross-endorsement.

## **PROPOSED MODEL**

We used the Pima Indians Diabetes Database [4], a freely downloadable dataset, to carry out our research. This dataset contains a distinct proportion of diabetic illness symptoms. The first information was supplied by the National Institute of Diabetes and Digestive and Kidney Diseases. Every single incidence that has been documented involves patients who are older than 21. The five steps in our proposed model are depicted in Fig. 2.

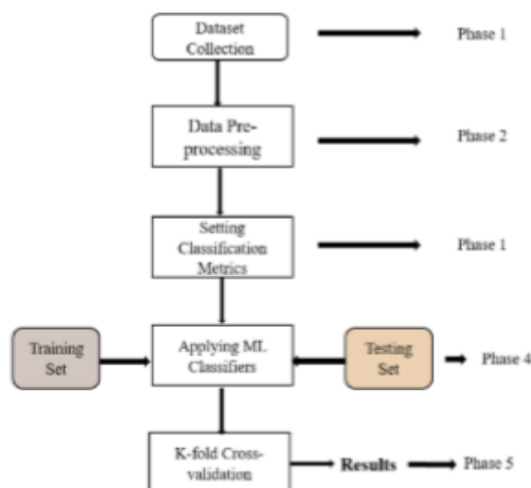


Fig. 2. Various periods of our trial.

**A. Information Collection**

Eight highlights from the dataset mentioned above [4] are listed in Table I.

**TABLE I  
LIST OF FEATURES PRESENT IN THE DATASET**

Features	Description
<b>Pregnancies</b>	Number of Pregnancies patients had earlier.
<b>Glucose</b>	Glucose level present in the patient.
<b>Blood Pressure</b>	Recorded blood pressure level at that particular time.
<b>Skin Thickness</b>	Skin thickness level of the patient.
<b>Insulin</b>	Amount of insulin present in the body.
<b>BMI</b>	Body Mass Index of the individual.
<b>Diabetes Pedigree Function</b>	Family history of Diabetes disease.
<b>Age</b>	Age of an individual.

The dataset also includes two names (0-No and 1-Yes) that are the consequence of the infection with diabetes in addition to the component. The supplementary sections go through the specifics of each quality or point of view.

- **Pregnancy:** Women with diabetes mellitus are destined to develop type 2 diabetes in the future. More deliveries increase the risk of developing diabetes in a person
- **Glucose:** The participants underwent an oral glucose tolerance test in which their blood sugar levels were controlled. They were instructed to check their plasma glucose fixation after two hours. Diabetes is certain to develop in people with higher significant levels.
- Circulatory strain:** A higher risk of developing diabetes is associated with having a heartbeat that is more than 140/90 mmHg of Mercury. However, some people who have a diastolic blood pressure < 70 mmHg may develop diabetes.
- **Collagen content** is a key factor in determining skin thickness, which is higher in diabetics with insulin-subordinate skin. The skin overlay on the subjects' back arm muscles was measured, and the results showed that individuals with skin thickness of at least 30mm are at a higher risk.
- **Insulin:** After 2 hours of glucose association, normal insulin levels range from 16 to 166 mIU/L. Subjects are more at danger depending on their level of self-esteem.
- People with a BMI (body mass index) greater than 25 have a marginally increased risk of acquiring diabetes.
- The diabetes pedigree function provides "a synthesis of family members' hereditary relationships to the subject as well as their history of insulin resistance." An individual is more prone to get diabetes the higher their DPF.
- **Age:** Although diabetes can afflict persons of any age, middle-aged adults are the most likely to be diagnosed with it (45 onwards). As a result, people who are older have a higher risk of developing diabetes.

**Pre-handling of Information**

The previously described dataset has vanished and revealed data. To make the dataset useable and extract information from it, we employed information preprocessing. We physically corrected the dataset's unexpected regions after breaking it apart to deal with inaccurate data. By computing the standard deviation of that particular element and applying it to the vacant slots, missing attributes are made up for. We used the Pandas [6] and NumPy [7] libraries to effectively handle the dataset and keep basic information all through the study in order to make the dataset valuable.

**C. Setting Classification Metrics**

We really want to set a few of measurements that will help us predict the Diabetes illness in order to describe the illness and achieve a prediction outcome. We used disarray grid as the classification measure measurements since we are doing our experiment with the scikit-learn (Sklearn) AI framework [8]. This page includes a list of all previously used metrics, including Accuracy, Recall, F1-Score, and Accuracy in our test.

- Precision (P) is defined as the quantity of genuine up-sides (Tp) over the quantity of genuine up-sides in addition to the quantity of bogus up-sides (Fp). Numerically,

$$P = \frac{Tp}{Tp + Fp} \quad (1)$$

- Review (R) is defined as the quantity of genuine up-sides (Tp) over the quantity of genuine up-sides in addition to the quantity of misleading negatives (Fn).

$$R = \frac{Tp}{Tp + Fn} \quad (2)$$

- F1-Score (F1) is defined as the symphonious mean of accuracy and review.

$$F1 = \frac{2 * P * R}{P + R} \quad (3)$$

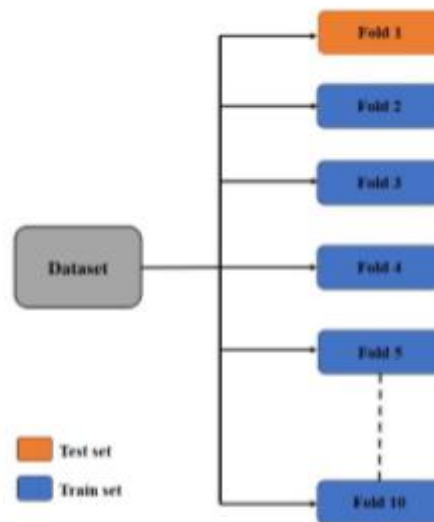
- Exactness (A) is defined as follows.

$$A = \frac{Tp + Tn}{Tp + Tn + Fp + Fn} \quad (4)$$

**B. Applying Machine Learning**

Algorithms On the pre-handled dataset, we will run 5 directed machine calculations for our investigation. The following are the calculations we used:

- 1) K-Nearest Neighbor (KNN) with K=10
  - 2) Naive Bayes (NB)
  - 3) Decision Tree (DT)
  - 4) Random Forest (RF)
  - 5) Support Vector Machine (SVM)
- E. K-Fold Cross Validation



**Fig. 3. 10-Fold Cross-Validation**

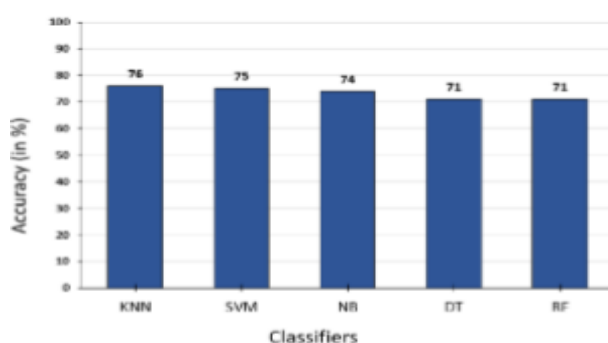
Finally, we used a different assessment technique, To guarantee that the dataset was handled appropriately and that the best accuracy results were produced, K-Fold cross-approval was performed [9]. The dataset is divided into K separate folds (in our case, K=10) for K fold cross-approval, and one of the overlaps, let's say Fold-1, is examined with the additional (k-1) folds in each cycle. Until every fold has been investigated, this cycle will go on forever. To further comprehend this assessment, we offer a visual representation of K-Fold Cross-approval measurements in Fig. 3.measurement.

## II. EXPLORATORY RESULTS

To conduct our investigation, we divided the dataset into two sections: preparation and testing, each with an 80:20 ratio. With a framework configuration of 8 GB RAM We applied all of our AI classifiers using Python 3.6 adaptation and an Intel i5 ninth generation portable CPU.

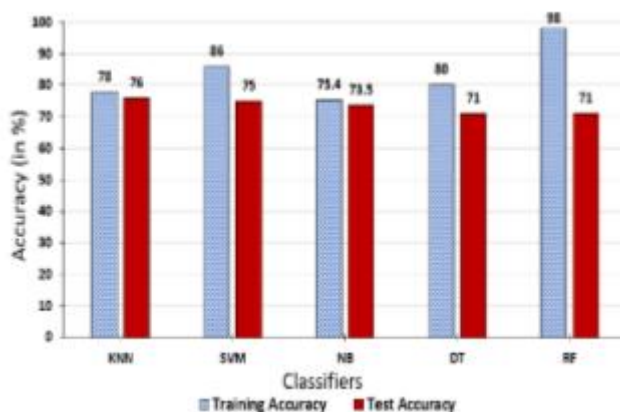
**TABLE II**  
**PERFORMANCE METRICS OF DIFFERENT CLASSIFIERS MODELS**

Classifier	P	R	F1	A(10-fold)
KNN	0.76	0.73	0.75	0.76
SVM	0.73	0.74	0.73	0.75
NB	0.74	0.74	0.74	0.74
DT	0.72	0.71	0.71	0.71
RF	0.70	0.71	0.71	0.71



**Fig. 4. Exactness of various classifiers after 10-overlap cross-approval.**

In Fig. 4, we can see that after carrying out 10-overlap cross-approval, KNN classifiers acquired the maximum precision of 76% and that other classifiers were able to reach excellent correctness of over 70% as well. The itemised data for each measurement that we collected during our inquiry are presented in Table II. The precision of the preparation and testing were then carefully compared. These correctness are crucial in figuring out whether our model can handle overfitting or underfitting. The specific classifier models are improving more from the test data than from the training data, which implies that our model is overfitting. Our algorithm, however, experiences model underfitting if the testing precision is higher, and the models are gaining far fewer standards from the preparation data, and they are unable to construct associations with test data. We need identical planning and testing exactness for a good result. We witnessed a preparation and testing exactness bring about Fig. 5, where all the preparation correctnesses were more significant than testing exactnesses, indicating that our classifiers are suffering from model overfitting and are advancing more instructions based on testing data.



**Fig. 5. Exactnesses of various classifiers after 10-overlap cross-approval.**

To achieve higher precision in order to prevent model overfitting and obtain better execution in terms of exactness and accuracy, we need additional cases (patient information) in the dataset that can establish connections between the data preparation and testing. Additionally, we can observe that the KNN has the lowest preparation and testing exactness distinction, indicating that it experiences the least structural model overfitting and will be the most likely classifier to predict diabetes infection, taking all relevant parameters into account. Of all AI classifiers, the SVM and RF classifiers experience model overfitting because of their extensive precision training and testing. Finally, we carefully examine the computing times of each classifier. The preparation and execution times were computed using the Python time module constructing (gaining knowledge from available data) an AI computation or specific model. On the other hand, testing time is the period of time used to examine the outcomes after contrasting newly generated data with previously generated data. The recipes listed below are used to determine computational time (Ct).

$Ct = Tt + Ts$  (5) where  $Tt$  and  $Ts$  addresses preparing time and testing time, individually.

**TABLE III**  
**Supercomputing TIME (IN SEC.) AND TRAINING AND TESTING TIME (IN SEC.) OF DIFFERENT CLASSIFIERS**

Classifier	Training Time	Testing Time	Computational Time
<b>KNN</b>	2.7549004457	5.1542658899	7.90916633
<b>NB</b>	0.18233000002	0.0567890879	0.23911908
<b>DT</b>	0.77532997643	0.0478995599	0.82322953
<b>RF</b>	12.6574437890	0.9546667586	13.6121105
<b>SVM</b>	670.49820000	15.954679321	686.452879

Table III shows that the SVM classifier takes up the majority of the computing time and is relatively resource intensive in terms of memory and CPU use. Additionally, it does not offer higher precision in comparison to other KNN classifiers, which use less CPU and memory and require a lot less computation time. Also cannot be regarded as a fair classifier because it requires the most computation time when compared to the SVM classifier and experiences model overfitting for predicting illness.

### III. CONCLUSION

The early detection of a condition, one of the main barriers to the development of innovation and therapy in this scenario is diabetes. However, major efforts were undertaken in this review to develop a model that is accurate enough to pinpoint the illness's beginning. The tests done on the Pima Indians Diabetes Database allowed us to predict this illness. Furthermore, using K-Nearest Neighbors classifiers, the results obtained demonstrated the system's suitability with a 76 percent accuracy. With this in mind, we are sure that we can use our model as a foundation for forecasting more fatal illnesses. The mechanisation of diabetes or other illness tests may yet have space for development in the future. In the future, we'll work with an urgent care facility or a clinical foundation to develop a diabetic dataset in an effort to achieve better outcomes. To get better results, we'll also add more machine learning and deep learning models.

### REFERENCES

- [1]. P. Saeedi, I. Petersohn, P. Salpea, B. Malanda, S. Karuranga, N. Unwin, S. Colagiuri, L. Guariguata, A. A. Motala, K. Ogurtsova, J.E.Shaw,D.Bright,andR.Williams,“Globalandregional diabetes prevalence estimates for 2019 and projections for 2030 and 2045: Results from the international diabetes federation diabetes atlas, 9th edition,” Diabetes Research and Clinical Practice, vol. 157, p. 107843, 2019.
- [2]. A. Mir and S. N. Dhage, “Diabetes disease prediction using machine learning on big data of healthcare,” in 2018 Fourth International Conference on Computing Communication Control and Automation (ICCUBEA), 2018, pp. 1–6.
- [3]. D. Sisodia and D. S. Sisodia, “Prediction of diabetes using classification algorithms,” Procedia Computer Science, vol. 132, pp. 1578 – 1585, 2018, international Conference on Computational Intelligence and Data Science. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S1877050918308548>
- [4]. J. Smith, J. Everhart, W. Dickson, W. Knowler, and R. Johannes, “Using the adap learning algorithm to forecast the onset of diabetes mellitus,” Proceedings - Annual Symposium on Computer Applications in Medical Care, vol. 10, 11 1988.
- [5]. P. S. Kohli and S. Arora, “Application of machine learning in disease prediction,” in 2018 4th International Conference on Computing Communication and Automation (ICCCA), 2018, pp. 1–4.
- [6]. Wes McKinney, “Data Structures for Statistical Computing in Python,” in Proceedings of the 9th Python in Science Conference, Stéfan van der Walt and Jarrod Millman, Eds., 2010, pp. 56 – 61.
- [7]. C. R. Harris, K. J. Millman, S. J. van der Walt, R. Gommers, P. Virtanen, D. Cournapeau, E. Wieser, J. Taylor, S. Berg, N. J. Smith, R. Kern, M. Picus, S. Hoyer, M. H. van Kerkwijk, M. Brett, A. Haldane, J. F. del Río, M. Wiebe, P. Peterson, P. G'erard-Marchant, K. Sheppard, T. Reddy, W. Weckesser, H. Abbasi, C. Gohlke, and T. E. Oliphant, “Array programming with NumPy,” Nature, vol. 585, no. 7825, pp. 357–362, Sep. 2020. [Online]. Available: <https://doi.org/10.1038/s41586-020-2649-2>

- [8]. F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and ' Edouard Duchesnay, "Scikit-learn: Machine Learning in Python," *Journal of Machine Learning Research*, vol. 12, no. 85, p. 28252830, 2011.
- [9]. S. Yadav and S. Shukla, "Analysis of k-fold cross-validation over hold-out validation on colossal datasets for quality classification," in *2016 IEEE 6th International Conference on Advanced Computing (IACC)*, 2016, pp. 78–83.