

Manual Diabetes Prediction System.

Saurabh Bisht

^{1,2} PG Student, Dept. of School of Computing, Graphic Era Hill University, Dehradun, Uttarakhand.

³ Assistant Professor, Dept. of School of Computing, Graphic Era Hill University, Dehradun, Uttarakhand.

Abstract- *Diabetes is a common health condition that can cause serious problems if it isn't caught early. In this research, we developed a system to predict the risk of diabetes using machine learning. The system looks at important factors like blood sugar levels, BMI, age, and family history to figure out if someone might be at risk. We made sure the data used in the system is clean and consistent to give accurate results. We also tested different machine learning techniques, such as logistic regression, decision trees, and neural networks, to find the best one. The result is a user-friendly tool that provides accurate predictions to help doctors and patients take action sooner. This system aims to make early detection easier, improve health outcomes, and reduce the impact of diabetes in communities..*

INTRODUCTION

Diabetes is a serious health problem that affects millions of people around the world. It happens when blood sugar levels are too high, which can lead to major health issues like heart disease, kidney damage, nerve problems, and vision loss. The sooner diabetes is detected and treated, the better the chances of preventing these complications.

Currently, diabetes is usually diagnosed through regular check-ups and lab tests, but these can sometimes be expensive or not easily accessible. With modern technology and the availability of health data, machine learning has become a powerful tool for predicting health risks. By analyzing patterns in health information, machine learning can help predict if someone is at risk of developing diabetes. This allows for earlier action, which can make a big difference in managing the disease.

In this research, we focus on creating a system that predicts the risk of diabetes using machine learning. The system looks at important details like blood sugar levels, BMI, age, and family history to provide a reliable risk assessment. The goal is to make a tool that is easy for doctors and patients to use, helping them identify people who might need attention before their condition worsens.

To achieve this, we tested different machine learning methods, such as logistic regression, decision trees, and neural networks, to find the best one for accurate predictions. By building this system, we aim to make diabetes detection faster and more accessible, improving health outcomes and reducing the complications caused by the disease.

LITERATURE SURVEY

Diabetes prediction has been a big consciousness in healthcare research, with many research exploring approaches to improve early detection the use of advanced technologies. Traditional techniques for predicting diabetes usually involve reading scientific and demographic factors through statistical fashions. While these methods are dependable, they frequently lack the precision and flexibility wanted for big-scale or personalised predictions.

Increased dynamics in system gaining knowledge of tactics for diabetes prediction. Researchers have started making use of logistic regression, choice bushes, SVM, and neural networks in analyzing fitness statistics to hit upon diabetes patterns. These fashions can system huge data sets whilst uncovering relationships that would be hard to find otherwise, for this reason offering more accurate outcomes than traditional strategies.

For instance, a few have used logistic regression because of its simplicity and interpretability as a primary step inside the evaluation. However, fashions consisting of selection timber and random forests were located to have better accuracy in reading the facts because they could deal with non-linear relationships. Deep mastering models, in particular neural networks, had been of interest due to their capability to investigate big complex facts and supply splendid predictive overall performance. These fashions are regularly computationally pricey and require sizable expertise to be implemented. Apart from algorithm selection, information fine is the maximum essential aspect in a diabetes prediction machine. Researchers emphasize preprocessing statistics to address missing values, normalize variables, and pick out relevant features. Techniques which includes function engineering and dimensionality discount have been studied to improve version overall performance and make predictions easier.

Besides, real-world usability has been diagnosed to design structures which might be accurate in addition to user-friendly. Integrating prediction gear in fitness care regularly calls for a layout that focuses on how a healthcare issuer can work efficiently with the interface in question, as a consequence making these systems useful for early prognosis and custom designed treatment.

PROPOSED SYSTEM

The Diabetes prediction system's motivation is to enhance the accuracy while shortening the time spent on execution of the fatigue and fatigue Prediction system.

Classifiers.py

file contains implementations for at least two classifiers:

3.1. KNNClassifier (K-Nearest Neighbors): -

The KNNClassifier is a distance-based algorithm that predicts a class by analyzing the closest neighbors to a data point using Euclidean distance. It uses a parameter k to determine how many neighbors influence the prediction.

3.2. SVMClassifier (Support Vector Machine):-

Support Vector Machines (SVM) are powerful supervised learning algorithms commonly used for classification and regression tasks. The SVMClassifier in this project is implemented with the following key features:

Objective:

The most important goal of the SVM is to search for an ideal hyper-plane that splits the data to different instructions with the maximum-margin. Margin refers to the gap between the hyper-aircraft and the nearest records points coming from exceptional training, regularly called support vectors.

Mathematical Foundation:

Here, the classifier optimizes a cost function that now not best minimizes the mistakes in category however additionally maximizes the margin.

A regularization parameter (λ) is used to govern the exchange-off between attaining a large margin and minimizing misclassification at the schooling statistics.

Algorithmic Approach:

Initialization: The weights (w) and bias (b) are initialized to 0.

Optimization: The version iterates over the dataset a couple of instances ($n_{\text{iterations}}$) to replace the weights and bias the use of a gradient-primarily based approach.

The gaining knowledge of price (learning_rate) controls the step length of weight updates for the duration of each iteration.

Advantages:

SVM is very effective in high-dimensional areas and is reminiscent-green as it best requires using a subset of the education records, specifically the assist vectors, for making predictions.

It works nicely on linearly separable and non-linearly separable facts; the use of kernel capabilities further extends its potential to work with complex datasets.

Finally, the fully connected layers are in charge of classifying the input picture based on the retrieved characteristics. The layers enable the model to learn and make educated guesses about high-level representations by connecting every neuron from the layer before to every neuron in the current layer. A SoftMax layer typically receives the output of the fully linked layers and uses it to assign probability to each category (drowsy or alert) based on learned attributes.

CNN models may efficiently extract hierarchical representations from pictures and produce accurate assessments of driver sleepiness based on the recovered characteristics by mixing these different sorts of layers. [20] Because of its capacity to learn complicated properties from pictures Haar classifier is a programme that has been taught to recognise specific objects in photographs. It works by superimposing a "perfect" image on top of several unfavorable pictures. Usually, the training takes place on a computer server in stages. This device can determine a person's level of fatigue or drowsiness instantly, more accurately, and in less time. This implies that it can improve road safety by reducing collisions brought on by tired drivers.

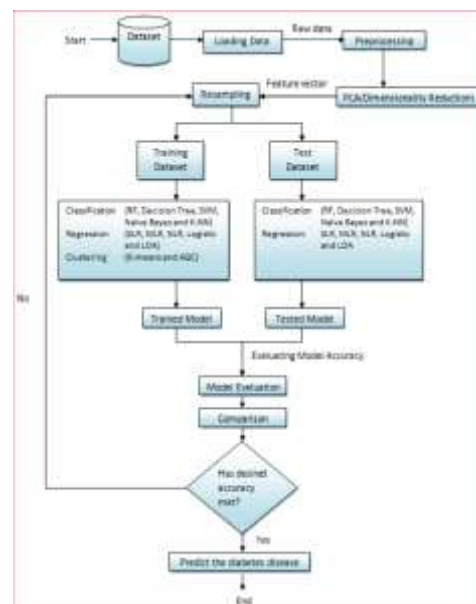


Figure 1. Diagram of Drowsiness Detection System.

IMPLEMENTATION

Data Processing-

- >The diabetes.csv dataset is loaded and explored to understand its structure.
- >The dataset is split into features (X) and labels (y), and further divided into training and testing sets.
- >Features are scaled using MinMaxScaler to normalize them between 0 and 1, ensuring better performance for machine learning models

.Classifiers Implementation-

Multiple classifiers are implemented, including KNNClassifier, SVMClassifier, and GradientDescentClassifier. These provide a comparative approach for predicting diabetes..

-KNNClassifier- The KNNClassifier is a distance-based algorithm that predicts a class by analyzing the closest neighbors to a data point using Euclidean distance. uses a parameter k to determine how many neighbors influence the prediction.

-SVMClassifier- Support Vector Machines (SVM) are powerful supervised learning algorithms commonly used for classification and regression tasks.

Training and Prediction-

The training process (fit method) adjusts model parameters based on the input data.

The prediction (predict method) uses trained parameters to classify new data points.

How It Works

Initialization:

The learning rate (learning_rate), regularization parameter (lambda_param), and iteration count (n_iterations) are set.

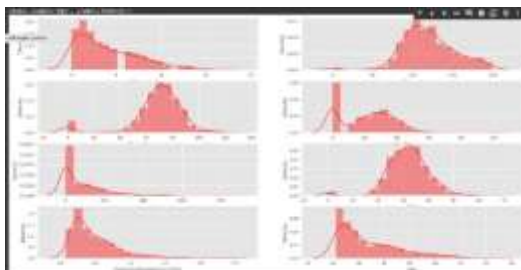
The weights (w) and bias (b) start as zeros.

Training:

- For each data point, calculate whether it satisfies the margin condition.
- If the margin condition is violated, update weights and bias to correctly classify the point.
- The model penalizes weights when misclassification occurs to maintain generalization.

Prediction:

- For new data points, calculate the linear combination ($w \cdot X + b$) to predict the class.
- Points above the hyperplane (linear combination ≥ 0) are classified as 1; otherwise, as 0.



Integration

- The SVMClassifier is integrated with the preprocessed training and testing data.
- Accuracy and performance metrics (e.g., accuracy score, confusion matrix) are used to evaluate predictions.
- Results are compared with other classifiers like KNN and Gradient Descent for final selection.

TESTING

The testing phase of the model evaluates its ability to generalize by using predicting effects for brand new, unseen information (the check dataset). Here's a detailed breakdown of how the testing phase works within the context of the diabetes prediction device:

Step-with the aid of-Step Testing Explanation

Step 1: Load and Scale Data

The take a look at dataset (X_{test}) undergoes the identical preprocessing steps as the schooling statistics to make sure consistency. For instance, features like glucose tiers and BMI are scaled between 0 and 1.

Step 2: Make Predictions

The model uses the are expecting function:

For each test pattern, it calculates the weighted sum of functions ($w \cdot X$) plus the prejudice (b).

If the end result is non-bad, the version outputs 1 (diabetic); in any other case, it outputs zero (non-diabetic).

Step three: Evaluate Predictions

Once predictions are made:

Accuracy measures the percentage of accurate predictions.

Confusion Matrix highlights the distribution of accurate and wrong classifications across each classes.

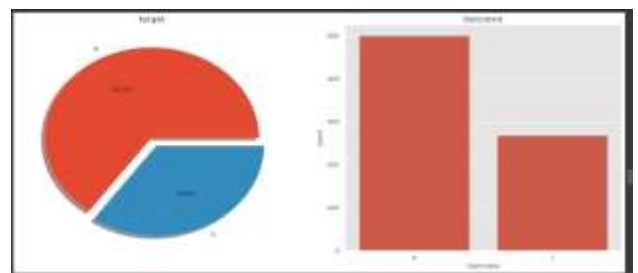
Precision and Recall examine the model's reliability in predicting wonderful cases (diabetes).

Step four: Interpret Results

The testing results screen how properly the version plays on unseen data:

A excessive accuracy indicates good overall overall performance.

A specific analysis of fake positives and false negatives helps refine the version if wanted.



Result

Model training result

The screenshot shows a Jupyter Notebook with the following content:

```
Support Vector Machine Evaluation
Training model...
Making predictions...

Results:
Accuracy: 78.78%
Precision: 88.88%
Recall: 8.52%
F1 Score: 15.76%

Confusion Matrix
True Positives: 3 False Positives: 0
False Negatives: 43 True Negatives: 106

Gradient Descent Evaluation
Training model...
Making predictions...

Results:
Accuracy: 86.52%
Precision: 78.88%
Recall: 88.88%
F1 Score: 88.52%

Confusion Matrix
True Positives: 28 False Positives: 12
False Negatives: 18 True Negatives: 98

Notice saved to model_metrics.json

Model Comparison Summary
Model Accuracy F1 Score
Support Vector Machine 78.78% 8.52%
Gradient Descent 86.52% 88.52%

Diabetes Prediction System
1. Train and evaluate models
2. Make new predictions
3. Exit

Enter your choice (1-3): 1
Loading training data...
Training models...
```

The notebook also shows a file explorer on the left with files like `diabetes.csv`, `diabetes_predict.py`, `main.py`, `model_metrics.json`, `package-lock.json`, `package.json`, `prediction.py`, `README`, `requirements.txt`, `utils.py`, and `visualization.py`.

Prediction Result

The screenshot shows a terminal window with the following output:

```
Enter your choice (1-3): 1
Please enter patient information:
Pregnancies: 2
Glucose: 125
BloodPressure: 23
SkinThickness: 5
Insulin: 32
BMI: 12
DiabetesPedigreeFunction: 35
Age: 34

Prediction Results:
XGB: Diabetic
SVM: Diabetic
GD: Diabetic

Consensus Prediction:
The patient is likely to be Diabetic

Diabetes Prediction System
1. Enter patient data manually
2. Exit

Enter your choice (1-2):
```

The terminal window also shows a file explorer on the left with files like `run.py` and `visualization.py`.

