

[Diabetes Prediction using Machine Learning]

Ebubechukwu Okonkwo
Bubebran@gmail.com

I. INTRODUCTION

The document begins by stating the importance of importing necessary libraries such as NumPy, Pandas, and scikit-learn's modules for preprocessing, model selection, support vector machine (SVM), and performance evaluation. It also mentions the use of matplotlib for inline plotting and warnings to manage warnings during execution. The core focus is on the "Diabetes dataset" and its analysis.

II. PROPOSED SOLUTION WITH JUSTIFICATIONS

The proposed solution involves using a Support Vector Machine (SVM) classifier to predict diabetes. The justification for using SVM, it was used because it is well-suited for high-dimensional datasets and works effectively in cases where the decision boundary between classes is complex. In medical diagnosis, where patterns in data might be non-linear, SVM can efficiently find the optimal boundary to separate diabetic and non-diabetic cases. Additionally, SVM performs well with smaller datasets and can handle outliers, making it a strong candidate for medical predictions. SVM is implemented in the "Training a model" section. The process includes:

- Loading the dataset into a Pandas DataFrame.
- Analyzing the dataset's shape, basic statistics, and class distribution.
- Preprocessing the data by standardizing it using StandardScaler.
- Splitting the dataset into training and testing sets.
- Training the SVM classifier on the training data.

III. RESULTS

The results section focuses on evaluating the model's performance:

- The classification report, including precision, recall, and f1-score, is presented.
- The accuracy of the model on both the training (78.66) and test data (77.27) is calculated and displayed.
- A confusion matrix is generated to visualize the model's predictions

IV. DISCUSSION

The results section focuses on evaluating the model's performance: The classification report, including precision, recall, and f1-score, is presented. The accuracy of the model on both the training and test data is calculated and displayed. A confusion matrix is generated to visualize the model's predictions.

V. CONCLUSION

What was achieved and how it can be applied to the real world

- What was achieved: A machine learning model using the Support Vector Machine (SVM) algorithm was developed to predict diabetes. The model's performance was evaluated, showing an accuracy of 78.66% on the training data and 77.27% on the test data.

- How it can be applied to the real world: The document illustrates a predictive system where, based on certain health parameters, the model predicts whether a person is diabetic or not. This system can be used in real-world healthcare settings for early diabetes detection, potentially enabling timely interventions and better patient outcomes.

A. Figures

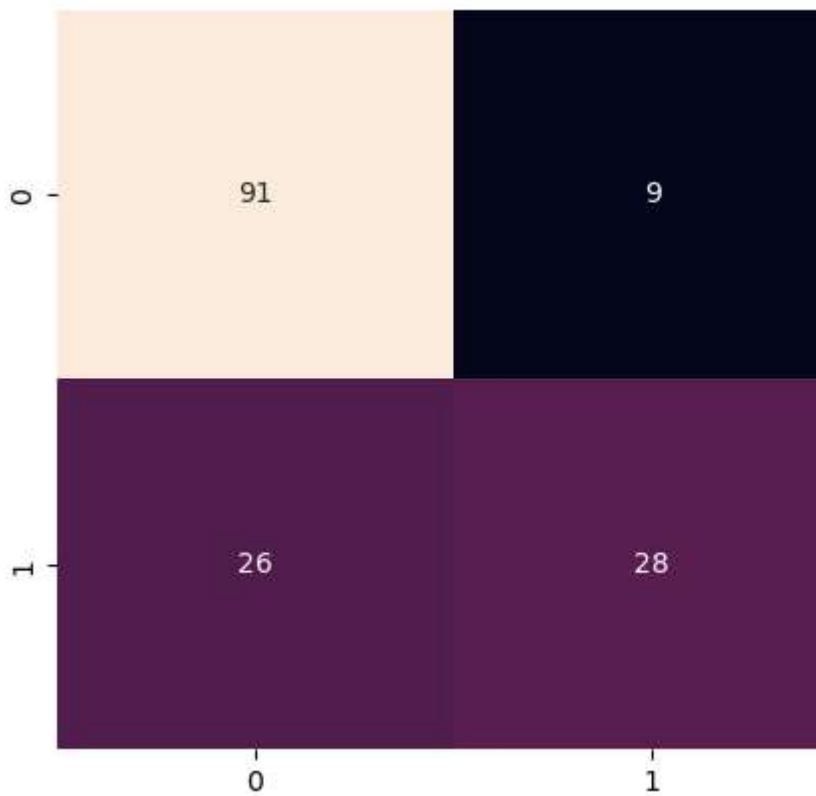


Fig. 1. heatmap representation of a confusion matrix

