# Predicting Diabetes Using Perceptron: A Comparative Study with SVM, Logistic Regression, and Random Forest

Rahul Kamarthi

The University of Adelaide

The University of Adelaide, South Australia 5005 Australia

a1938009@adelaide.edu.au

October 6, 2024

## Abstract

This study investigates the performance of a Single Layer Perceptron (SLP) for predicting diabetes and compares its results to well-established models such as Support Vector Machine (SVM), Logistic Regression, and Random Forest. Using the Pima Indians Diabetes dataset, we evaluated each model based on metrics like accuracy, precision, recall, and F1-score. The SLP demonstrated competitive performance, achieving an accuracy of 75.97%, making it a viable tool for early diabetes detection. Our feature importance analysis revealed that glucose levels and BMI were the most significant predictors of diabetes. The study highlights the potential of simpler models like SLP in medical diagnostics, while also recognizing its limitations compared to more complex models.

## 1 Introduction

Diabetes is a global health issue that, when detected early, can be managed to reduce long-term complications. Machine learning models have proven useful in predicting the onset of diseases like diabetes by analyzing clinical features. Traditional models such as Support Vector Machine (SVM), Logistic Regression, and Random Forest have consistently shown strong performance in medical predictions. However, there is still untapped potential in simpler models like the Single Layer Perceptron (SLP), which is a type of neural network often overlooked in this domain.

In this study, we aim to evaluate the SLP's effectiveness in predicting diabetes and compare it with more complex, traditional models. By utilizing the Pima Indians Diabetes dataset, we want to see if the SLP can perform on par with these models in terms of accuracy and other key metrics. Additionally, we aim to identify the most important features influencing diabetes prediction using feature importance analysis.

## 2 Related Work

Previous studies have explored various machine learning models for predicting diabetes, particularly with datasets like the Pima Indians Diabetes dataset. SVM is known for handling high-dimensional data effectively, while Logistic Regression remains popular due to its simplicity and interpretability. Random Forest, with its ensemble of decision trees, has also proven to handle non-linear data well.

Despite these models' strengths, less research has been done on neural network models like the Single Layer Perceptron for this particular task. The SLP's simplicity makes it a potentially powerful tool, particularly when computational efficiency is essential. In this study, we aim to bridge this gap by comparing the SLP's performance with these more established models.

# 3 Dataset and Data Preprocessing

## 3.1 The Dataset

We used the Pima Indians Diabetes dataset from Kaggle, which contains 768 samples, each with 8 features. These features include factors such as Pregnancies, Glucose, Blood Pressure, Skin Thickness, Insulin, BMI, Diabetes Pedigree Function, and Age, all of which are used to predict the likelihood of a patient developing diabetes.

## 3.2 Data Preprocessing

To ensure reliable model performance, we followed several preprocessing steps:

- **Handling Missing Values:** The dataset had some missing values for features like insulin levels. We addressed these missing values using a KNN imputer, which fills in missing data based on the average values of the nearest neighbors.

- **Scaling the Features:** Since models like SLP and SVM are sensitive to the scale of input features, we used StandardScaler to standardize the dataset. This ensured that all features had a mean of 0 and a variance of 1.

- **Splitting the Data:** The data was split into a training set (80%) and a test set (20%) to assess the models' performance on unseen data.

# 4 Methodology

We implemented four machine learning models—Single Layer Perceptron (SLP), SVM, Logistic Regression, and Random Forest—using Python's scikit-learn library. The goal was to compare these models based on their ability to predict diabetes, considering both their simplicity and predictive power.

## 4.1 Single Layer Perceptron (SLP)

The SLP is a basic neural network that works by learning linear decision boundaries. It uses a weighted sum of the inputs to produce a prediction, which is passed through a sigmoid activation function to map the result to a probability.

**Formula for Prediction:**

$$z = \sum_{i=1}^{n} w_i x_i + b \tag{1}$$

Where $w_i$ is the weight for each feature $x_i$, and $b$ is the bias term.

**Sigmoid Activation Function:**

$$\sigma(z) = \frac{1}{1 + e^{-z}} \tag{2}$$

The output of the sigmoid function is interpreted as a probability, with values above 0.5 classified as diabetic and values below 0.5 as non-diabetic.

**Training the SLP:** The SLP uses gradient descent to minimize the error between predicted and actual values. The weights are updated iteratively using:

$$\Delta w_i = \eta (y_{\text{true}} - y_{\text{pred}}) x_i \tag{3}$$

Here, $\eta$ is the learning rate, which controls how quickly the model adjusts during training.

## 4.2 Support Vector Machine (SVM)

SVM is a powerful model for binary classification that works by finding the hyperplane that best separates the two classes. We used the RBF kernel to handle non-linear relationships in the data.

## 4.3 Logistic Regression

Logistic Regression models the probability of a binary outcome based on the features provided. It uses the log-odds transformation to produce a probability estimate:

$$P(y = 1|x) = \frac{1}{1 + e^{-(\mathbf{w}^T \mathbf{x} + b)}} \tag{4}$$

## 4.4 Random Forest

Random Forest is an ensemble learning method that constructs multiple decision trees and combines their predictions. This approach reduces the risk of overfitting and improves accuracy, especially when dealing with non-linear patterns in the data.

### 4.5 Hyperparameter Tuning and Feature Importance

To optimize each model, we used GridSearchCV to find the best hyperparameters. We also conducted feature importance analysis using Mutual Information to identify which features had the most influence on the predictions. Our analysis revealed that Glucose and BMI were the most significant factors in predicting diabetes.

### 4.6 Code Repository

The complete codebase for this study is available on GitHub. You can access the repository

https://github.com/varma1825/Predicting-Diabetes-Using-Perceptron.git

## 5 Experimental Analysis

### 5.1 Motivation for Testing

We conducted various tests to evaluate the SLP's performance and compare it to the more complex models:

1. **Benchmarking the SLP:** We wanted to see how the SLP performs relative to traditional models like SVM and Logistic Regression in predicting diabetes.

2. **Impact of Feature Scaling:** Since scaling affects models like SLP and SVM, we wanted to measure its impact on performance.

3. **Regularization Impact:** We applied L2 regularization to the SLP to see how it affected overfitting and overall performance.

### 5.2 Tests Conducted

1. **Baseline Performance:** We compared all four models in terms of accuracy, precision, recall, and F1-score.

2. **Feature Scaling Impact:** We tested how feature scaling improved the SLP's performance.

3. **Regularization Impact:** L2 regularization was applied to see if it improved the SLP's ability to generalize to new data.

| Model | Accuracy | Precision | Recall | F1-Score |
|---|---|---|---|---|
| SLP | **75.97%** | 0.82 | 0.81 | 0.81 |
| SVM | 73.38% | 0.77 | 0.83 | 0.80 |
| Logistic Regression | 75.32% | 0.81 | 0.80 | 0.81 |
| Random Forest | 72.08% | 0.79 | 0.78 | 0.78 |

Table 1: Performance Comparison of Different Models

### 5.3 Results

## 6 Discussion

The experiments revealed that the SLP performed well relative to more complex models, especially in terms of computational efficiency and simplicity. However, its reliance on linear decision boundaries was a limitation, particularly when the data was non-linearly separable. SVM, with its RBF kernel, and Random Forest handled non-linearity better, resulting in better recall for these models.

Regularization helped reduce overfitting for the SLP, improving its ability to generalize to unseen data. Additionally, feature scaling proved essential for both the SLP and SVM, as unscaled features led to poorer performance.

## 7 Conclusion

This study has shown that the Single Layer Perceptron is a viable model for predicting diabetes, particularly when computational efficiency is important. While the SLP performed comparably to Logistic Regression and SVM, it struggled with non-linear relationships in the data. Regularization and feature scaling significantly improved its performance, but future research could explore using Multi-Layer Perceptrons (MLPs) to handle more complex data patterns.

### 7.1 Future Work

- **Multi-Layer Perceptrons (MLP):** Using MLPs could improve the SLP's ability to handle non-linear data.

- **Regularization Improvements:** Exploring more

advanced regularization techniques like dropout or L1 regularization could enhance the SLP's generalization capabilities.

- **Ensemble Methods:** Incorporating ensemble methods like bagging or boosting could further improve model performance and stability.

# References

[1] Pima Indians Diabetes Dataset. Available at: https://www.kaggle.com/uciml/pima-indians-diabetes-database

[2] J. Brownlee, *Deep Learning with Python*, 2nd ed. O'Reilly Media, 2021.

[3] D. A. Freedman, *Statistical Models: Theory and Practice*. Cambridge University Press, 2005.

[4] G. E. Hinton, *Neural Networks for Machine Learning*, University of Toronto, Lecture Series, 2012.

[5] F. Pedregosa, et al., "Scikit-learn: Machine Learning in Python," *Journal of Machine Learning Research*, vol. 12, pp. 2825-2830, 2011.

[6] A. Géron, *Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow*, O'Reilly Media, 2019.

[7] K. P. Murphy, *Machine Learning: A Probabilistic Perspective*. MIT Press, 2012.