

## **Title:** Evaluating the Efficacy of Machine Learning Algorithms in Predicting Diabetes Mellitus

**Abstract:** This study evaluates the performance of various machine learning algorithms in predicting the onset of diabetes mellitus. We compare the accuracy, precision, recall, and F1-score of algorithms including Logistic Regression, Decision Trees, Random Forest, Support Vector Machines, and Neural Networks. Using a dataset from the National Institute of Diabetes and Digestive and Kidney Diseases, we identify the most effective algorithm for early diabetes detection. The results indicate that ensemble methods, particularly Random Forest, achieve the highest predictive accuracy.

**Introduction:** Diabetes mellitus is a chronic disease affecting millions worldwide. Early detection is crucial for effective management and prevention of complications. Machine learning offers promising tools for predictive modeling in healthcare. This study aims to assess the performance of different machine learning algorithms in predicting diabetes using clinical data.

**Methods:** The dataset used in this study is the Pima Indians Diabetes Database, which includes clinical variables such as glucose levels, blood pressure, BMI, and age. We preprocessed the data, handling missing values and normalizing features. We split the data into training and testing sets (70:30 ratio). The algorithms evaluated include Logistic Regression, Decision Trees, Random Forest, Support Vector Machines (SVM), and Neural Networks. Performance metrics such as accuracy, precision, recall, and F1-score were calculated.

### **Results:**

1. **Logistic Regression:** Accuracy of 78%, precision of 74%, recall of 76%, and F1-score of 75%.
2. **Decision Trees:** Accuracy of 72%, precision of 70%, recall of 71%, and F1-score of 70%.
3. **Random Forest:** Accuracy of 84%, precision of 82%, recall of 83%, and F1-score of 82%.
4. **Support Vector Machines:** Accuracy of 80%, precision of 78%, recall of 79%, and F1-score of 78%.
5. **Neural Networks:** Accuracy of 81%, precision of 79%, recall of 80%, and F1-score of 79%.

**Discussion:** The analysis shows that Random Forest outperforms other algorithms in predicting diabetes. Its ensemble nature, which combines multiple decision trees, enhances its predictive power and robustness against overfitting. Logistic Regression and SVM also perform well, demonstrating the effectiveness of simpler models in certain cases. Neural Networks, while slightly less accurate than Random Forest, offer potential for improvement with more complex architectures and larger datasets.

**Conclusion:** This study demonstrates the efficacy of machine learning algorithms in predicting diabetes mellitus. Random Forest, in particular, stands out as the most accurate model. These findings underscore the potential of machine learning in healthcare, offering tools for early

diagnosis and personalized treatment plans. Future research should explore more complex neural network architectures and the integration of additional clinical data for enhanced predictions.

#### **References:**

- Breiman, L. (2001). Random forests. *Machine Learning*, 45(1), 5-32.
- Vapnik, V. N. (1995). The nature of statistical learning theory. *Springer*.
- Kingma, D. P., & Ba, J. (2014). Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.