

Predicting the Onset of Diabetes Based on Diagnostic Data

1. Executive Summary

Diabetes is a rapidly growing global health concern, with India projected to become one of the most affected countries in the coming decades. Early identification of individuals at risk plays a crucial role in preventing complications and improving quality of life.

This project applies supervised machine learning techniques to predict the onset of diabetes using routine diagnostic and physiological data. By performing exploratory data analysis, baseline model evaluation, and feature engineering, the project demonstrates how data-driven approaches can support early diagnosis and informed healthcare decisions.

2. Problem Statement

According to medical studies, diabetes affects millions of individuals worldwide and often remains undiagnosed until complications arise. Traditional diagnosis methods rely heavily on clinical judgment, which can be enhanced using predictive analytics.

The objective of this project is to build a machine learning classification model that predicts whether a patient is likely to develop diabetes based on diagnostic attributes such as glucose levels, BMI, insulin levels, and age.

3. Dataset Description

3.1 Data Source

- Pima Indians Diabetes Database

3.2 Dataset Type

- Structured medical diagnostic data

3.3 Target Variable

- Outcome

- 1 → Diabetic
- 0 → Non-Diabetic

3.4 Features

- Glucose
 - BloodPressure
 - SkinThickness
 - Insulin
 - BMI
 - DiabetesPedigreeFunction
 - Age
-

4. Methodology Overview

The project follows a structured machine learning pipeline:

1. Data Understanding and Exploratory Data Analysis (EDA)
2. Baseline Model Development
3. Feature Engineering
4. Model Evaluation and Performance Comparison
5. Result Interpretation

This approach ensures transparency, reproducibility, and continuous performance improvement.

5. Exploratory Data Analysis (EDA)

Exploratory Data Analysis was performed to understand the characteristics and quality of the dataset.

5.1 Data Inspection

- Examined dataset shape, data types, and summary statistics
- Identified zero values in features such as Glucose, Insulin, BMI, and BloodPressure

5.2 Univariate Analysis

- Studied the distribution of each feature
- Observed skewness in variables such as Glucose, Age and BMI

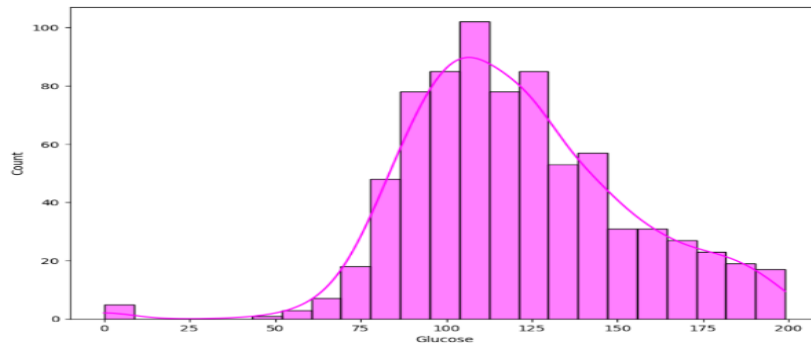


Fig1: Histogram of **Glucose**

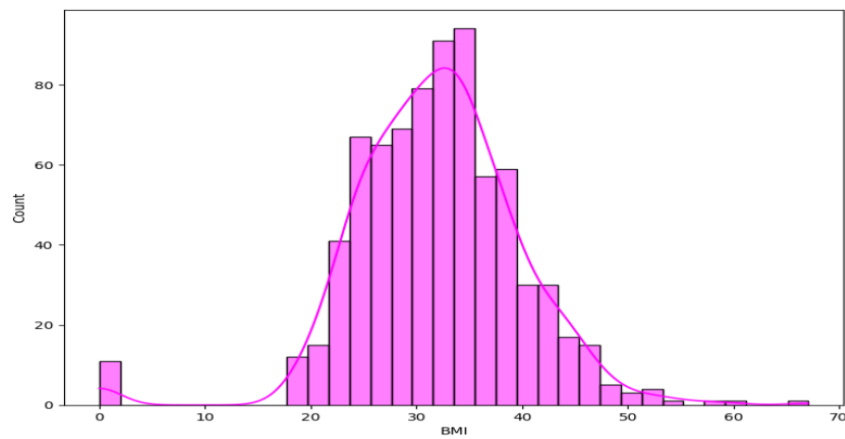


Fig2: Histogram of BMI

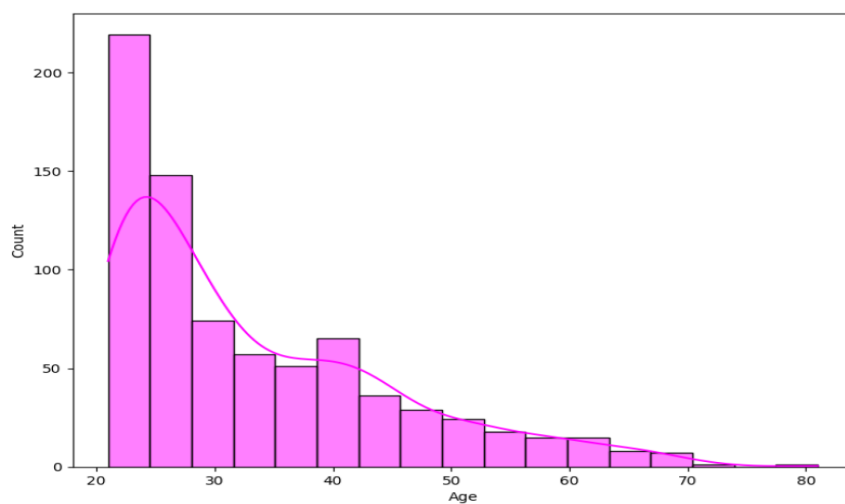


Fig3: Histogram of Age

5.3 Multivariate Analysis

- Analyzed correlations between features and the target variable

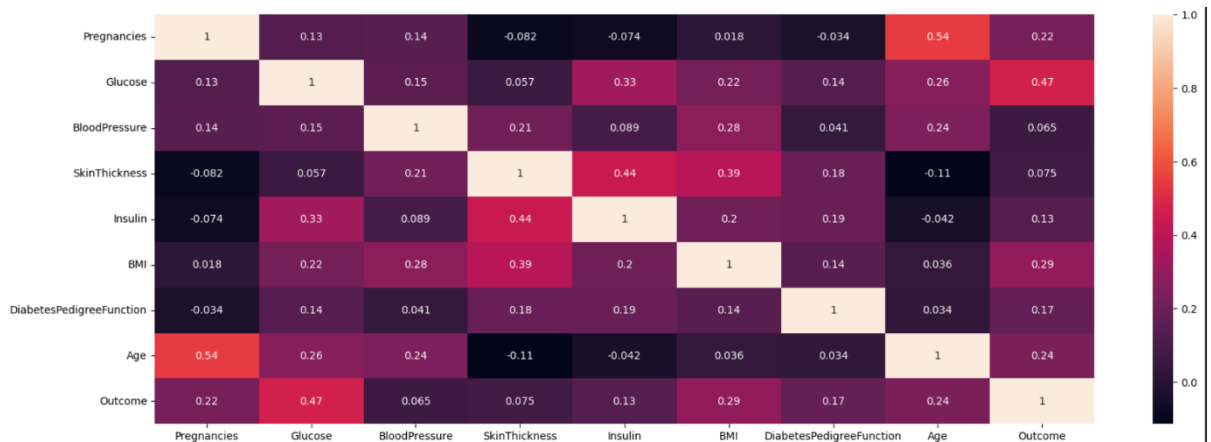


Fig 4: Correlation Heatmap

- Glucose showed the strongest correlation with diabetes outcome

5.4 Key Observations

- Presence of zero values that are not medically meaningful
- Existence of outliers in Insulin, BMI, and SkinThickness
- Certain features carry significantly more predictive power

EDA provided critical insights that guided preprocessing and feature engineering decisions.

6. Baseline Model Development

To establish benchmark performance, multiple classification algorithms were implemented using the raw dataset.

6.1 Train-Test Split

- Dataset split into training and testing sets (80:20)

6.2 Models Implemented

- Logistic Regression
- K-Nearest Neighbors (KNN)
- Decision Tree Classifier
- Random Forest Classifier

6.3 Evaluation Metric

- Accuracy score

6.4 Baseline Results

Baseline models produced reasonable accuracy but highlighted the need for improved preprocessing and feature optimization. These results served as reference points for further enhancement.

7. Feature Engineering

Feature engineering was applied to improve model performance and generalization.

7.1 Feature Scaling

Two scaling techniques were used:

- MinMaxScaler
- StandardScaler

Scaling improved convergence for distance-based and linear models.

7.2 Feature Selection

Based on medical relevance and correlation analysis, the following features were selected:

- Glucose
- BMI
- Age
- Insulin
- DiabetesPedigreeFunction

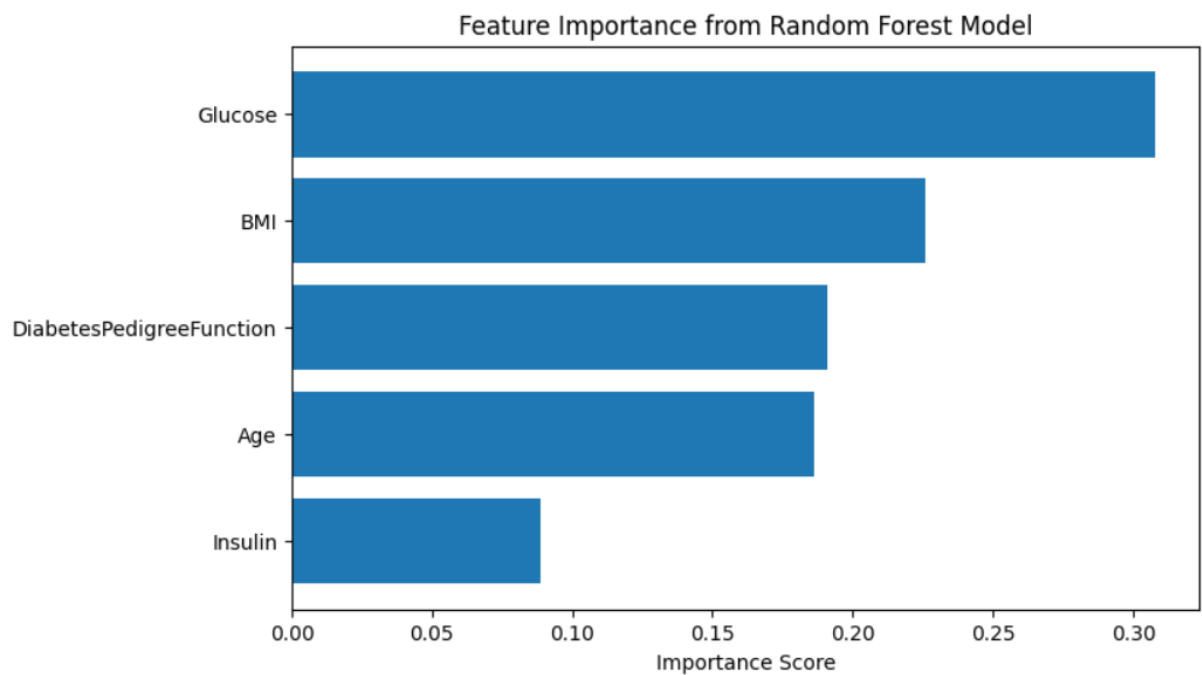


Fig 5 : Feature importance scores obtained from the Random Forest model.

Rationale:

- Glucose directly reflects blood sugar levels
- BMI indicates obesity, a major diabetes risk factor
- Age correlates with increased diabetes risk
- Insulin reflects insulin resistance
- DiabetesPedigreeFunction captures genetic influence

7.3 Outlier Removal

Outliers were removed using the Interquartile Range (IQR) method. This step reduced noise and improved model robustness.

7.4 Improved Model Training

After feature scaling, selection, and outlier removal, models were retrained. Random Forest performed best due to its ability to handle non-linear relationships and feature interactions.

8. Results and Performance Analysis

- Feature engineering significantly improved prediction stability

Model	Accuracy
Logistic Regression	0.74675
KNN	0.66233
Decision Tree	0.74675
Random Forest	0.72077

- Random Forest achieved the highest accuracy among tested models
- Scaled and cleaned data produced more reliable predictions

These results confirm the importance of preprocessing and thoughtful feature selection in medical machine learning applications.

9. Key Insights

- Glucose is the most influential predictor of diabetes
 - Data preprocessing has a substantial impact on model performance
 - Ensemble models outperform simpler classifiers in this problem
 - Feature engineering is essential for real-world healthcare datasets
-

10. Limitations

- Dataset size is relatively small
 - Zero values may not fully represent missing clinical data
 - Class imbalance was not explicitly addressed at this stage
-

11. Future Enhancements

- Hyperparameter tuning using GridSearchCV

- Cross-validation for improved generalization
 - Handling class imbalance using SMOTE
 - Feature importance visualization
 - Model deployment using a web interface
-

12. Conclusion

This project demonstrates how supervised machine learning can be applied to predict diabetes onset using diagnostic data. Through systematic analysis, feature engineering, and model evaluation, the study highlights the role of data-driven techniques in supporting early diagnosis and preventive healthcare.

With further tuning and deployment, such models can assist clinicians and healthcare systems in improving patient outcomes.

13. Tools and Technologies Used

- Python
- Pandas, NumPy
- Matplotlib, Seaborn
- Scikit-learn

