# DIABETES PREDICTION USING MACHINE LEARNING
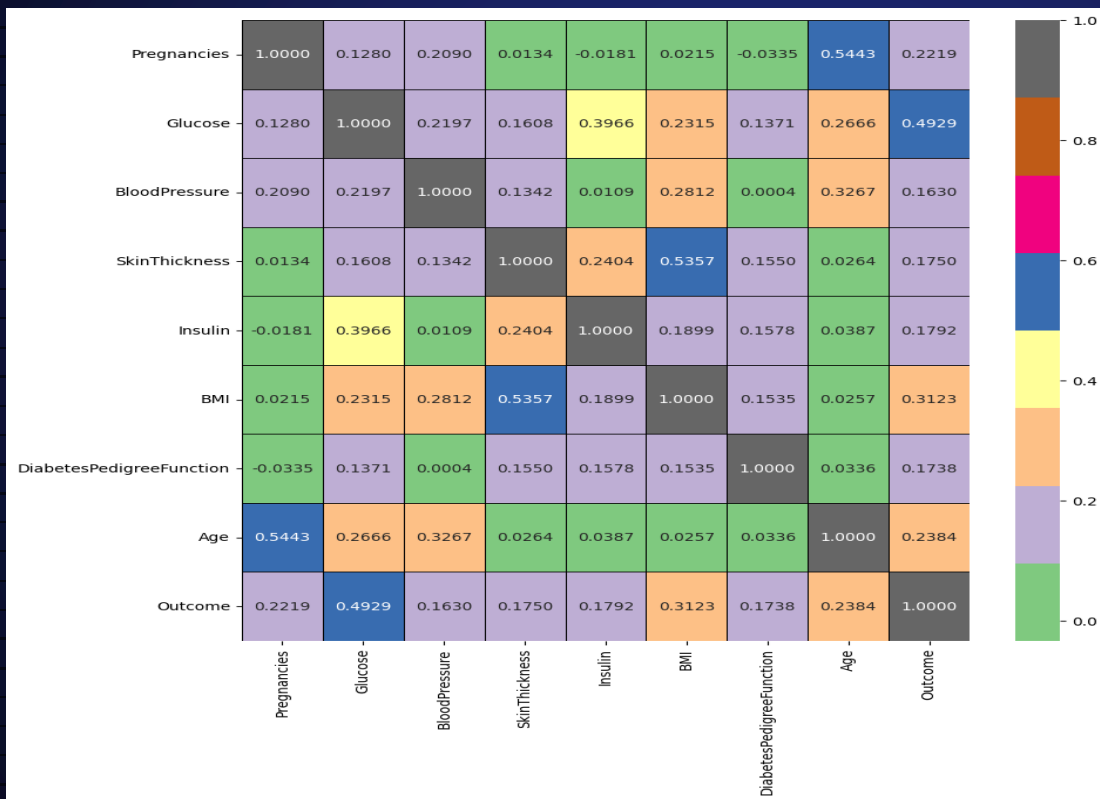
Kokoh Dwiko Listanto

# PROBLEM RESEARCH

The research objective is to develop a **predictive model** for diabetes diagnosis. This endeavor is motivated by the potential to **improve healthcare** by enabling early intervention and targeted care, contributing to public health insights, aiding clinical decision-making, and promoting preventative measures. Additionally, it provides an opportunity to **apply machine learning techniques** to a relevant healthcare challenge, ultimately benefiting both patients and the medical community.
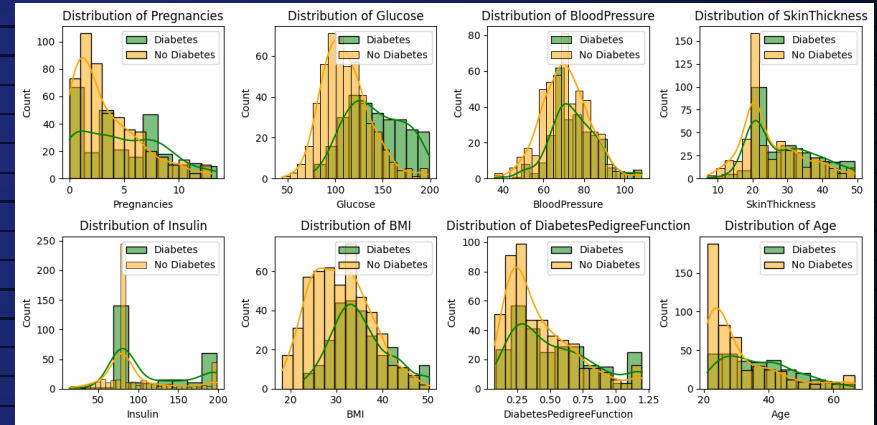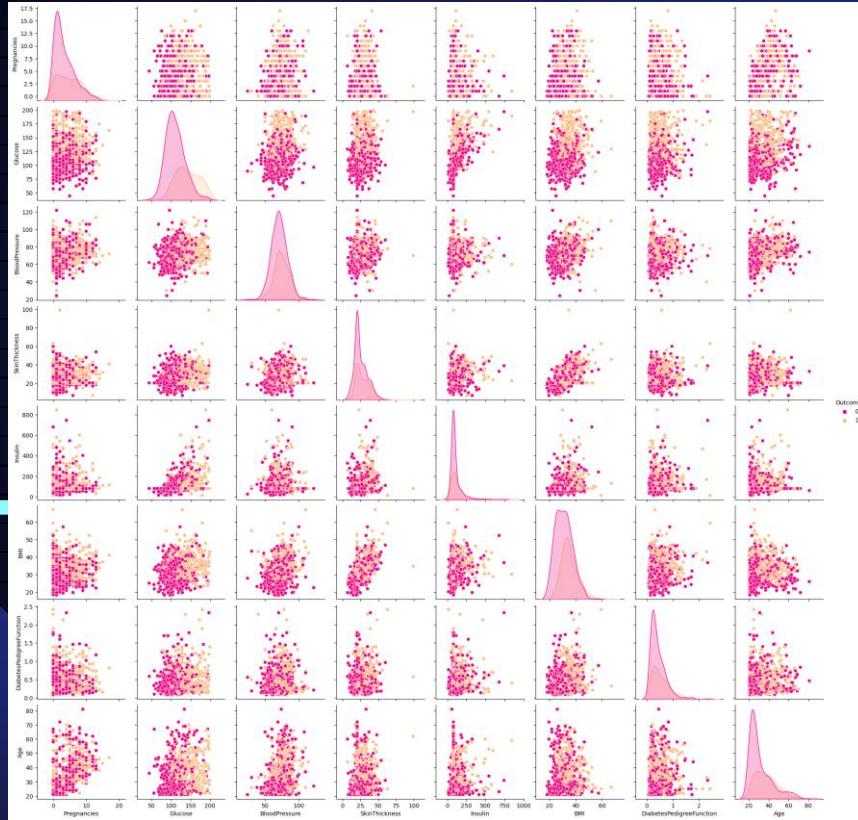
# DATASET INFORMATION

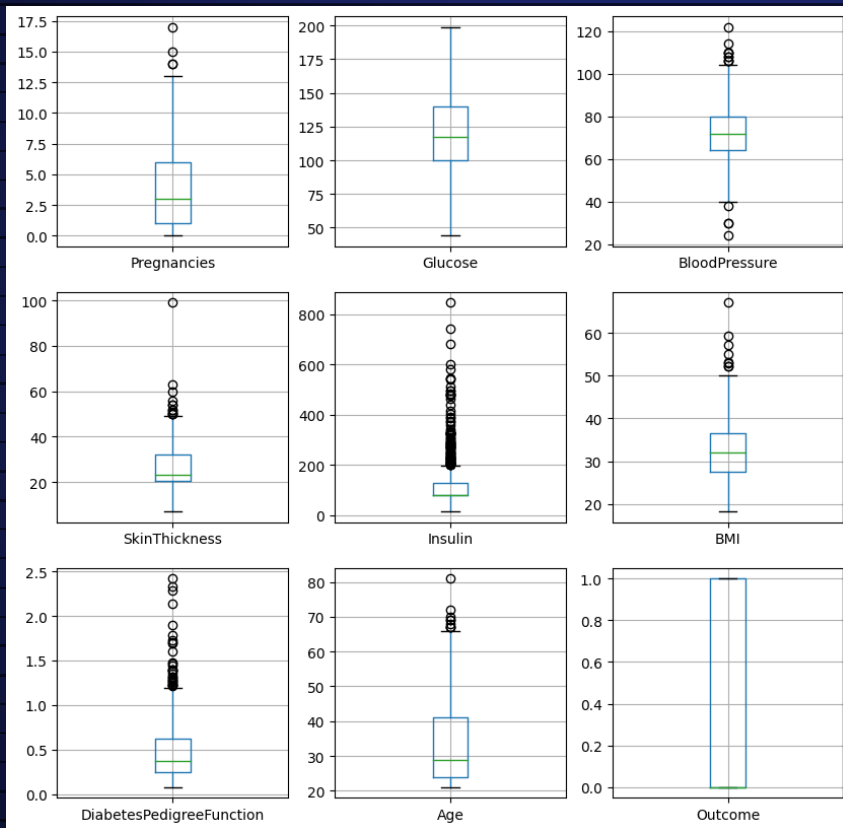| | |
|---|---|
| **Pregnancies:** | To express the Number of pregnancies |
| **Glucose:** | To express the Glucose level in blood |
| **BloodPressure:** | To express the Blood pressure measurement |
| **SkinThickness:** | To express the thickness of the skin |
| **Insulin:** | To express the Insulin level in blood |
| **BMI:** | To express the Body mass index |
| **DiabetesPedigreeFunction:** | To express the Diabetes percentage |
| **Age:** | To express the age |
| **Outcome:** | To express the final result 1 is Yes and 0 is No |

# Data Visualization



There are **no features that have a high positive correlation** with each other.
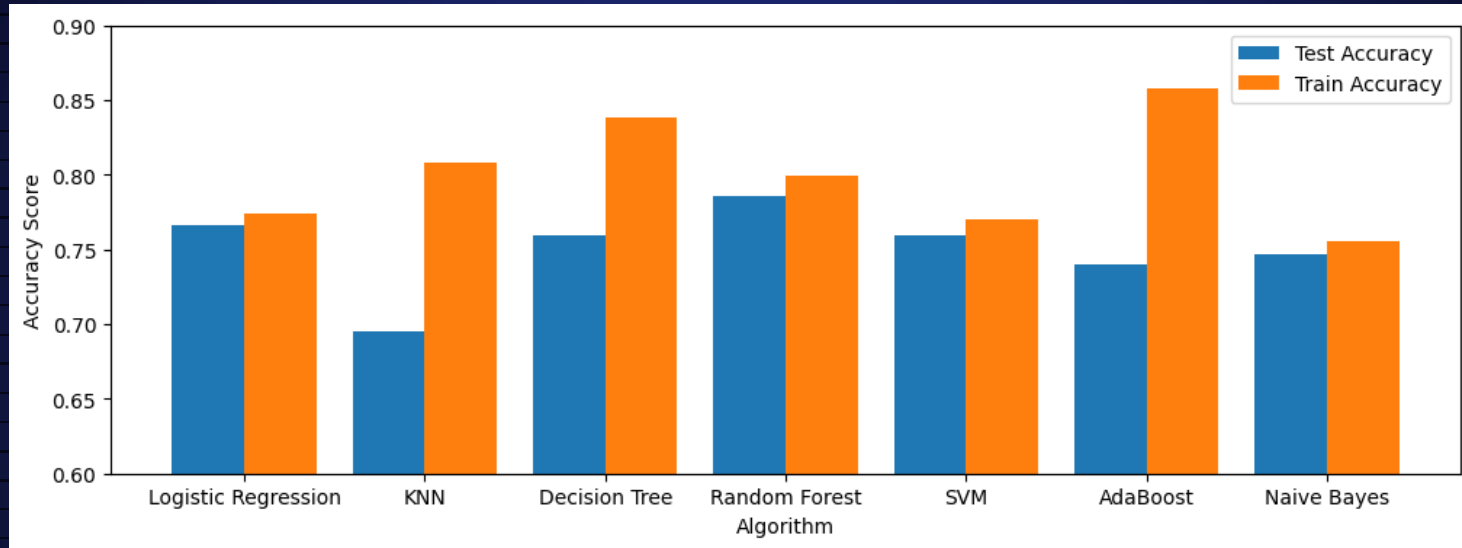
# Data Visualization

# Data Visualization



Some data has many outliers, so these **outliers must be removed**

# Machine Learning Output



Among the given machine learning algorithms, the **Random Forest algorithm** appears to perform the best with a test accuracy of 0.7857 and a relatively close train accuracy of 0.7997. Random Forest combines the strength of multiple decision trees and tends to be robust and accurate in a variety of scenarios, which **is reflected in its performance on the test data**. It strikes a **good balance** between generalization and overfitting, making it a strong candidate for the model of choice.