

# Supervised Learning Project

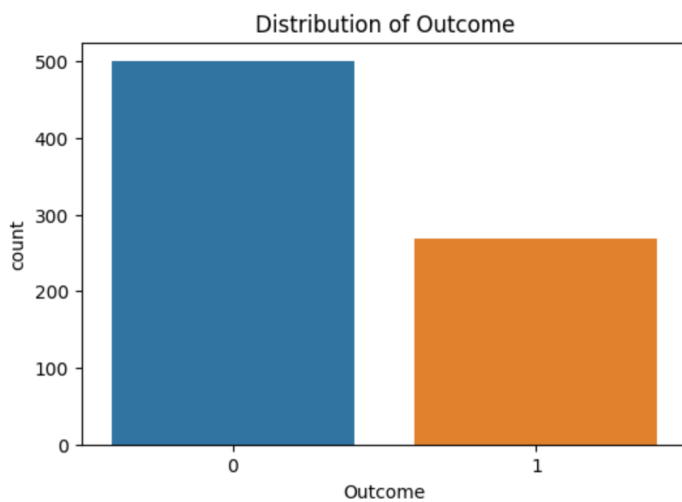
## — based on the “Diabetes” dataset

### Part1: Exploratory Data Analysis

1. Get some basic information about the dataset, including shape, columns, data types, null values.
2. Visualize the relationships between the different variables.

### Visualization

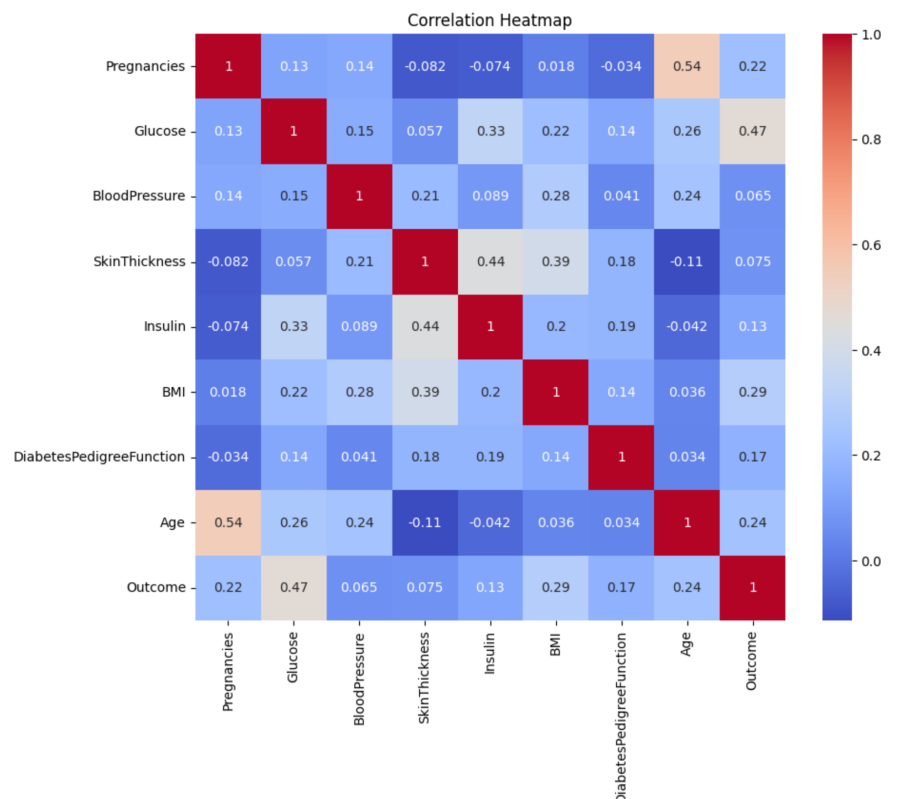
#### 2.1 Outcome Distribution



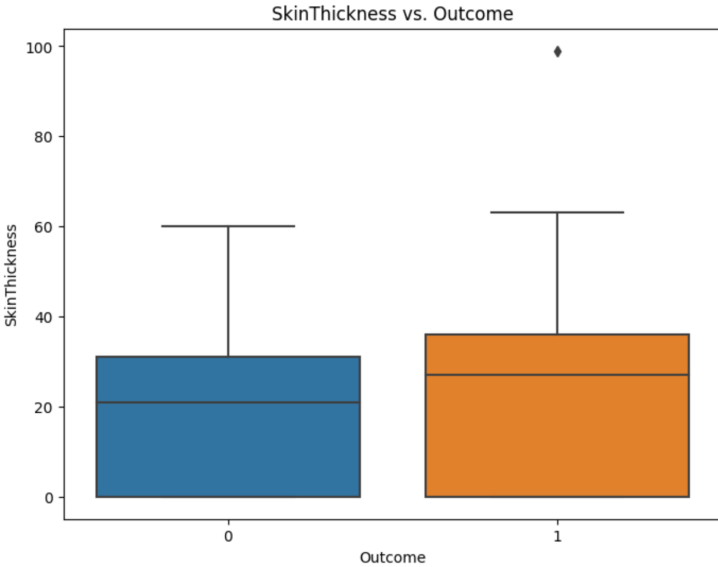
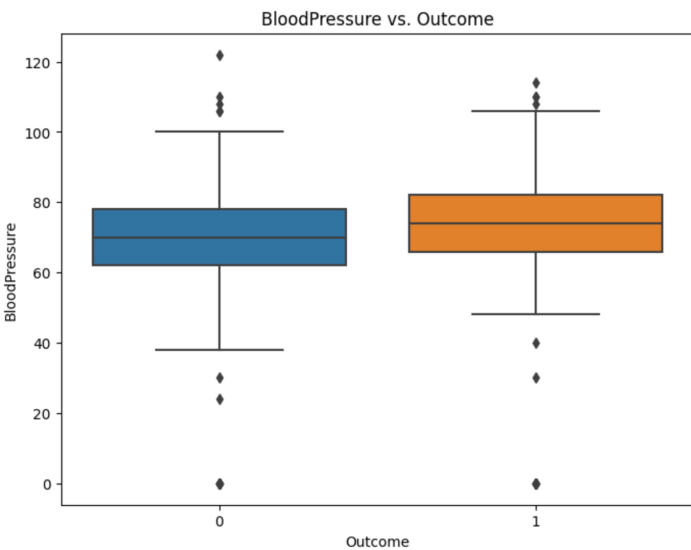
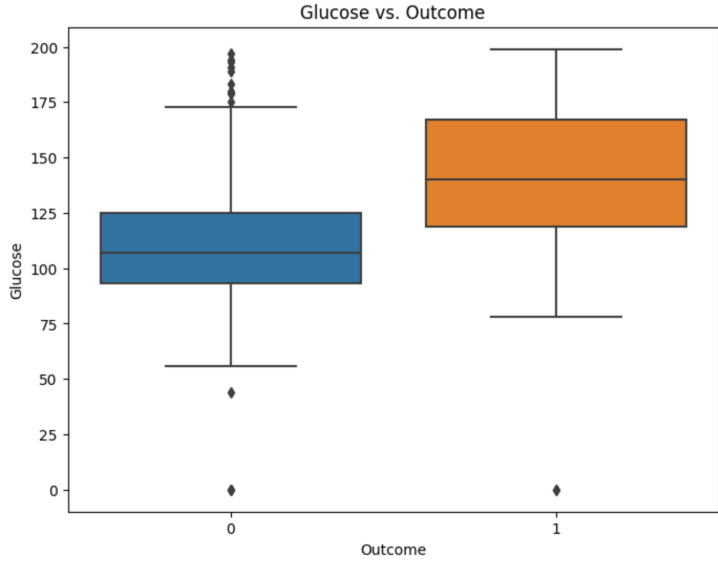
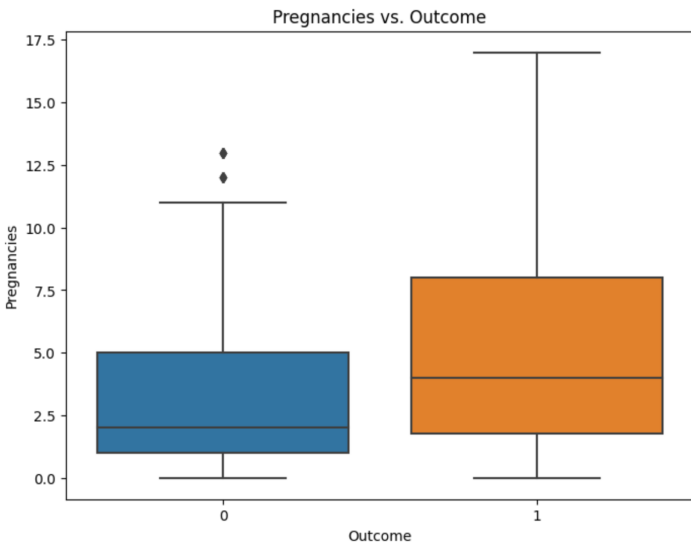
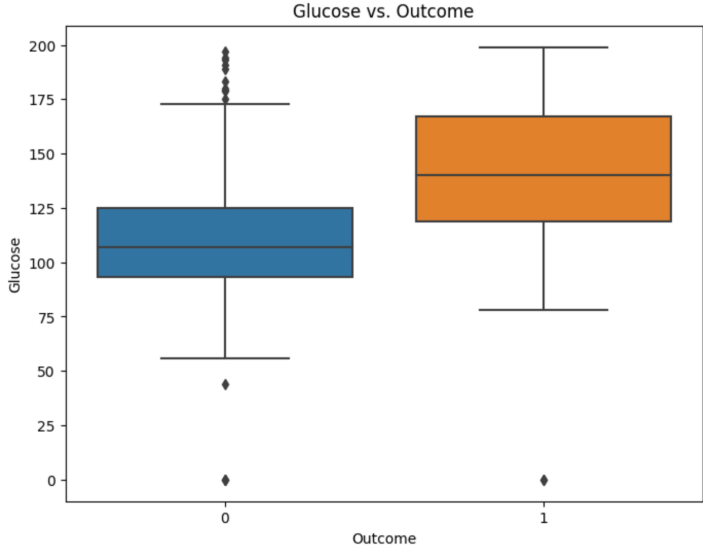
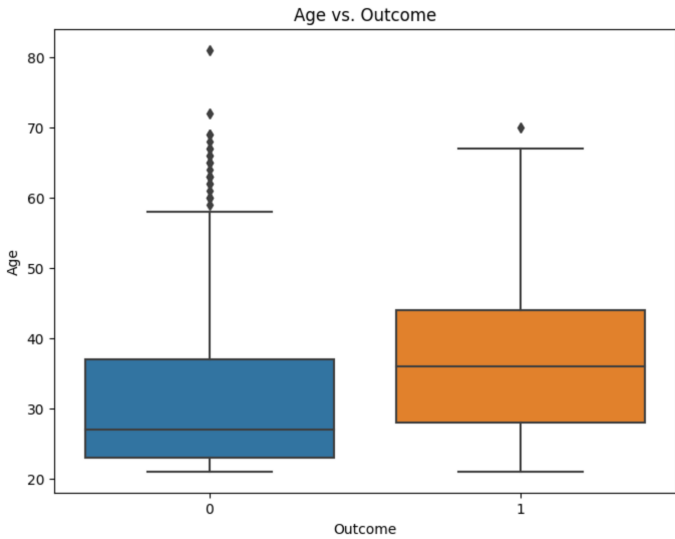
**Attention:** imbalanced data

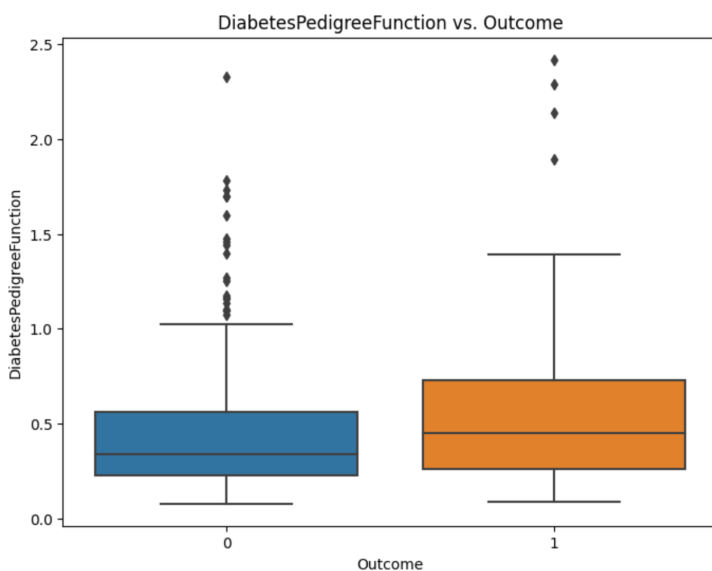
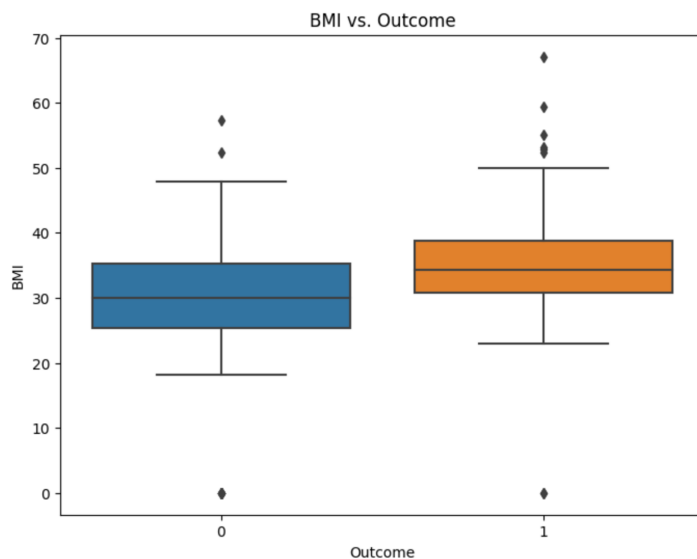
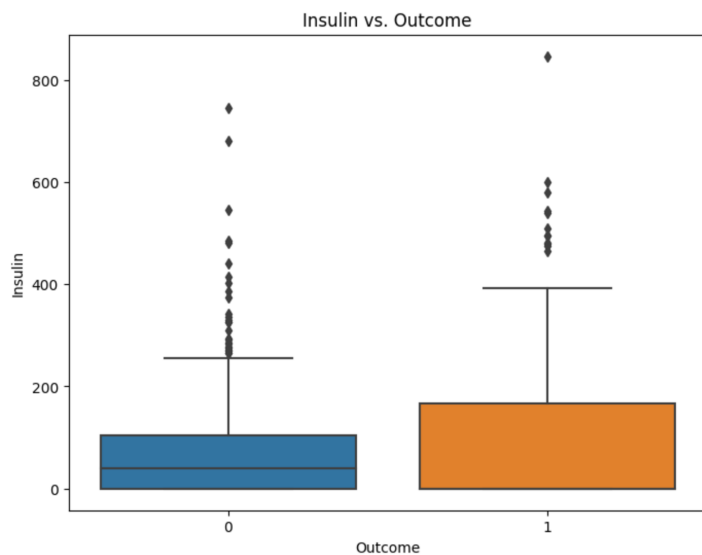
#### 2.2 Correlation Heatmap

Glucose - Outcome: 0.47  
BMI - Outcome: 0.29  
Pregancies - Outcome: 0.22  
Age - Outcome: 0.24



## 2.3 Boxplot for outliers





## Part2: Preprocessing & Feature Engineering

1. Handling missing values  
No missing values in this dataset.
2. Handling outliers  
Using IQR-based approach.
3. Feature engineering  
Create a new feature named 'BMI\_Category' with labels: Underweight, Normal, Overweight, Obese and applied label encoding.
4. Handling imbalanced data  
Apply random undersampling to balance the classes.

## Part3: Train ML Model

Select **Logistic Regression** and **Random Forest** models.

### 3.1 Train and test model with all features

	<b>Logistic Regression</b>	<b>Random Forest</b>
<b>accuracy</b>	0.6962025316455697	0.7341772151898734
<b>precision</b>	0.5476190476190477	0.6060606060606061
<b>recall</b>	0.8214285714285714	0.7142857142857143
<b>F-1 score</b>	0.6571428571428571	0.6557377049180327
<b>roc-auc</b>	0.7244397759103641	0.7296918767507004

**Random Forest WINS!!!**

### 3.2 Remove BMI\_Category feature

	<b>Logistic Regression</b>	<b>Random Forest</b>
<b>accuracy</b>	0.7341772151898734	0.7468354430379747
<b>precision</b>	0.5945945945945946	0.6111111111111112
<b>recall</b>	0.7857142857142857	0.7857142857142857
<b>F-1 score</b>	0.676923076923077	0.6875000000000001
<b>roc-auc</b>	0.745798319327731	0.7556022408963585

**Random Forest WINS!!!**

### 3.3 KFold cross-validation

	<b>Logistic Regression</b>	<b>Random Forest</b>
<b>Fold 1</b>	0.7722	0.7468
<b>Fold 2</b>	0.6835	0.7468
<b>Fold 3</b>	0.7722	0.7468
<b>Fold 4</b>	0.7468	0.7468
<b>Fold 5</b>	0.7564	0.7564
<b>Mean Accuracy</b>	0.7462	0.7488

**Random Forest WINS!!!**

## **Part4: Conclusion**

From the machine learning models developed and the exploratory data analysis (EDA) conducted, there are my findings:

1. Based on the correlation heatmap, Glucose is the most significant predictor of diabetes outcome. Also, age, BMI and pregnancy play important roles.
2. Proper preprocessing steps, including feature scaling, one-hot-encoding significantly improved the model's performance in this case.
3. The accuracy increased for both models after removing the new generated feature - BMI\_Category\_Encoded.
4. Logistic Regression and Random Forest were developed as predictive models for diabetes outcome. Random Forest has a better performance on the test dataset.
5. The dataset shows an imbalanced distribution, with a higher number of non-diabetic cases compared to diabetic cases.