



SOUMYAJITA BOSE

IITP000180

2312RES655

DATA SCIENCE

soumyajita_2312res655@iitp.ac.in

<https://github.com/23f1000193/Diabetes-Risk-Prediction-from-Health-Records.git>

Diabetes Risk Prediction from Health Records

This project aims to predict the likelihood of diabetes in patients using health indicators such as **Glucose, BMI, Age, Blood Pressure, Insulin, and Skin Thickness**. By applying **data science techniques and machine learning models**, this project demonstrates how early detection can help in prevention and effective healthcare planning.

Dataset Used: **PIMA Indian Diabetes Dataset**

Goals & Motivation

- Build a **predictive model** to identify individuals at risk of diabetes.
- Explore and analyze medical data through **data visualization**.
- Apply **machine learning algorithms** to achieve reliable predictions.
- Provide **insights** into key health factors contributing to diabetes.
- Showcase an **end-to-end data science pipeline** from raw data to results.

- **Models Implemented:**
 - Logistic Regression
 - Random Forest Classifier
 - Gradient Boosting (XGBoost)
 - Support Vector Machine (SVM)
- **Train-Test Split: 80:20**
- **Cross-validation** used to ensure reliability

```

✓ Dataset loaded. Shape: (768, 9)
Pregnancies  Glucose  BloodPressure  SkinThickness  Insulin  BMI  \
0           6      148           72           35         0  33.6
1           1       85           66           29         0  26.6
2           8      183           64            0         0  23.3
3           1       89           66           23        94  28.1
4           0      137           40           35       168  43.1

DiabetesPedigreeFunction  Age  Outcome
0           0.627      50         1
1           0.351      31         0
2           0.672      32         1
3           0.167      21         0
4           2.288      33         1

```

Cross-validation results:

	Model	CV ROC-AUC
0	LogisticRegression	0.849029
3	SVC	0.831892
1	RandomForest	0.830144
2	GradientBoosting	0.817014

✓ Best model: LogisticRegression

Test set evaluation:

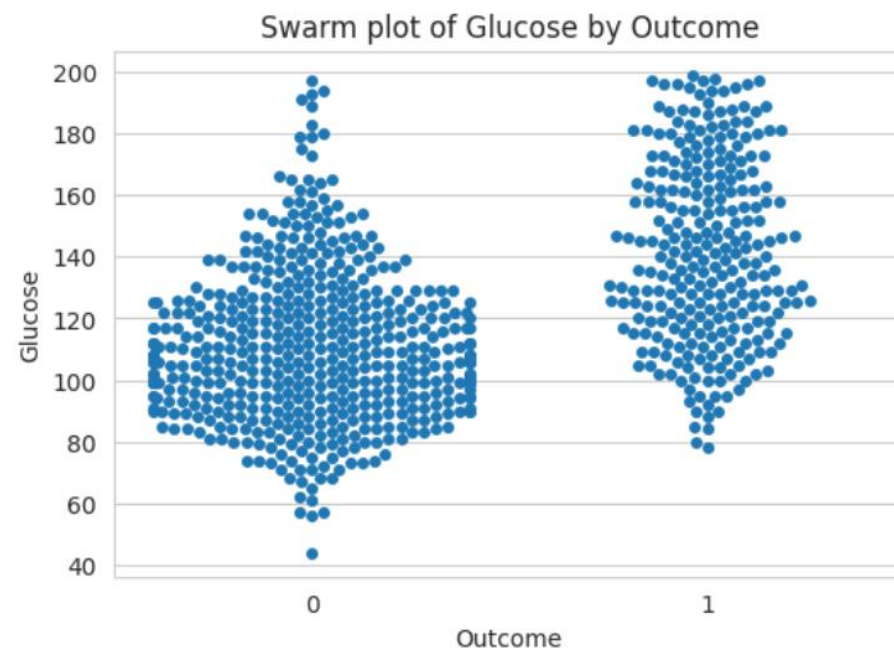
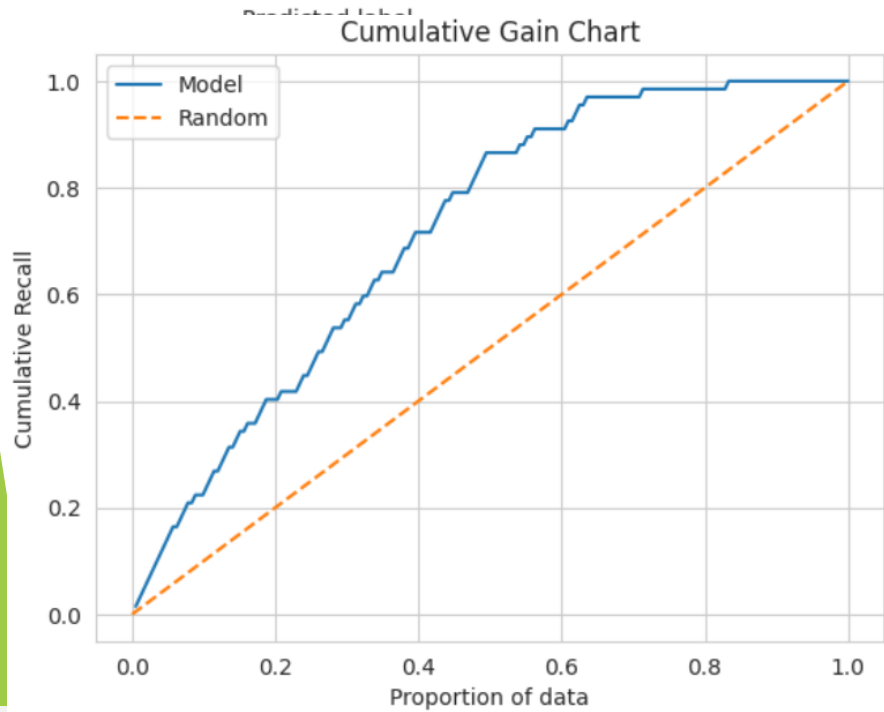
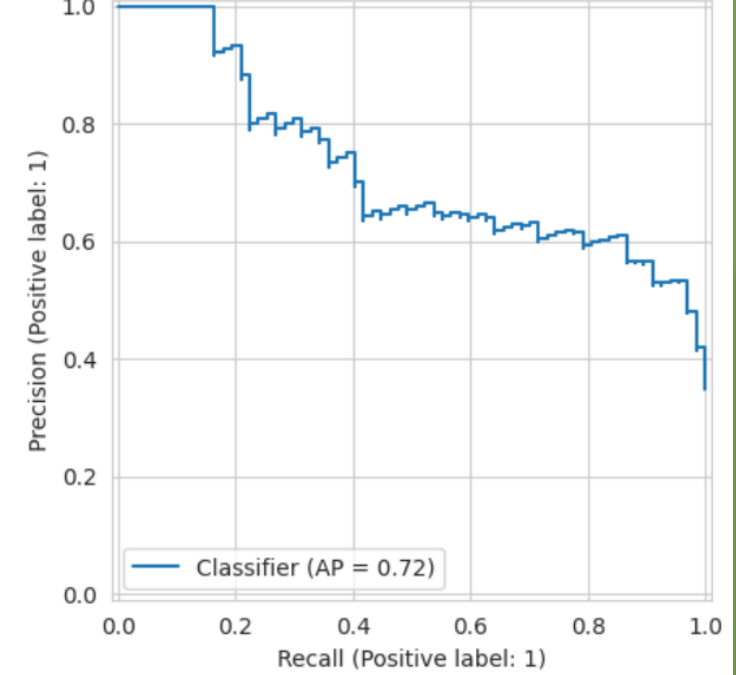
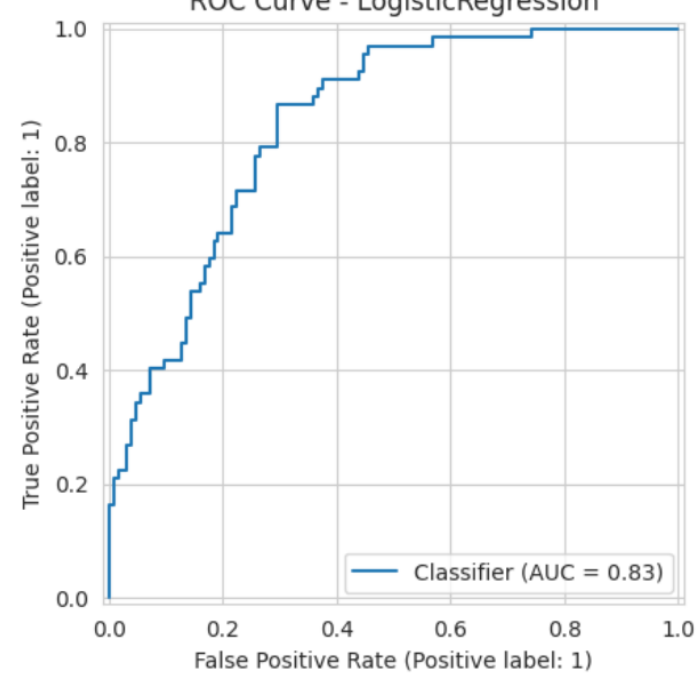
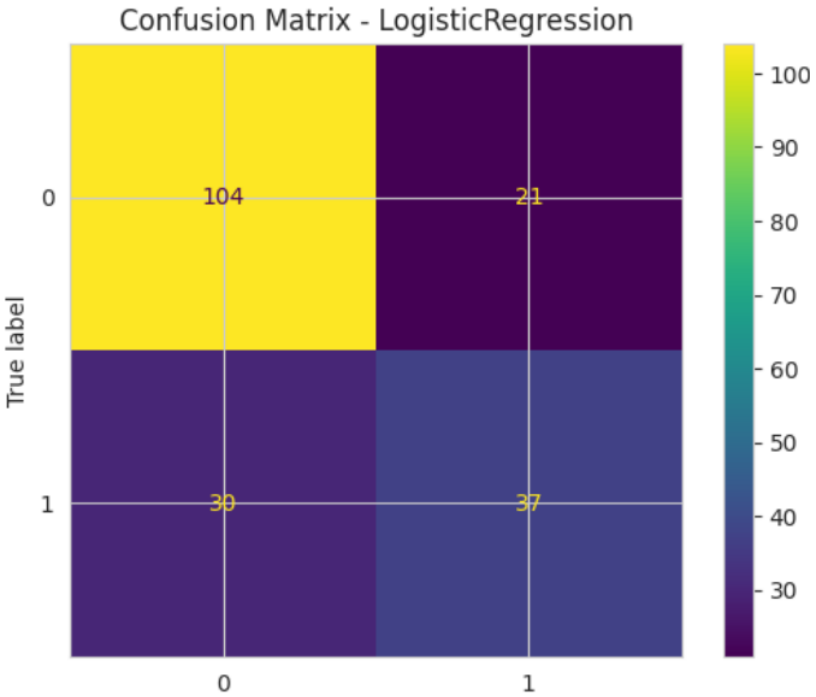
	Accuracy	Precision	Recall	F1	ROC_AUC
0	0.734375	0.637931	0.552239	0.592	0.83403

Classification Report:

	precision	recall	f1-score	support
0	0.78	0.83	0.80	125
1	0.64	0.55	0.59	67
accuracy			0.73	192
macro avg	0.71	0.69	0.70	192
weighted avg	0.73	0.73	0.73	192

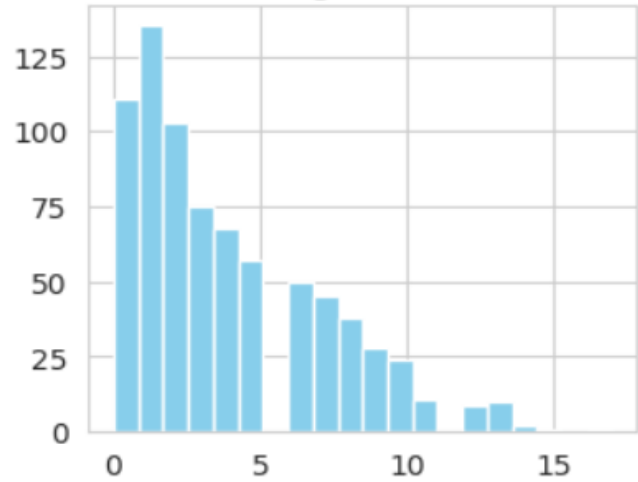
Summary statistics:

	count	mean	std	min	25%	\
Pregnancies	768.0	3.845052	3.369578	0.000	1.00000	
Glucose	768.0	121.656250	30.438286	44.000	99.75000	
BloodPressure	768.0	72.386719	12.096642	24.000	64.00000	
SkinThickness	768.0	27.334635	9.229014	7.000	23.00000	
Insulin	768.0	94.652344	105.547598	14.000	30.50000	
BMI	768.0	32.450911	6.875366	18.200	27.50000	
DiabetesPedigreeFunction	768.0	0.471876	0.331329	0.078	0.24375	
Age	768.0	33.240885	11.760232	21.000	24.00000	
Outcome	768.0	0.348958	0.476951	0.000	0.00000	
		50%	75%	max		
Pregnancies	3.0000	6.00000	17.00			
Glucose	117.0000	140.25000	199.00			
BloodPressure	72.0000	80.00000	122.00			
SkinThickness	23.0000	32.00000	99.00			
Insulin	31.2500	127.25000	846.00			
BMI	32.0000	36.60000	67.10			
DiabetesPedigreeFunction	0.3725	0.62625	2.42			
Age	29.0000	41.00000	81.00			
Outcome	0.0000	1.00000	1.00			

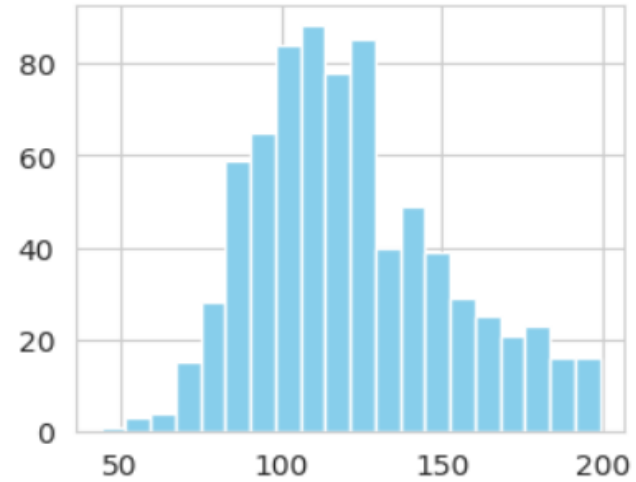


Histograms of Features

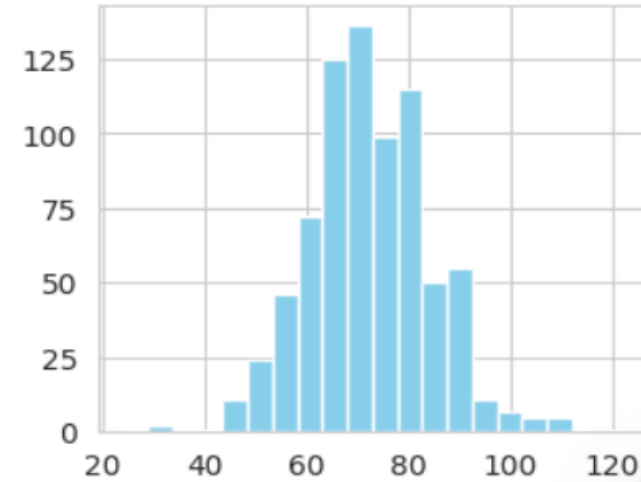
Pregnancies



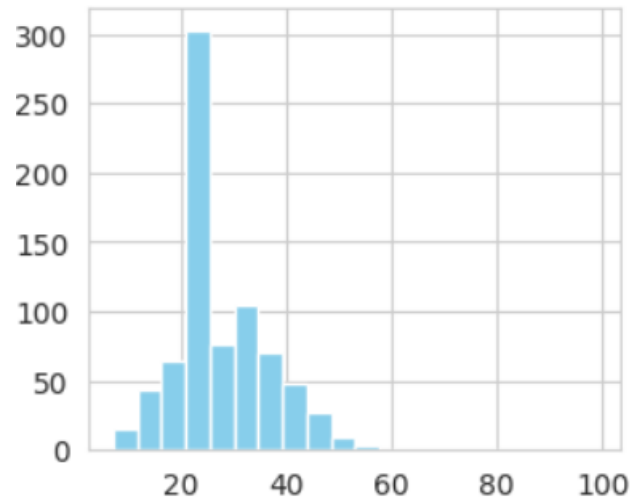
Glucose



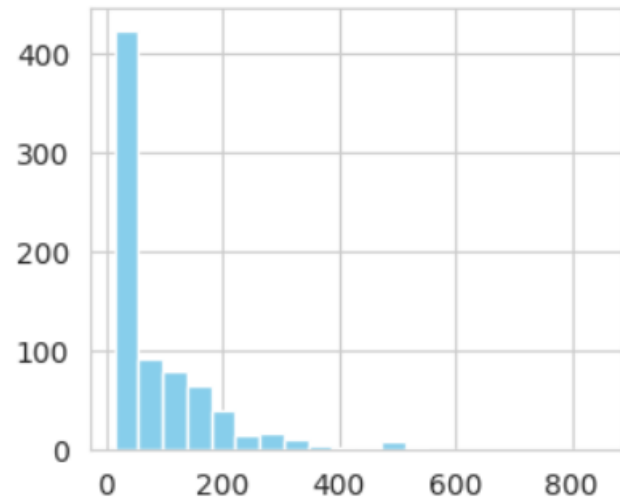
BloodPressure



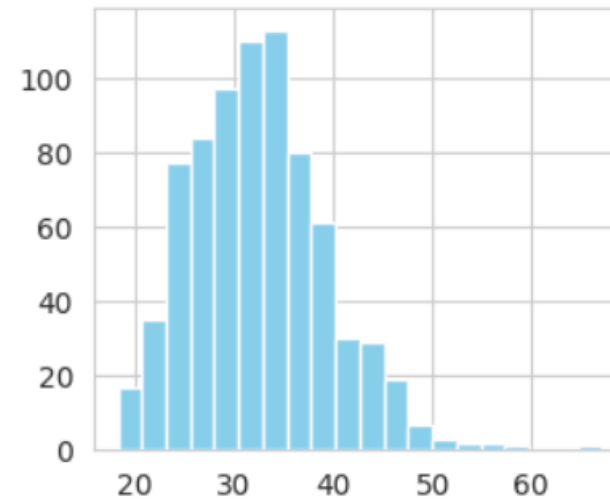
SkinThickness

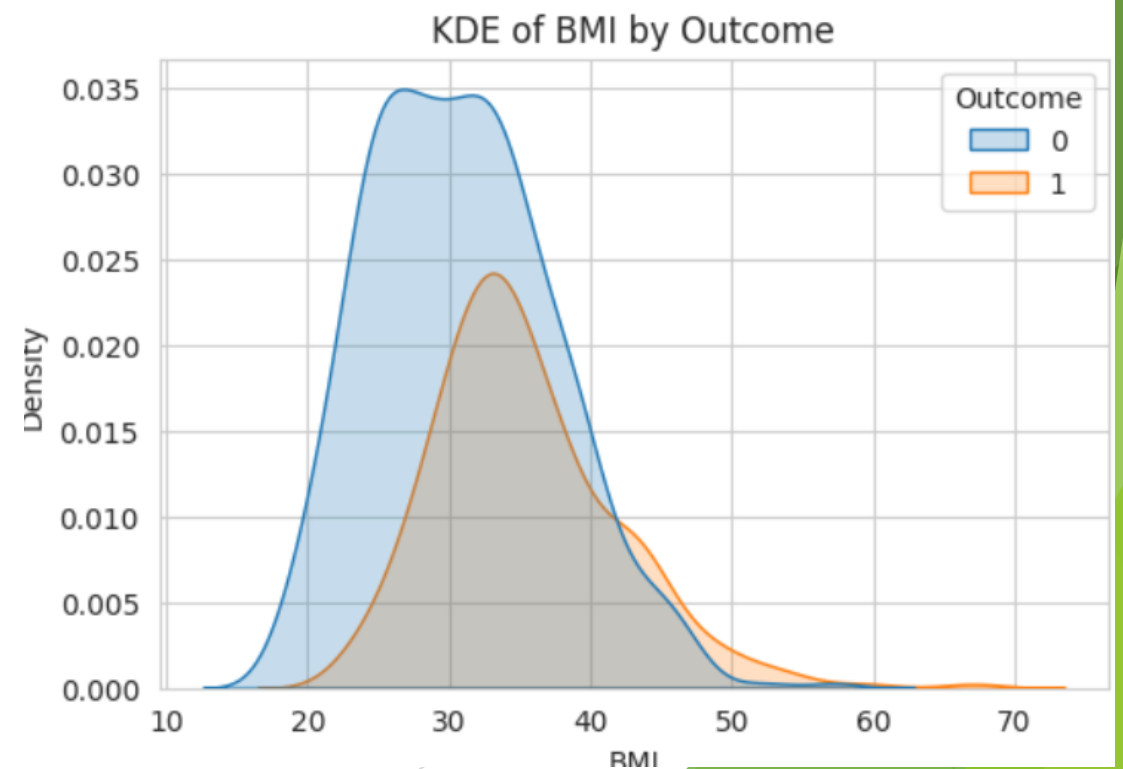
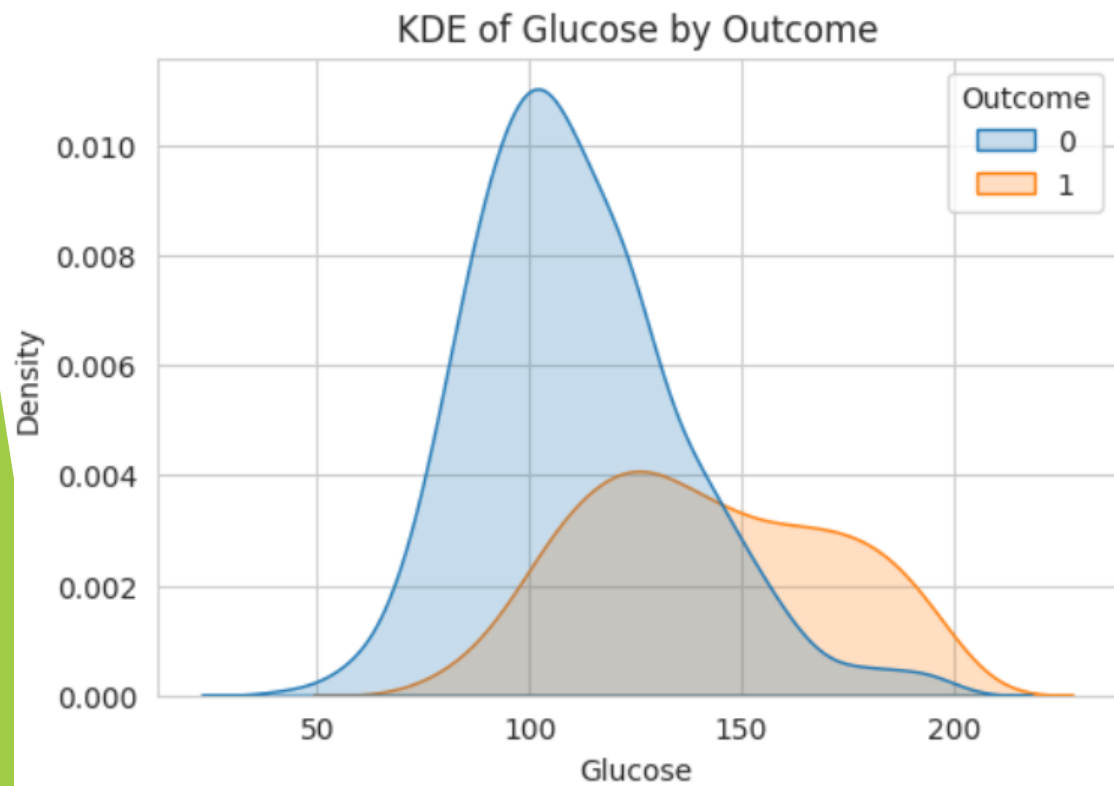
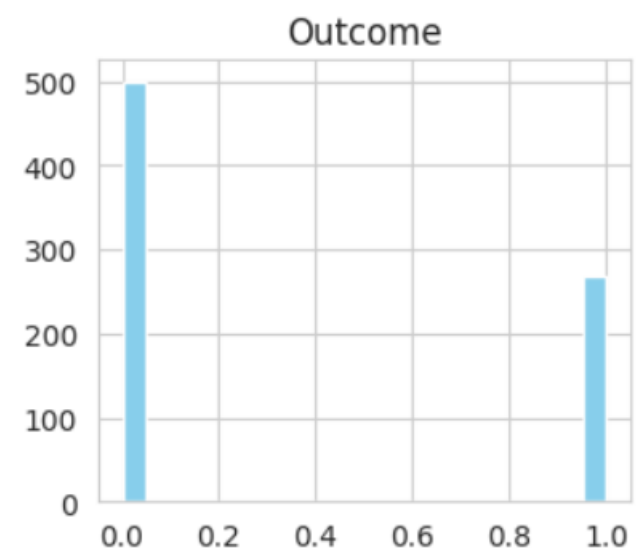
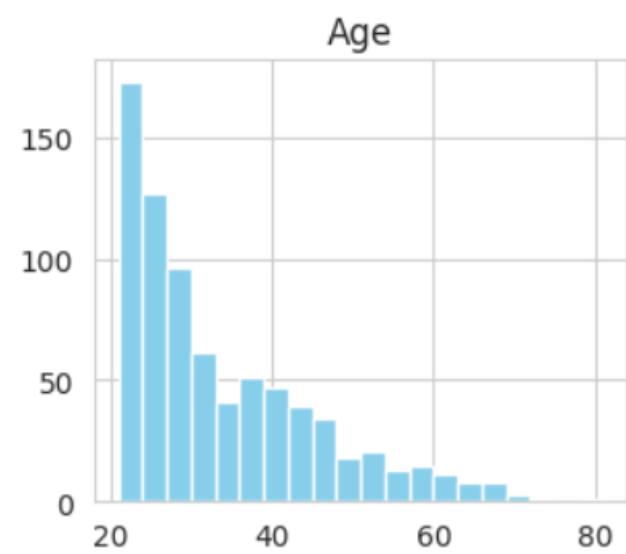
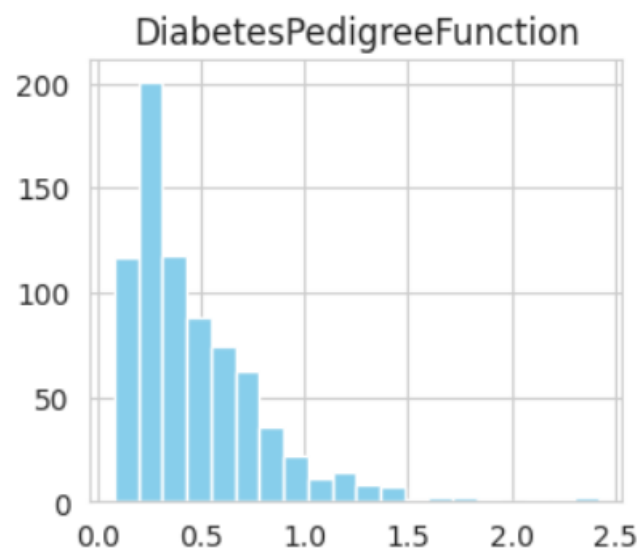


Insulin

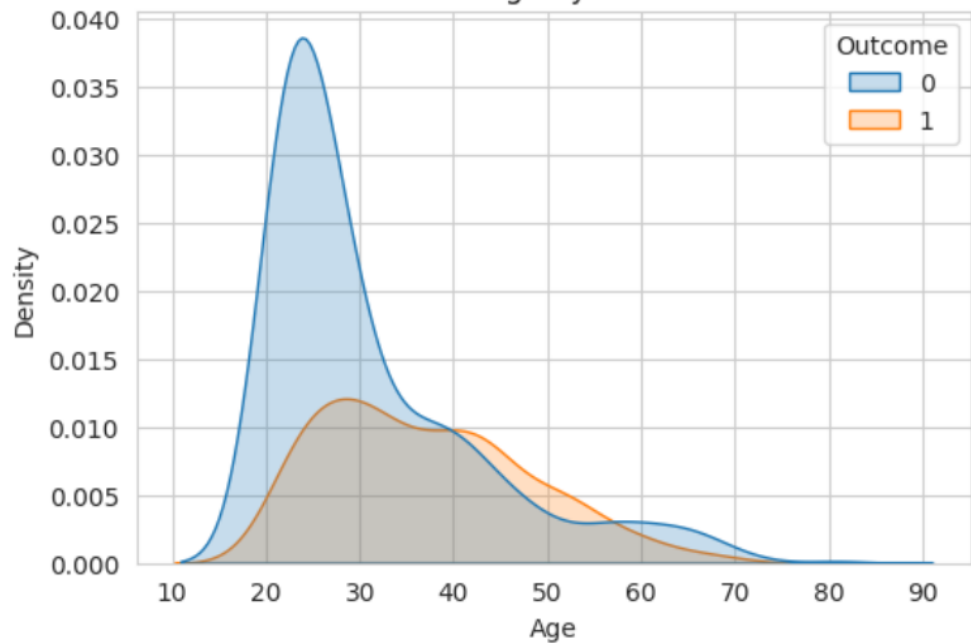


BMI

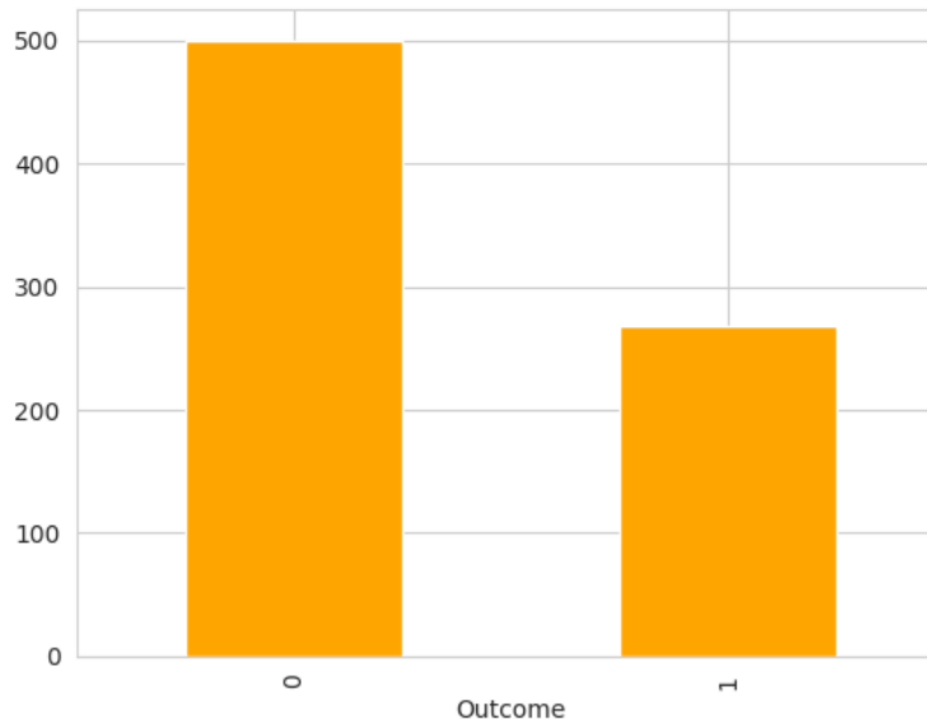




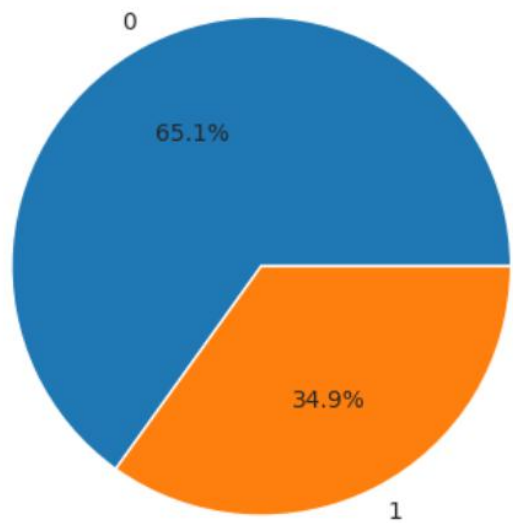
KDE of Age by Outcome



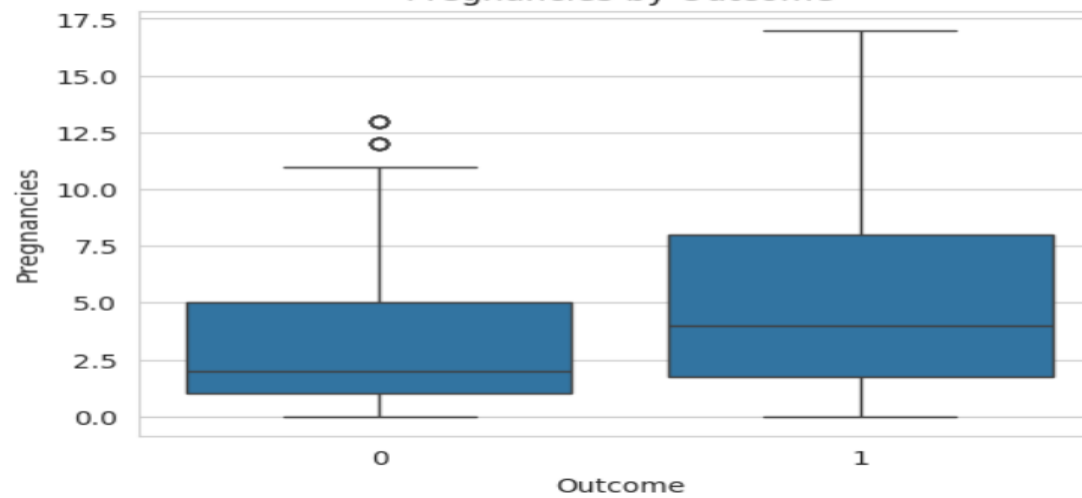
Outcome Distribution

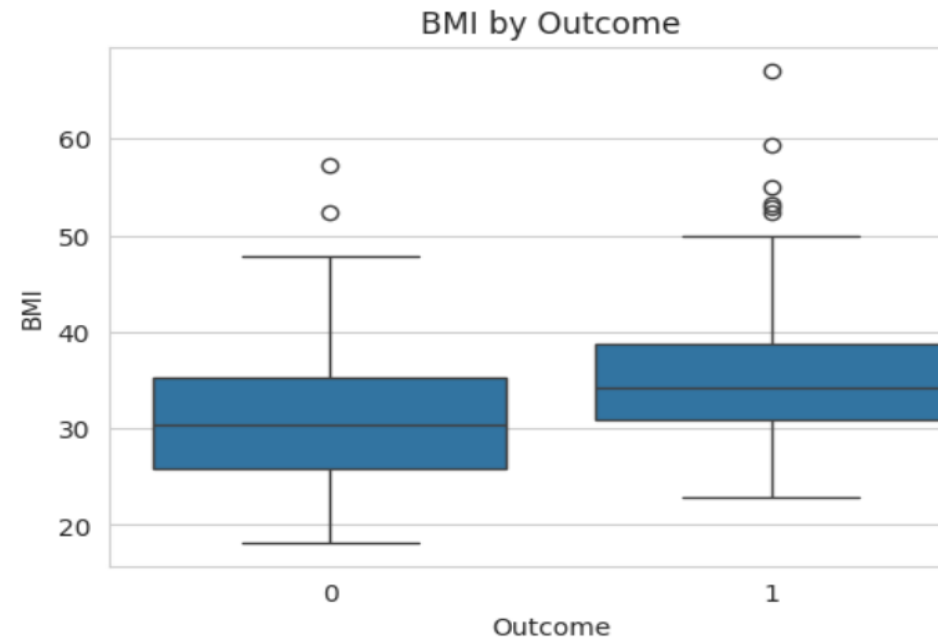
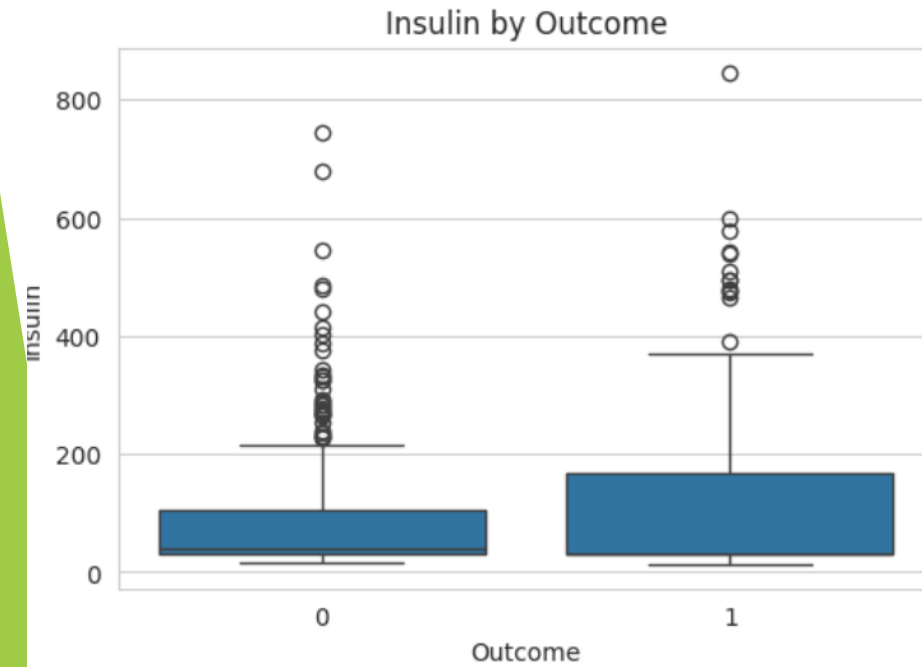
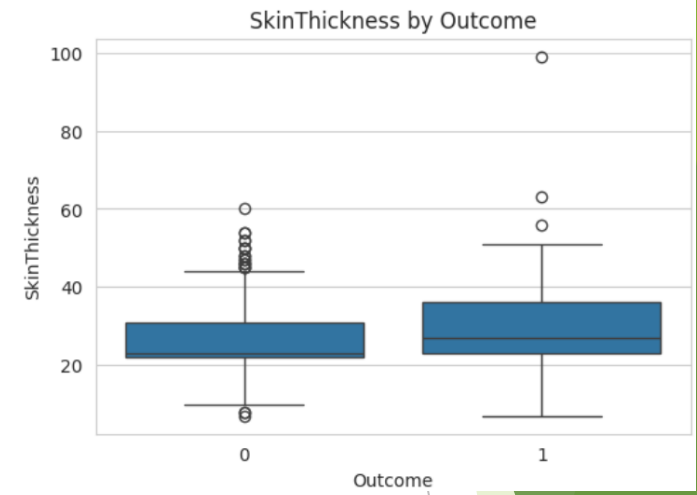
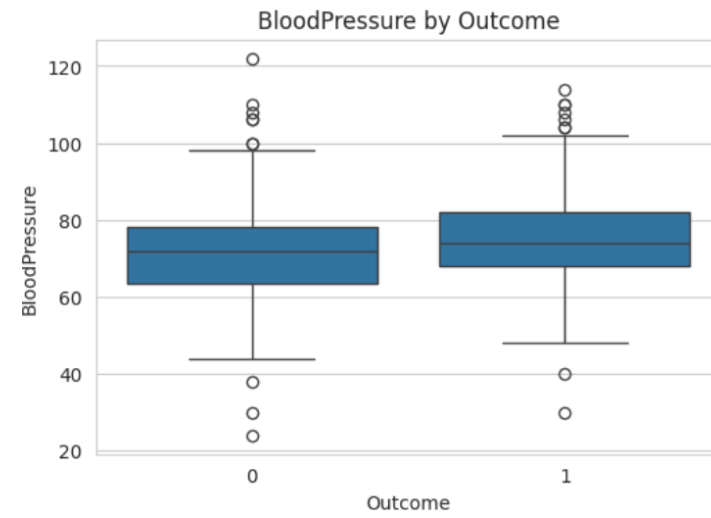
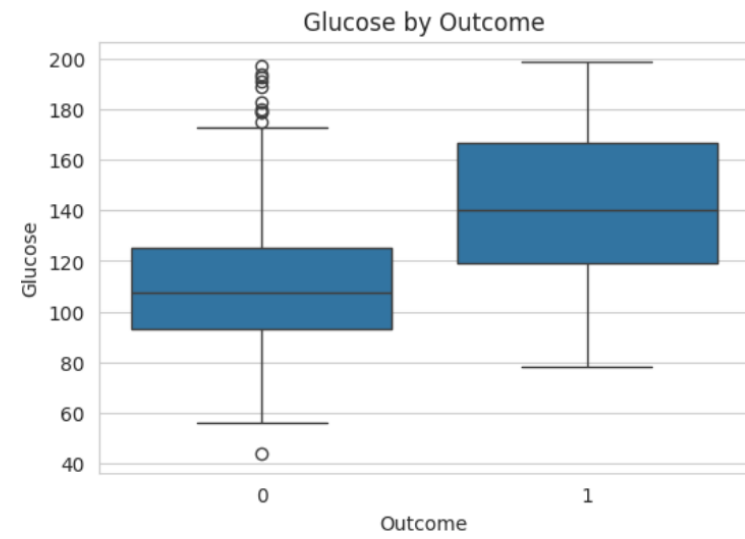


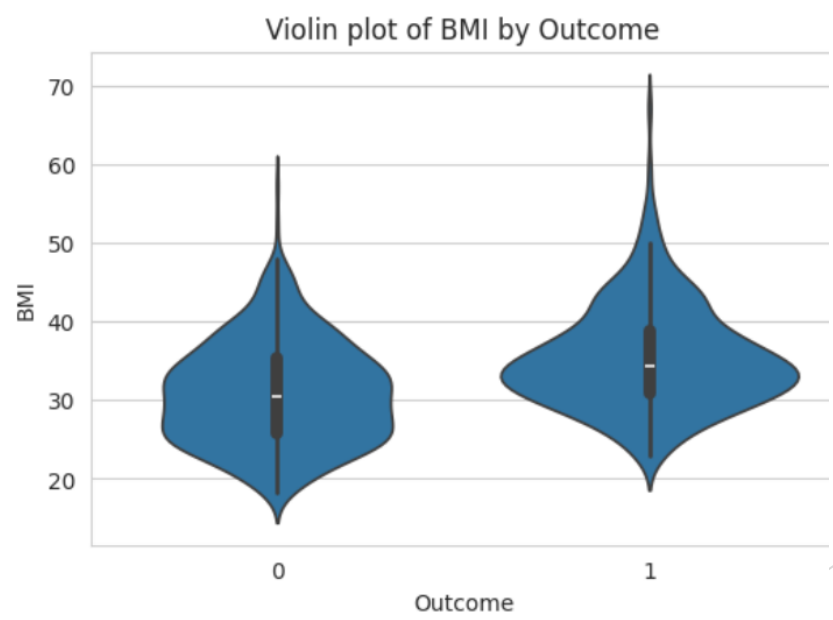
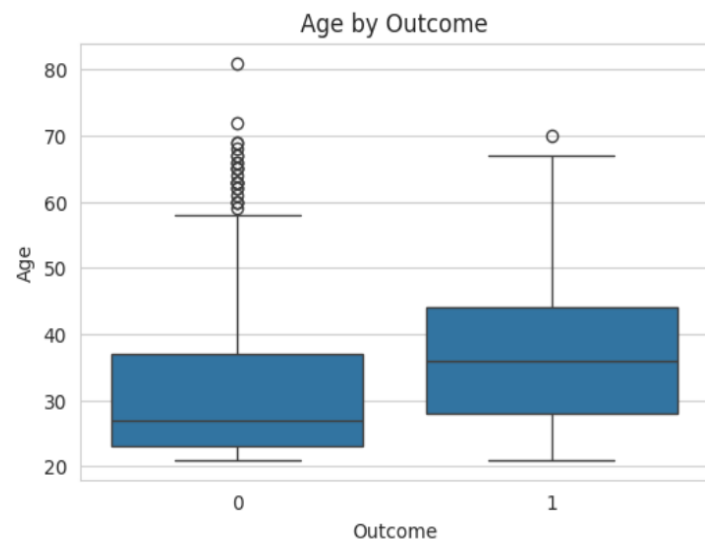
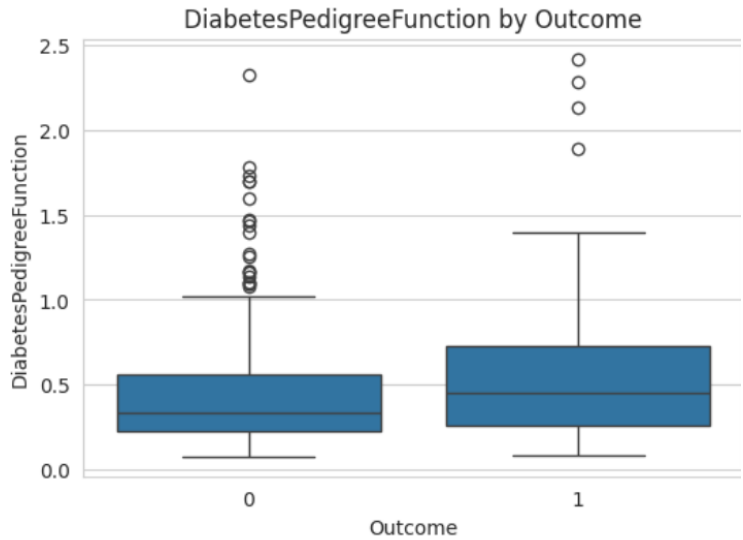
Outcome Distribution (Pie)

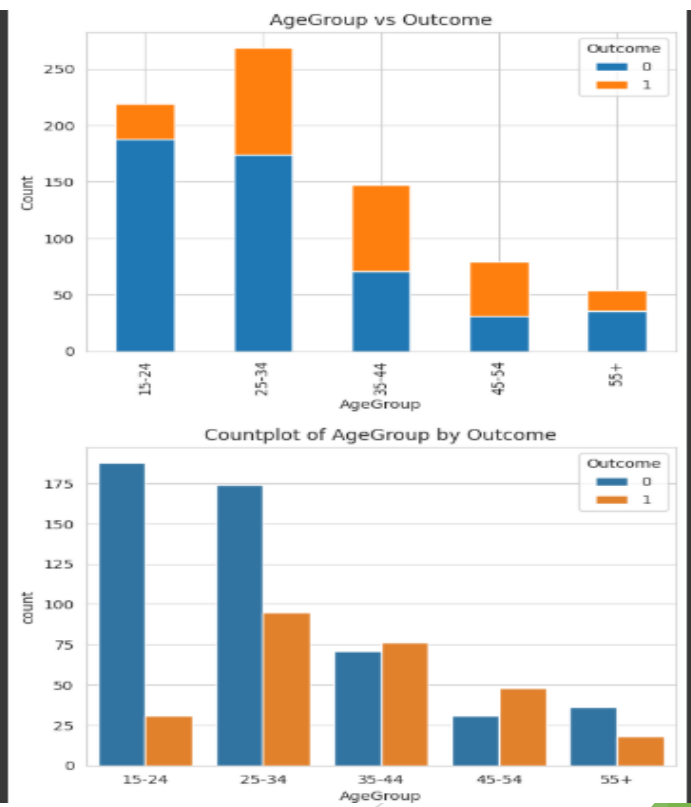
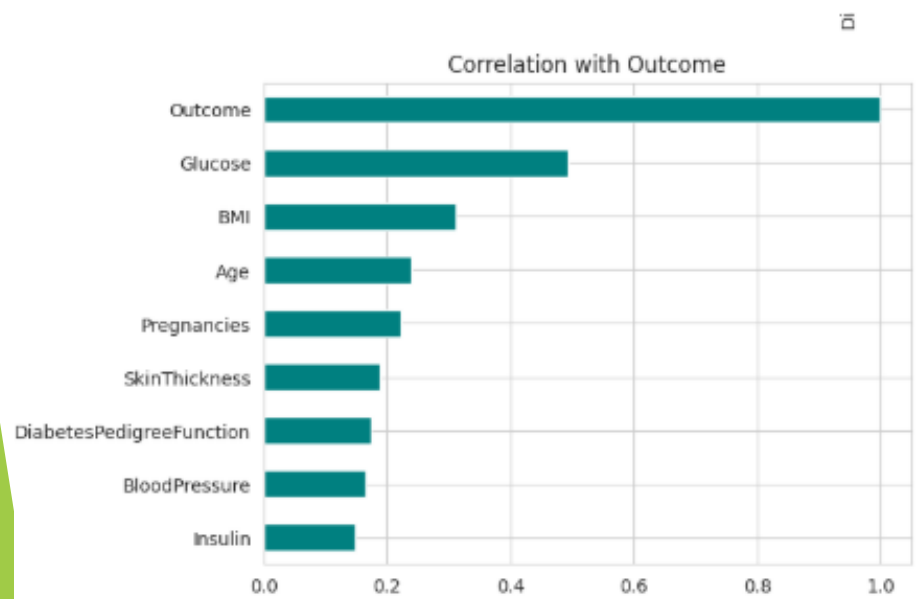
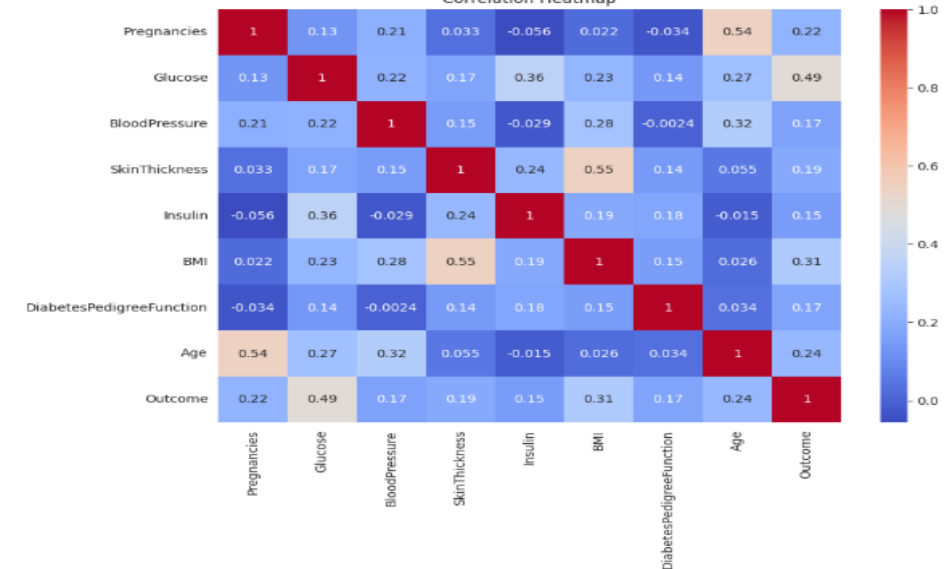
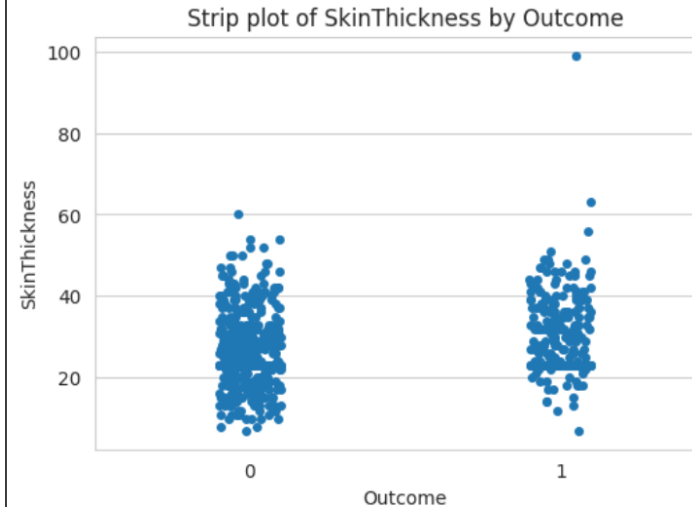
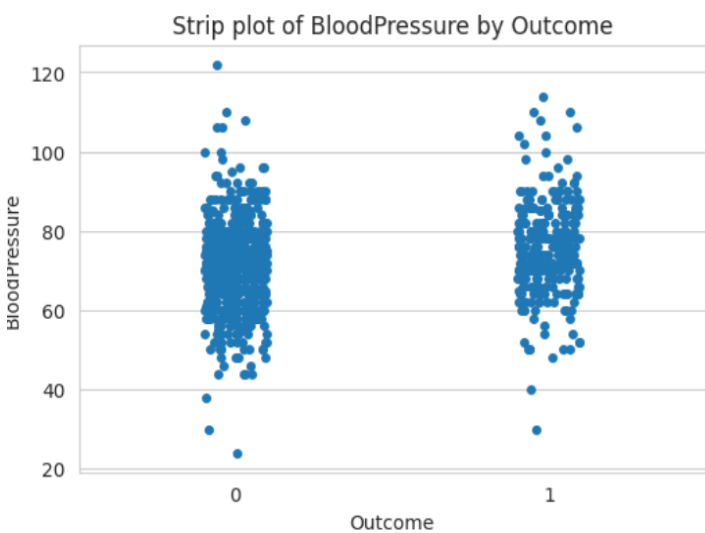


Pregnancies by Outcome











Thank You

SOUMYAJITA BOSE - IITP000180

✉ soumyajita_2312res655@iitp.ac.in

