

20C14001 - Le Duong Tuan Anh

Course: Data Mining

Homework 2

Student ID: 20C14001

Student Name: Le Duong Tuan Anh

Dataset: “diabetes.csv”

1) Describe the dataset

Dataset type: Record.

This dataset has 2,000 rows and 9 columns. The last column, “Outcome”, shows the result that the person has a diabetes or not.

Column	Type	Datatype	Has missing value?
Pregnancies	Interval-scaled	Int64	No
Glucose	Interval-scaled	Int64	No
BloodPressure	Interval-scaled	Int64	No
SkinThickness	Interval-scaled	Int64	No
Insulin	Interval-scaled	Int64	No
BMI	Interval-scaled	Float	No
DiabetesPedigreeFunction	Interval-scaled	Float	No
Age	Interval-scaled	Int64	No
Outcome	Categorical Data Binary (0/1)	Int64	No

Pregnancies: Number of times pregnant

Glucose: Plasma Glucose Concentration.

BloodPressure: Diastolic Blood Pressure.

SkinThickness: Estimate body fat.

Insulin: 2-Hour Serum Insulin.

BMI: Body Mass Index.

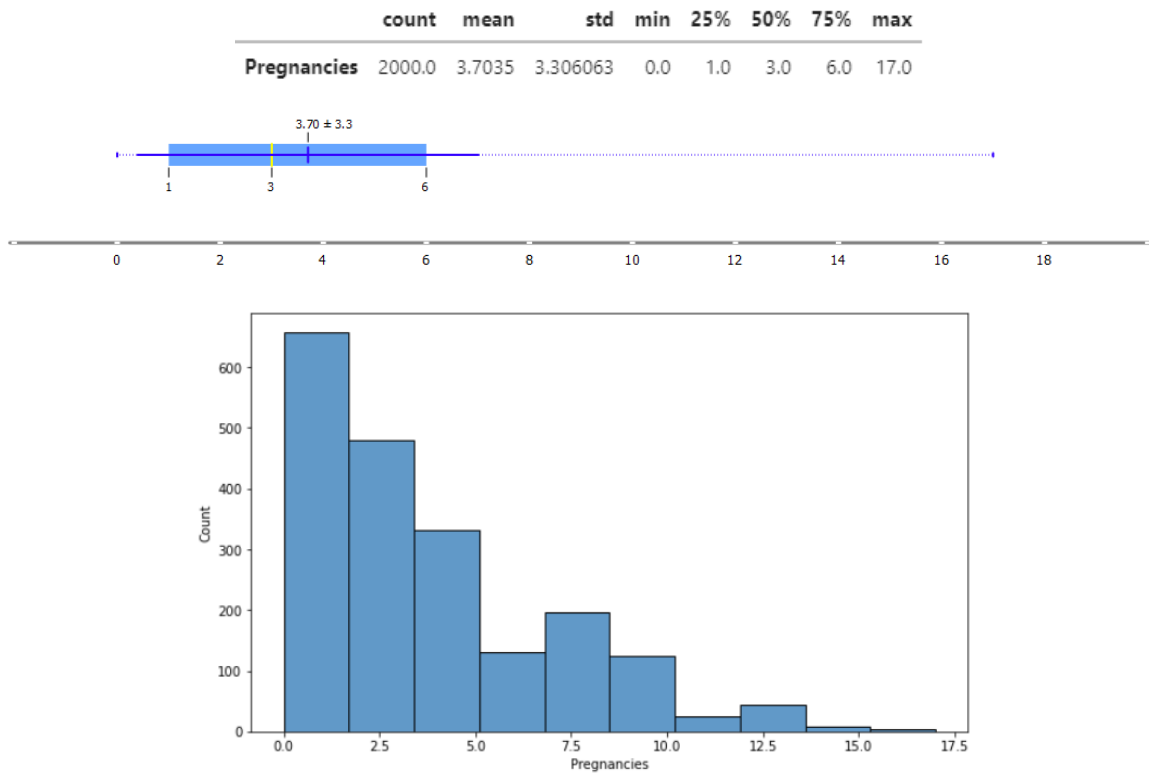
DiabetesPedigreeFunction: Information about diabetes history in relatives and genetics.

Age: Age (years).

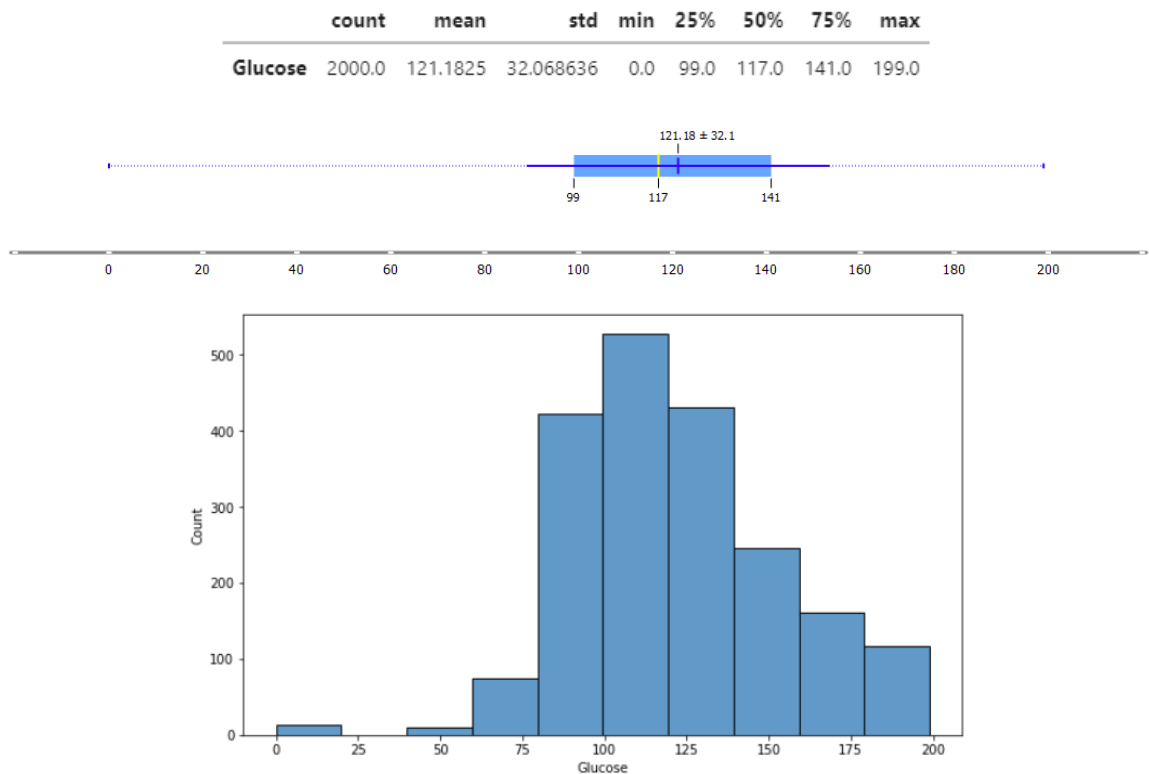
Outcome: 0 = Diabetic, 1 = Not Diabetic

2) Apply basic statistical descriptions for the dataset

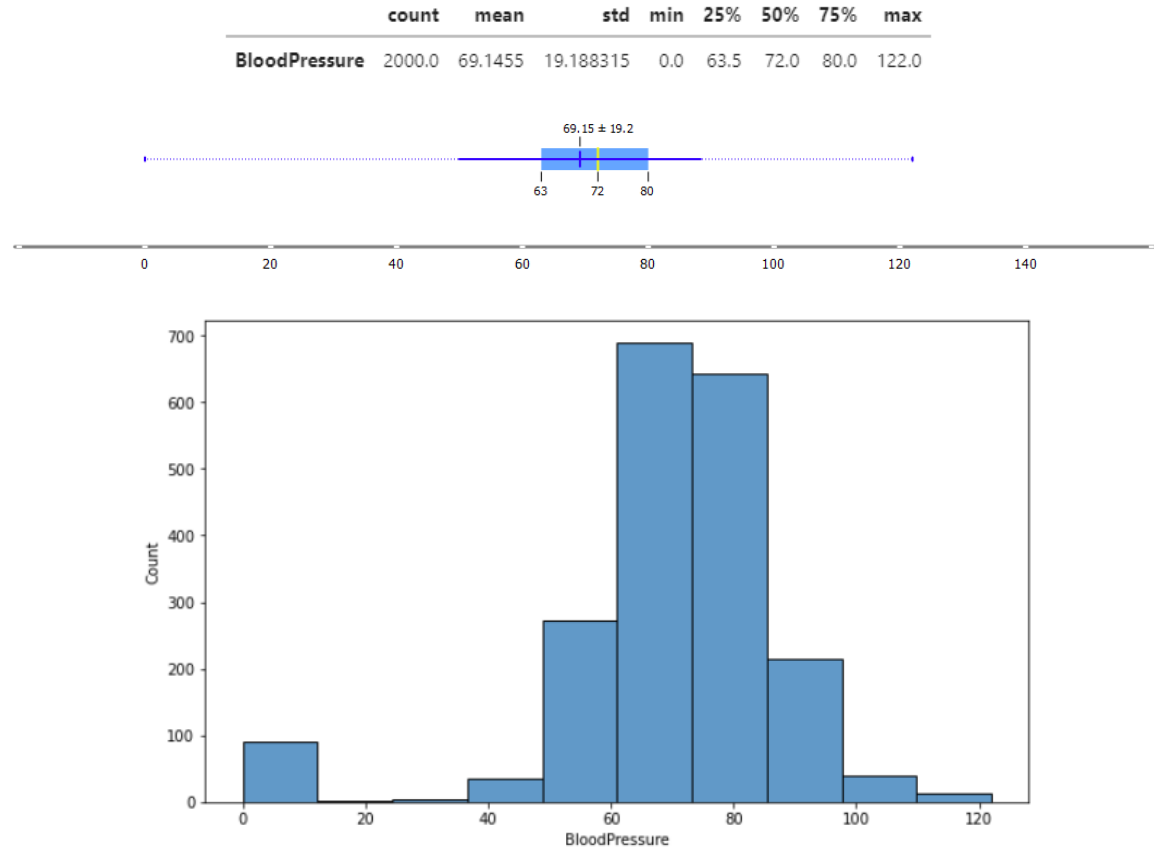
2.1. Pregnancies



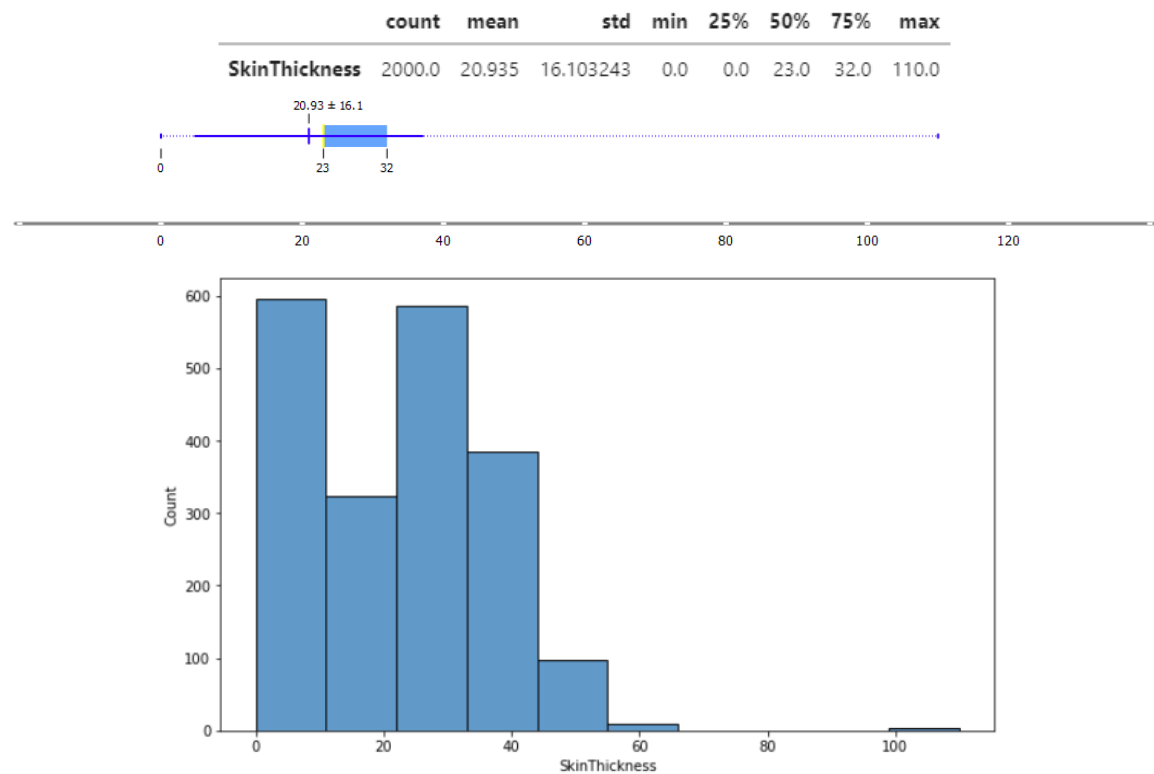
2.2. Glucose



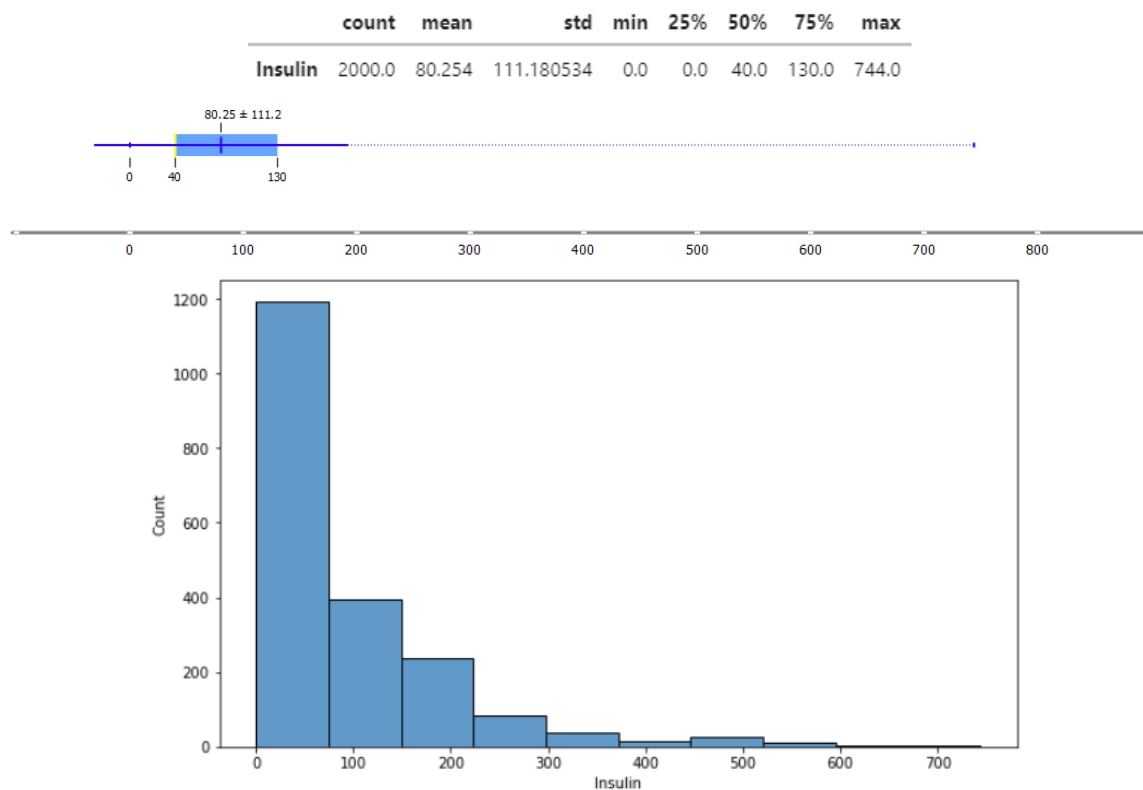
2.3. BloodPressure



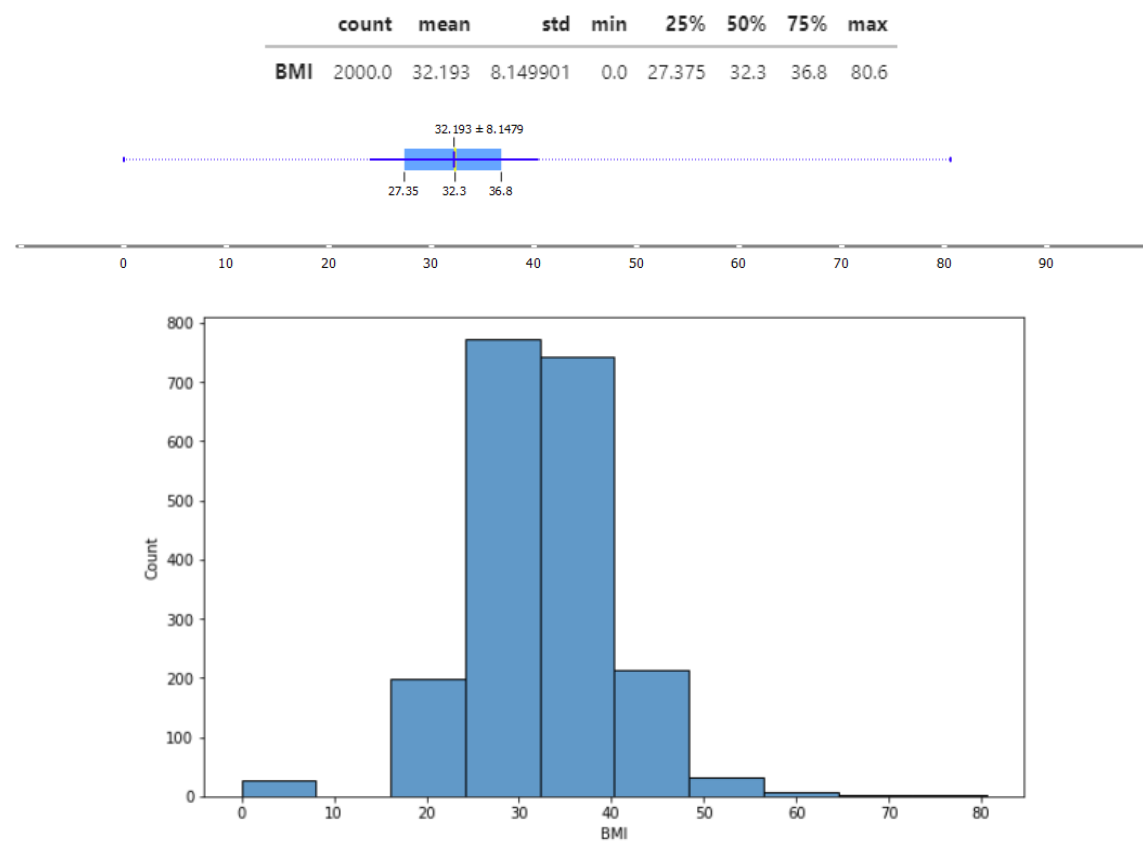
2.4. SkinThickness



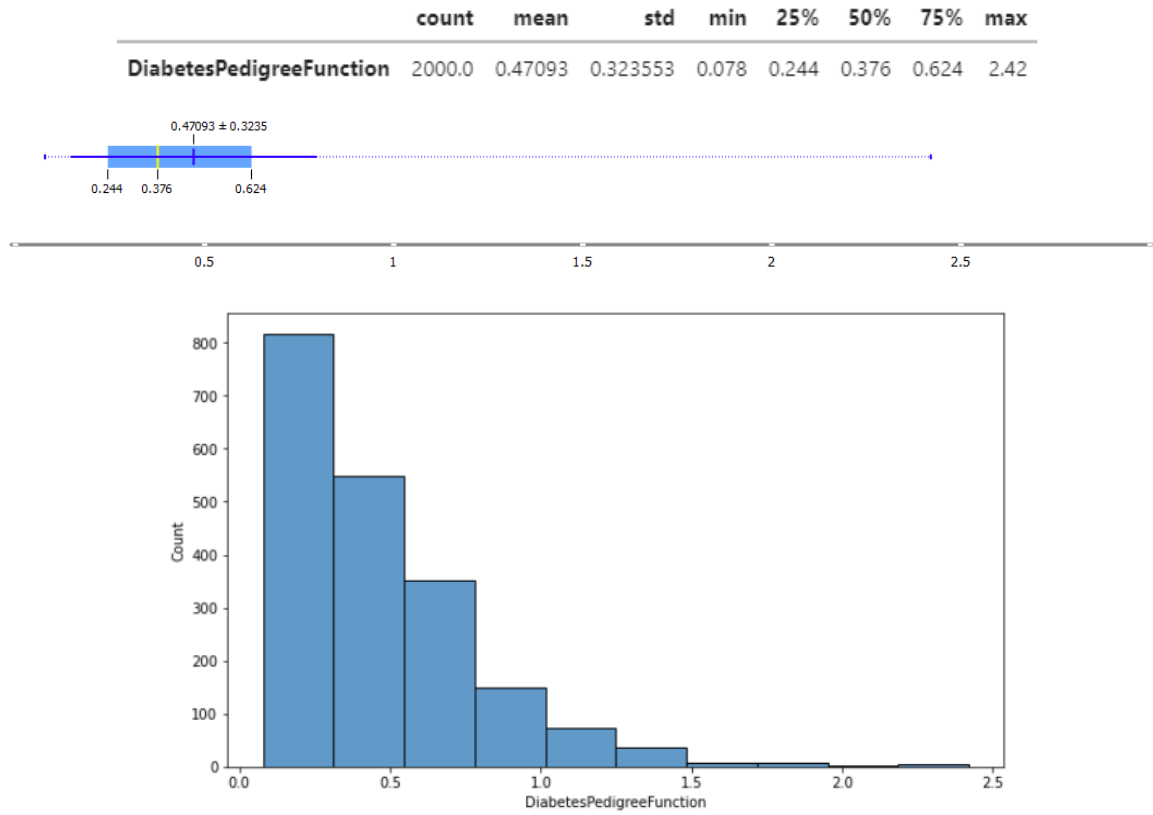
2.5. Insulin



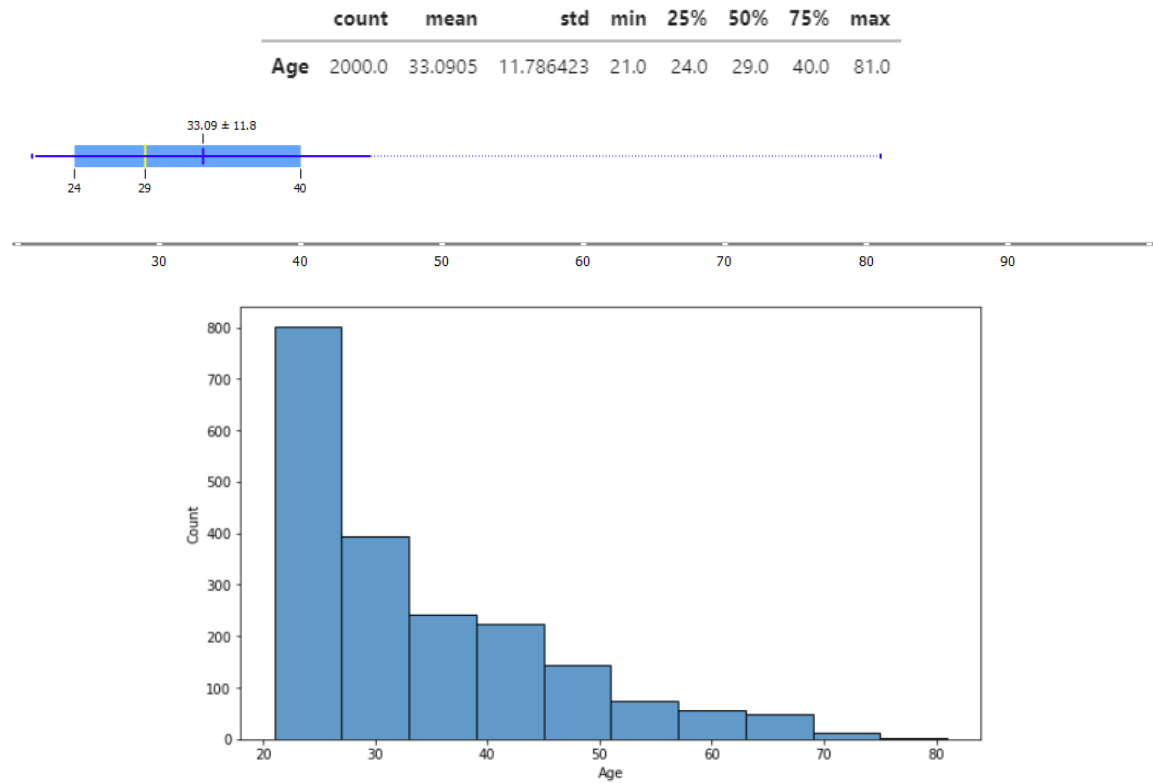
2.6. BMI



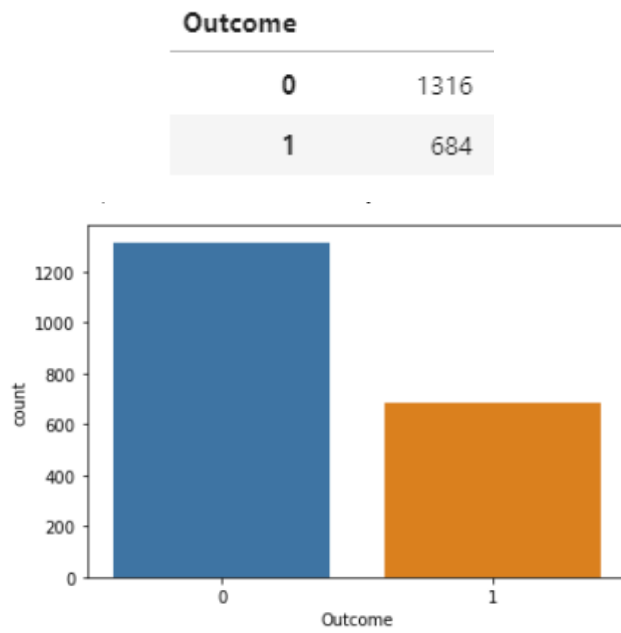
2.7. DiabetesPedigreeFunction



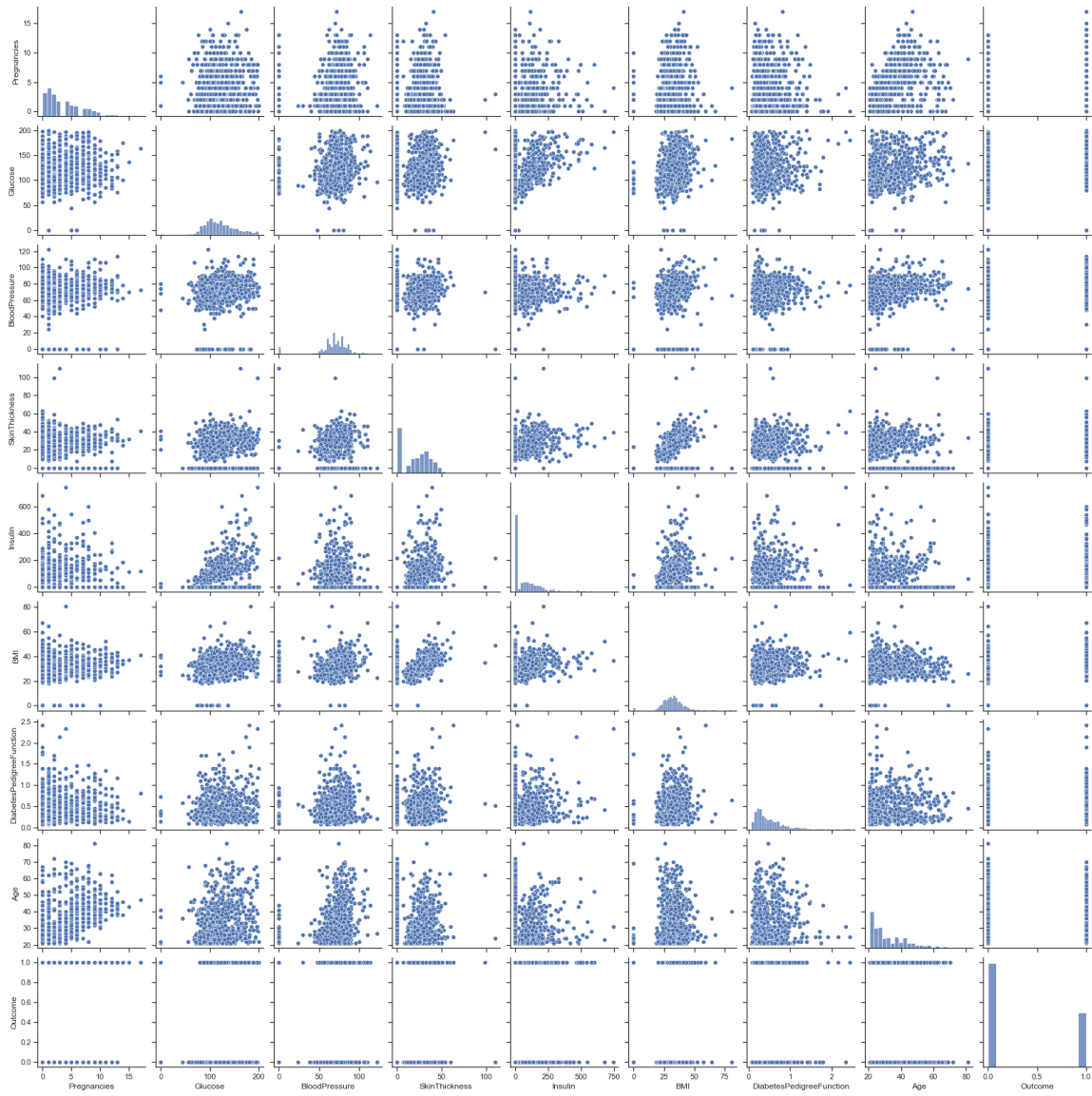
2.8. Age



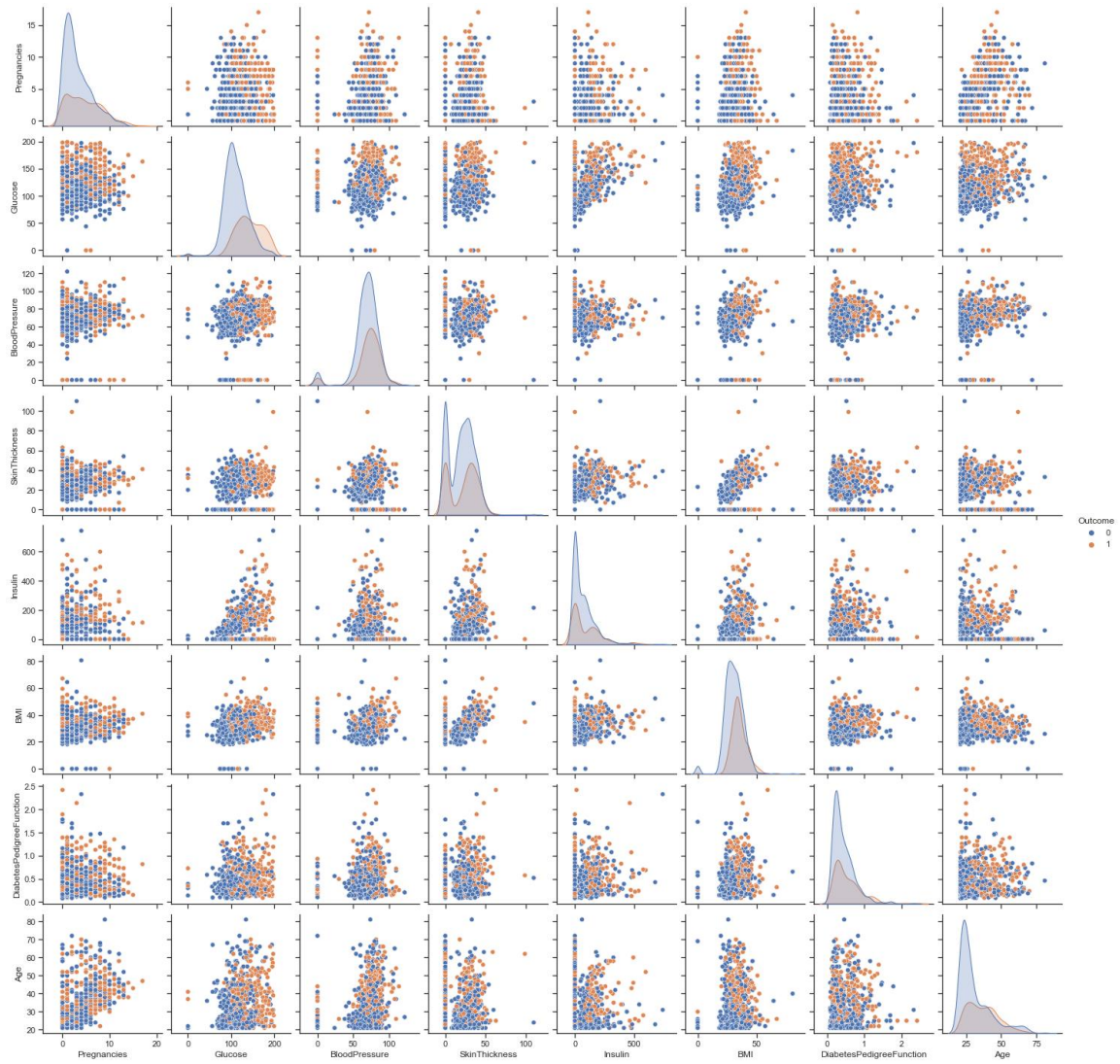
2.9. Outcome



3) Visualize this dataset by using scatterplot matrix.
Scatter Plot without distinct Outcome result.



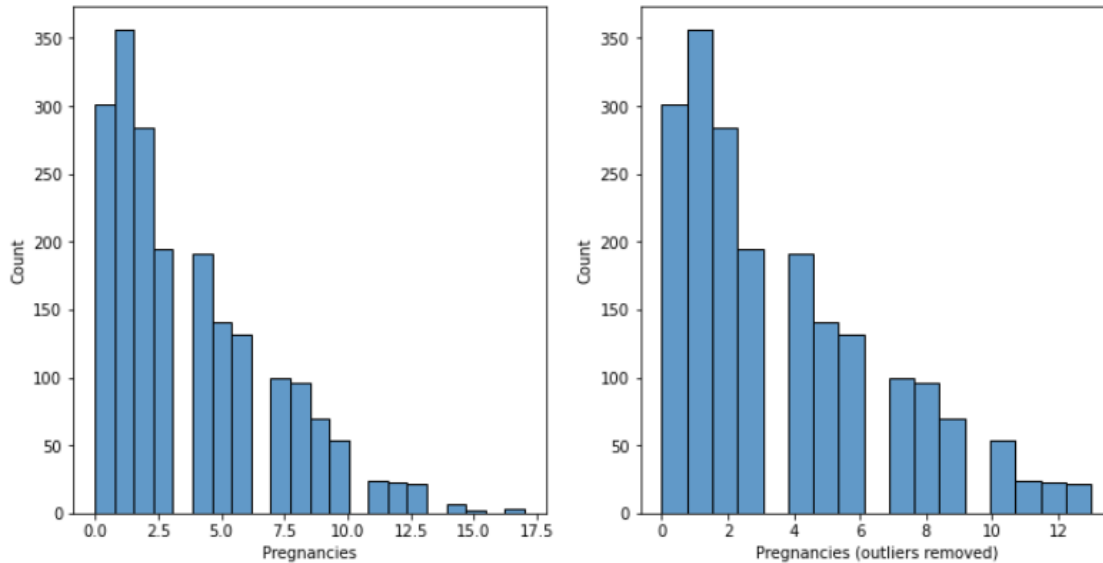
Scatter Plot with distinct Outcome result.



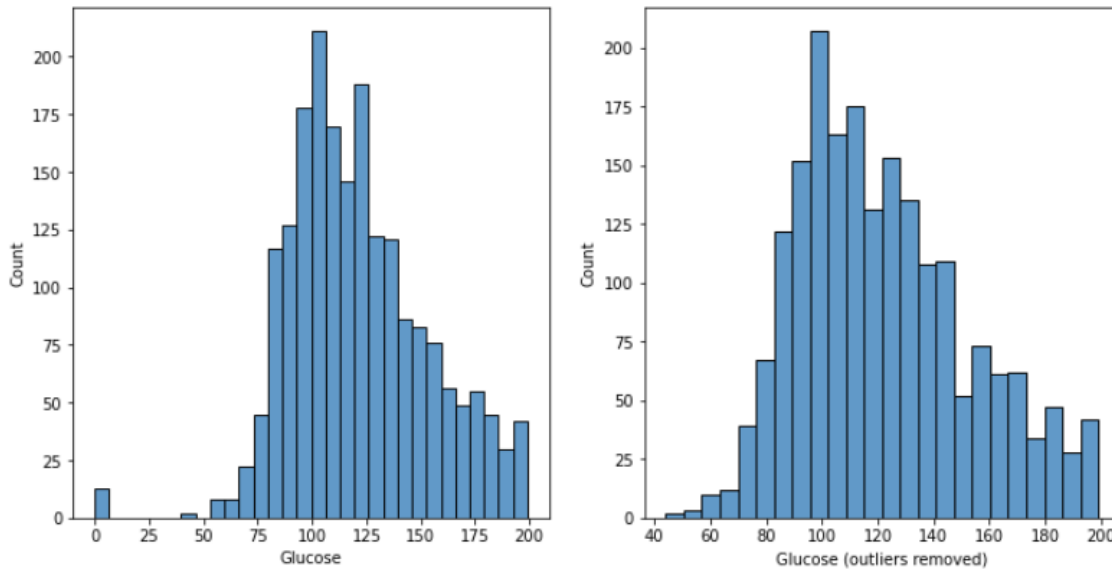
4) Do we have outliers in this dataset? For each attribute, list them out.

We would consider outliers by using Inter-quartile Range method. In case the data point is out of range $[Q1 - 1.5 \cdot IQR, Q3 + 1.5 \cdot IQR]$, it might be outliers.

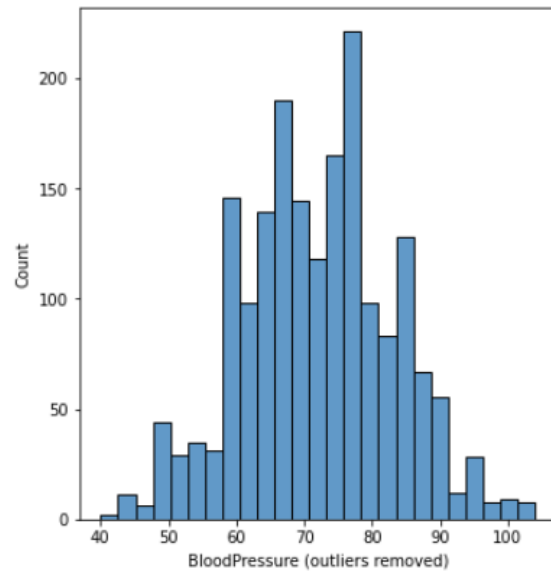
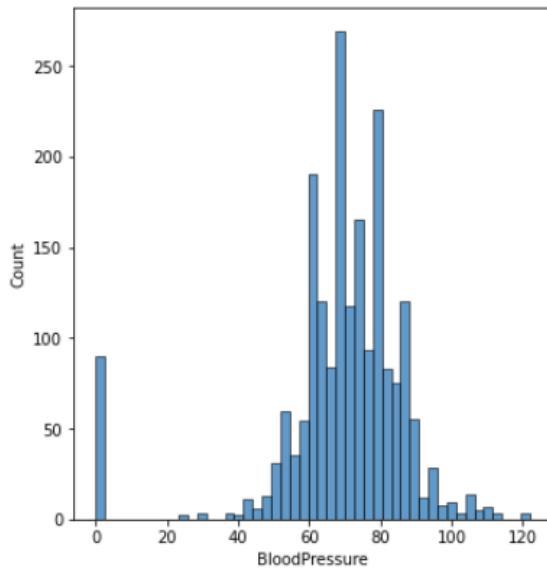
Attribute: Pregnancies. $Q1 = 1.0$; $Q3 = 6.0$; $IQR = 7.5$. Total outlier(s): 12



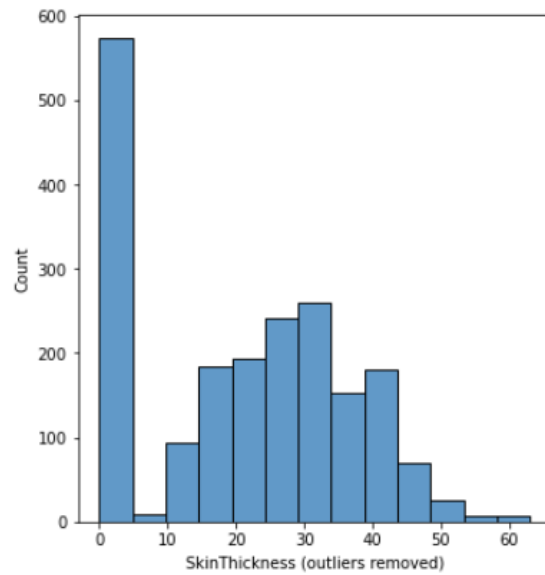
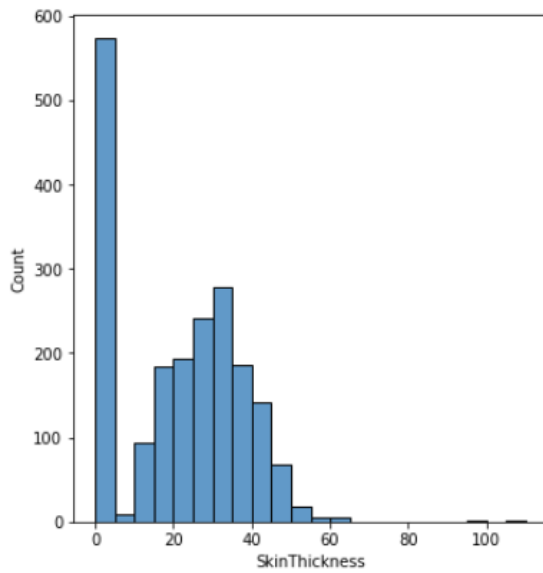
Attribute: Glucose. $Q1 = 99.0$; $Q3 = 141.0$; $IQR = 63.0$. Total outlier(s): 13



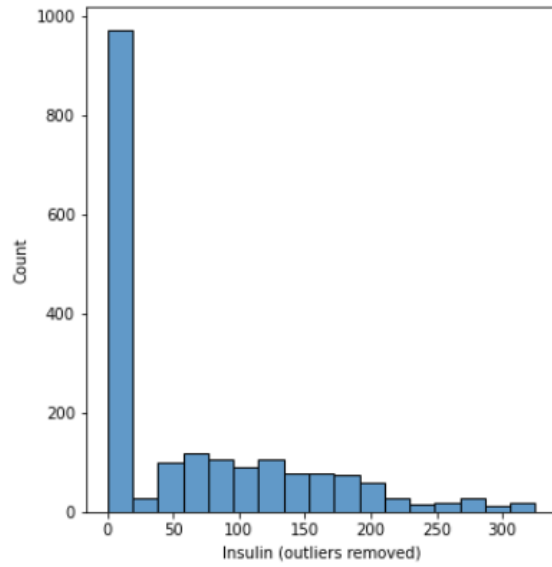
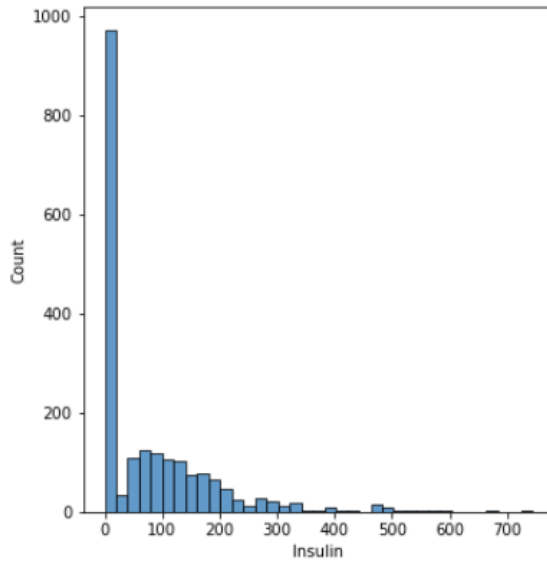
Attribute: BloodPressure. Q1 = 63.5 ; Q3 = 80.0 ; IQR = 24.75. Total outlier(s): 125



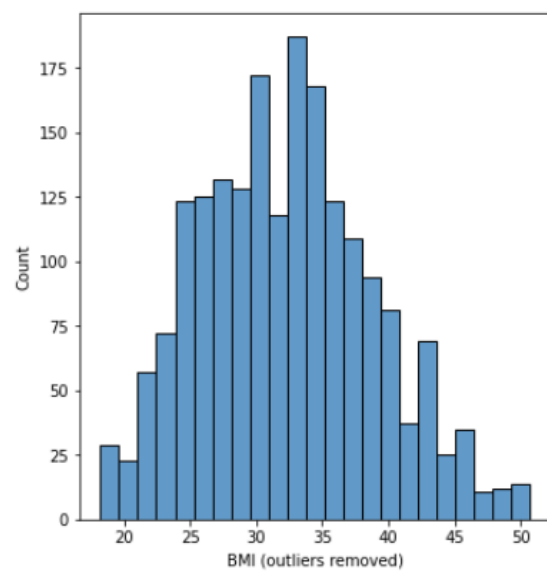
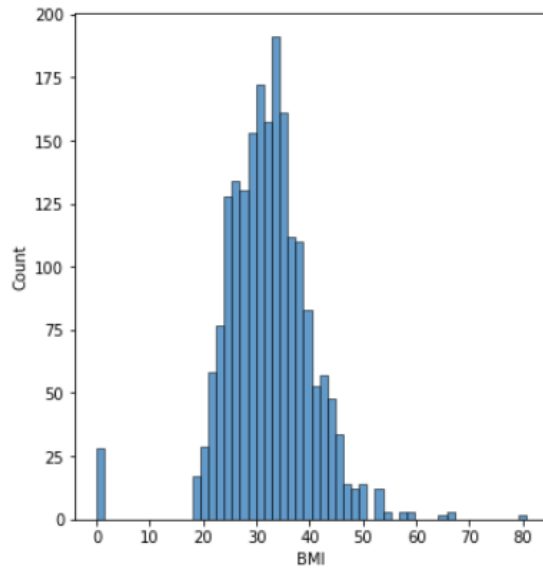
Attribute: SkinThickness. Q1 = 0.0 ; Q3 = 32.0 ; IQR = 48.0. Total outlier(s): 4



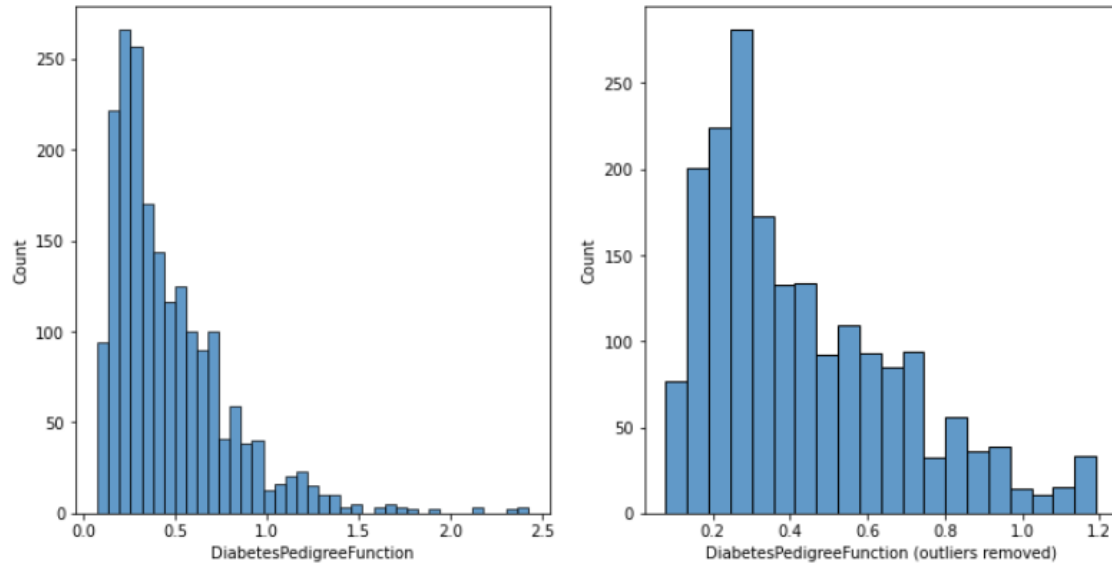
Attribute: Insulin. Q1 = 0.0 ; Q3 = 130.0 ; IQR = 195.0. Total outlier(s): 73



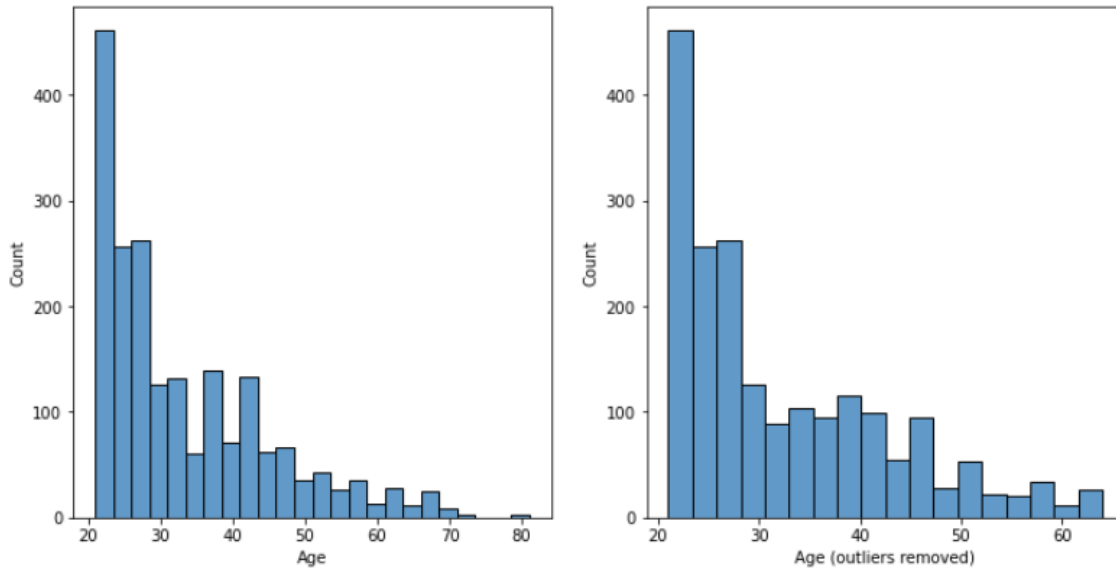
Attribute: BMI. Q1 = 27.375 ; Q3 = 36.8 ; IQR = 14.137499999999996. Total outlier(s): 56



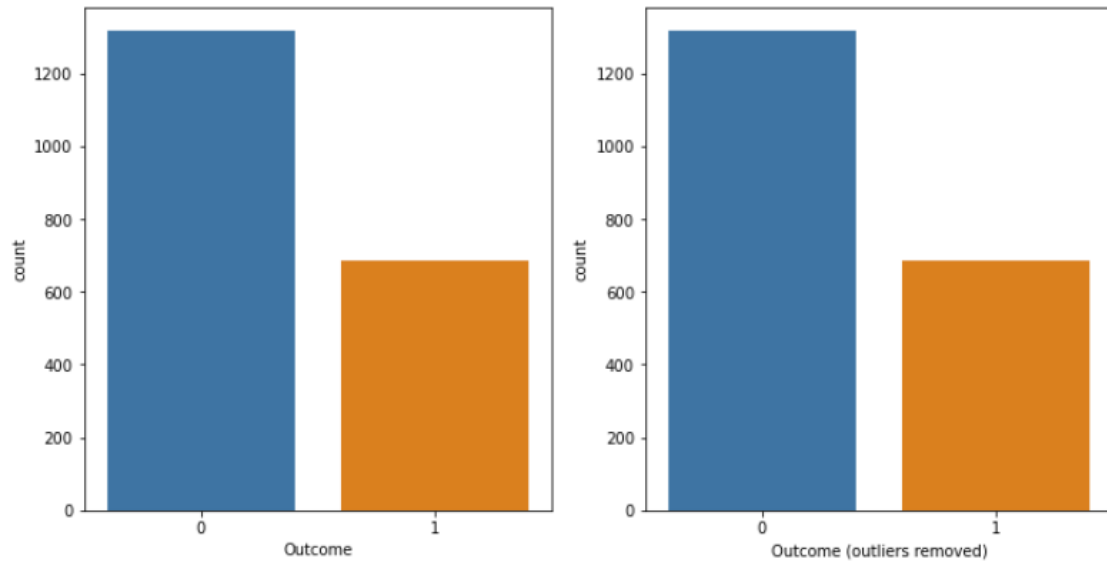
Attribute: DiabetesPedigreeFunction. Q1 = 0.244 ; Q3 = 0.624 ; IQR = 0.5700000000000001. Total outlier(s): 68



Attribute: Age. Q1 = 24.0 ; Q3 = 40.0 ; IQR = 24.0. Total outlier(s): 48



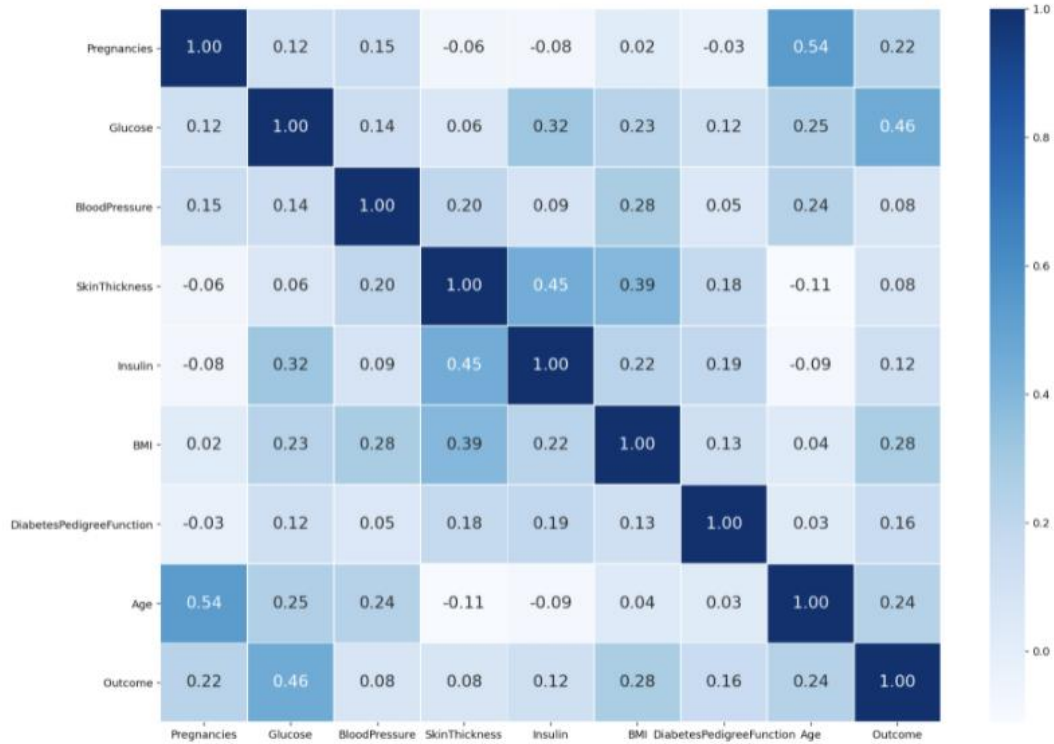
Attribute: Outcome. Q1 = 0.0 ; Q3 = 1.0 ; IQR = 1.5. Total outlier(s): 0



5, 6) Check the relationship of the two attributes

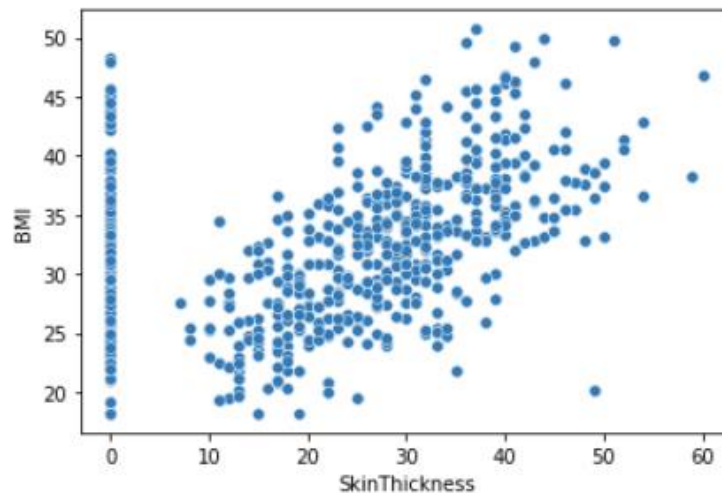
The correlation of dataset is show as below.

(Standard correlation co-efficient, **Pearson**)



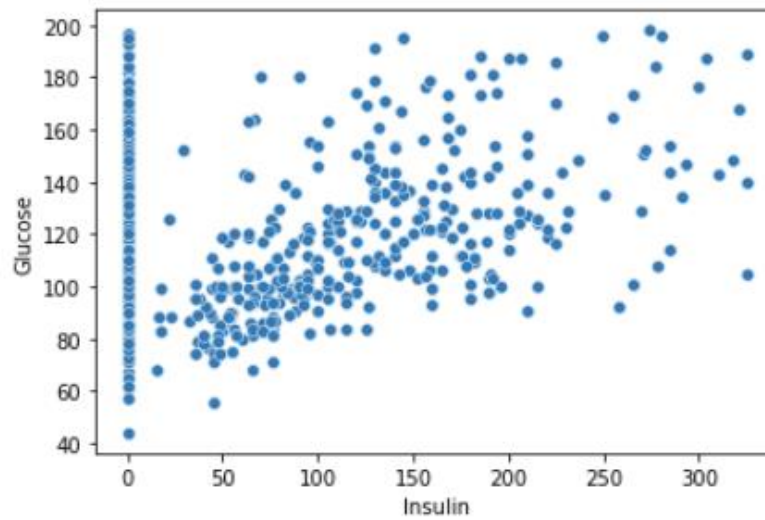
5) *SkinThickness* and *BMI*.

The correlation between SkinThickness and BMI is quite strong, 0.45.



Note that, the SkinThickness should not be 0 (as the thickness of skin is quite not reasonable when equal to 0).

6) *Insulin* và *Glucose*.



Note that, the Insulin should not be 0.

7) *Standardize this dataset such that all attributes have the same data unit.*

We standardize features by removing the mean, then scaling dataset to unit variance. Note that we will not standardize the **Outcome** column, as it is a **categorical column**.

The standardized score is calculated as below:

$$Z = (X - U) / S$$

Where:

- X: The value of sample in dataset.
- U: The mean of dataset.
- S: The standard deviation of dataset.

Code:

```
df = data.drop(['Outcome'], axis=1)
```

```
df.mean()
```

```
Pregnancies      3.70350
Glucose           121.18250
BloodPressure     69.14550
SkinThickness     20.93500
Insulin           80.25400
BMI               32.19300
DiabetesPedigreeFunction  0.47093
Age               33.09050
dtype: float64
```

```
df.std()
```

```
Pregnancies      3.306063
Glucose           32.068636
BloodPressure     19.188315
SkinThickness     16.103243
Insulin           111.180534
BMI               8.149901
DiabetesPedigreeFunction  0.323553
Age               11.786423
dtype: float64
```

Compare histogram result between and after standardized:

