# 1. Introduction

## DTSA 5510 - Unsupervised Algorithms in Machine Learning Final Project
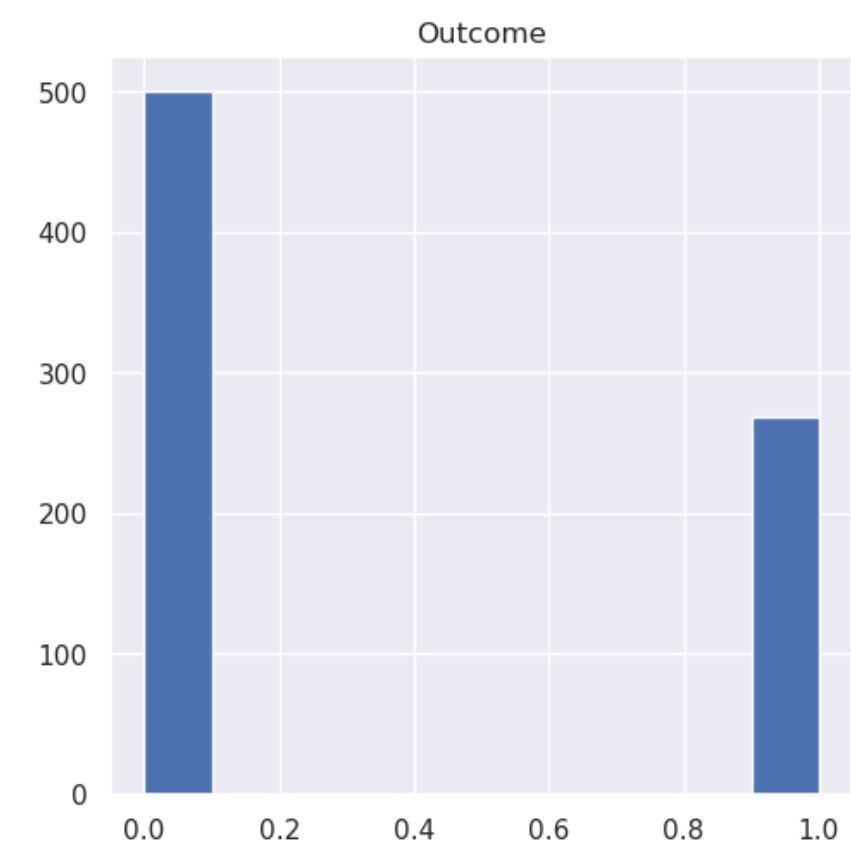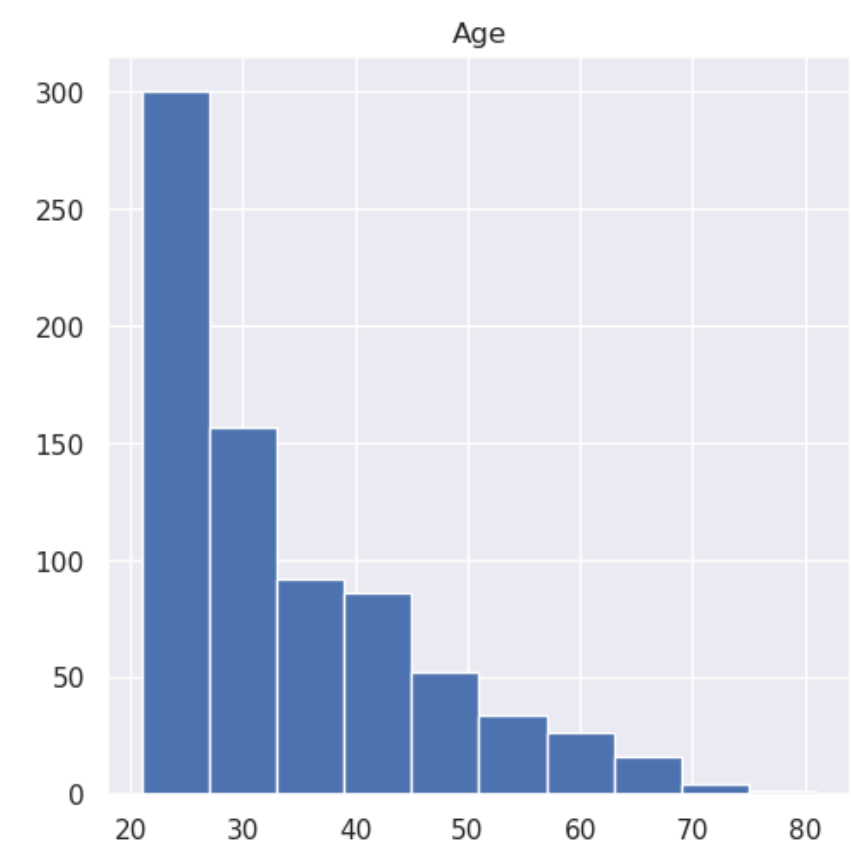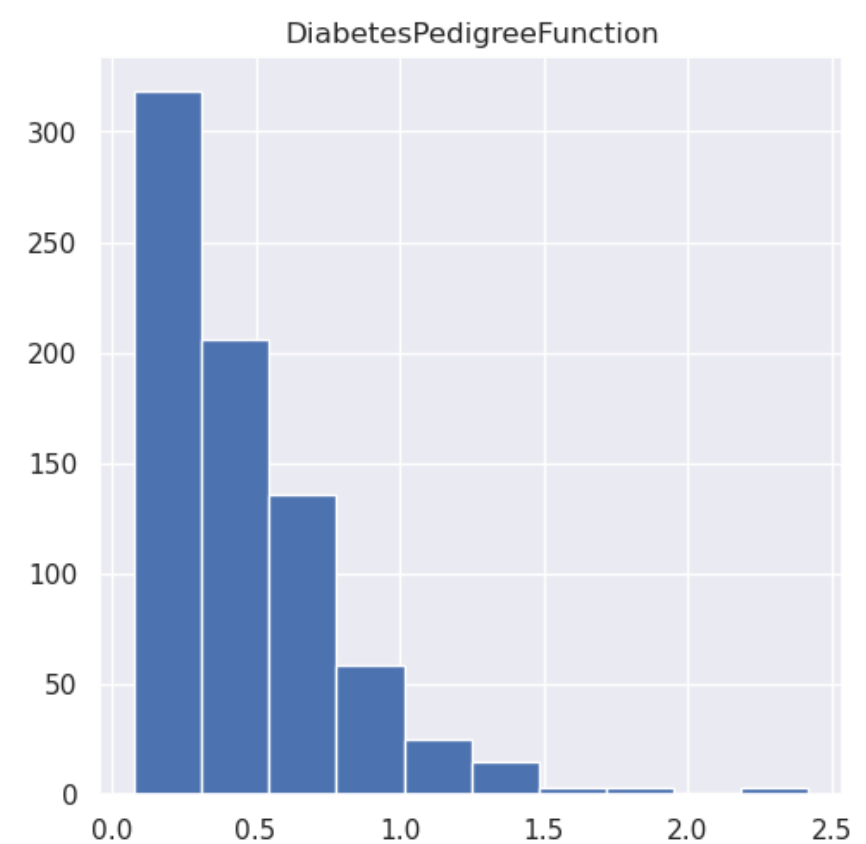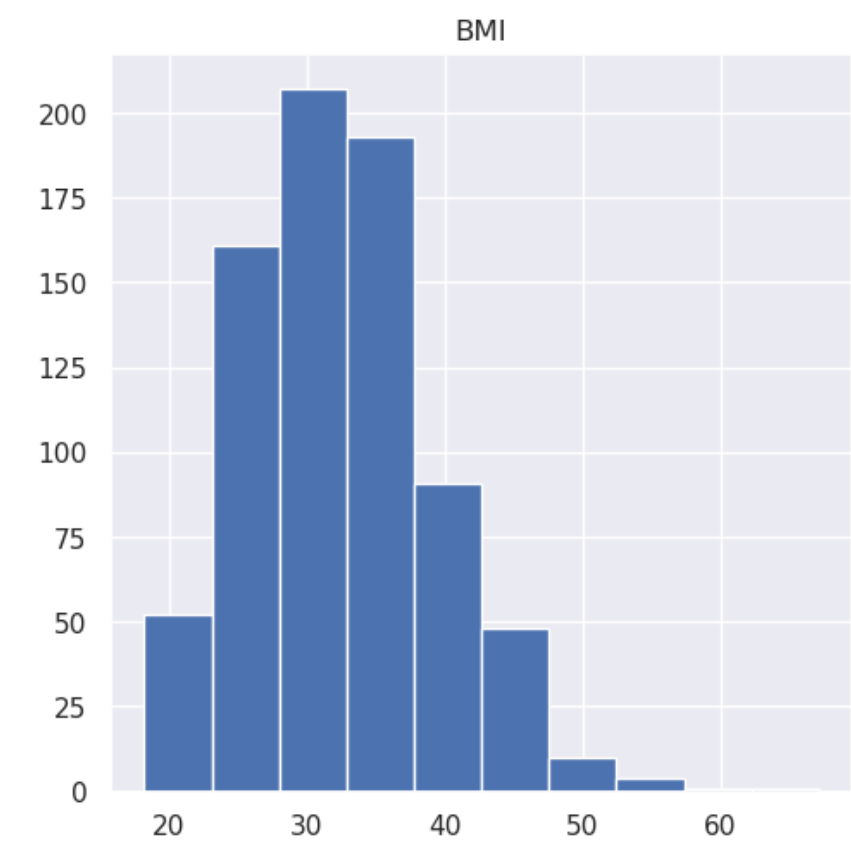
Feb, 22, 2023

# 1. Exploratory Data Analysis

```
In [55]: diabetes = pd.read_csv('/kaggle/input/pima-indians-diabetes-database/diabetes.csv')
         diabetes.head()
```

Out[55]:

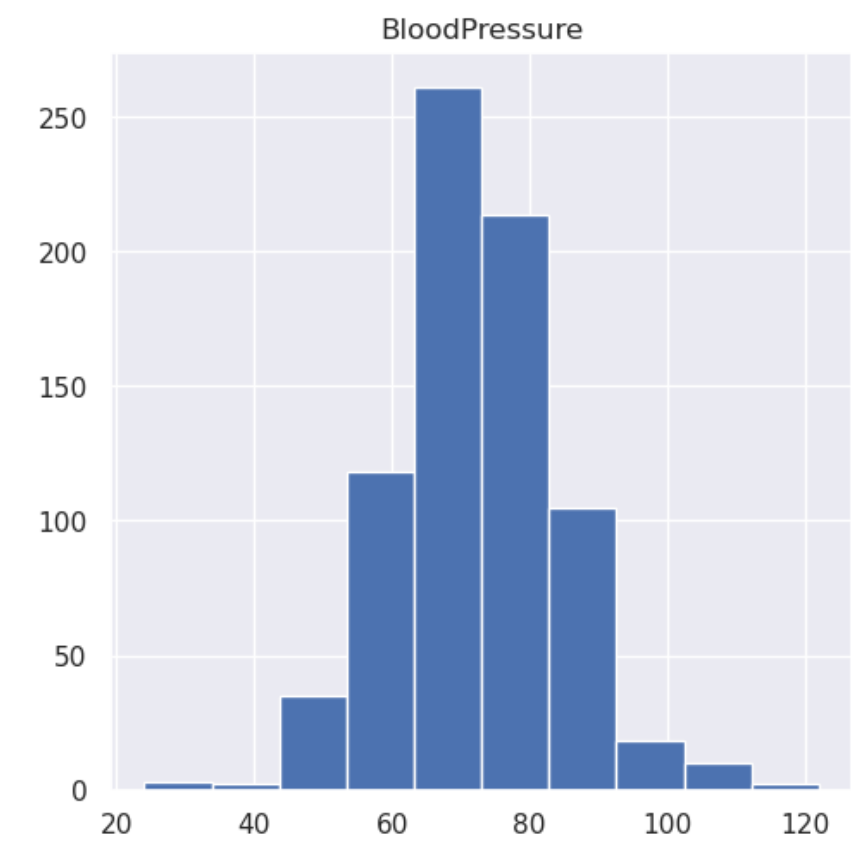| | Pregnancies | Glucose | BloodPressure | SkinThickness | Insulin | BMI | DiabetesPedigreeFunction | Age | Outcome |
|---|---|---|---|---|---|---|---|---|---|
| **0** | 6 | 148 | 72 | 35 | 0 | 33.6 | 0.627 | 50 | 1 |
| **1** | 1 | 85 | 66 | 29 | 0 | 26.6 | 0.351 | 31 | 0 |
| **2** | 8 | 183 | 64 | 0 | 0 | 23.3 | 0.672 | 32 | 1 |
| **3** | 1 | 89 | 66 | 23 | 94 | 28.1 | 0.167 | 21 | 0 |
| **4** | 0 | 137 | 40 | 35 | 168 | 43.1 | 2.288 | 33 | 1 |

# 2. Modeling

**Scaling**

$$z = \frac{x_i - \mu}{\sigma}$$

# Scaling

```python
StandardScaler = StandardScaler()
X = pd.DataFrame(StandardScaler.fit_transform(diabetes.drop(["Outcome"],axis = 1),),
       columns=['Pregnancies', 'Glucose', 'BloodPressure', 'SkinThickness', 'Insulin',
       'BMI', 'DiabetesPedigreeFunction', 'Age'])
```

```python
y = diabetes["Outcome"]
X.head()
```

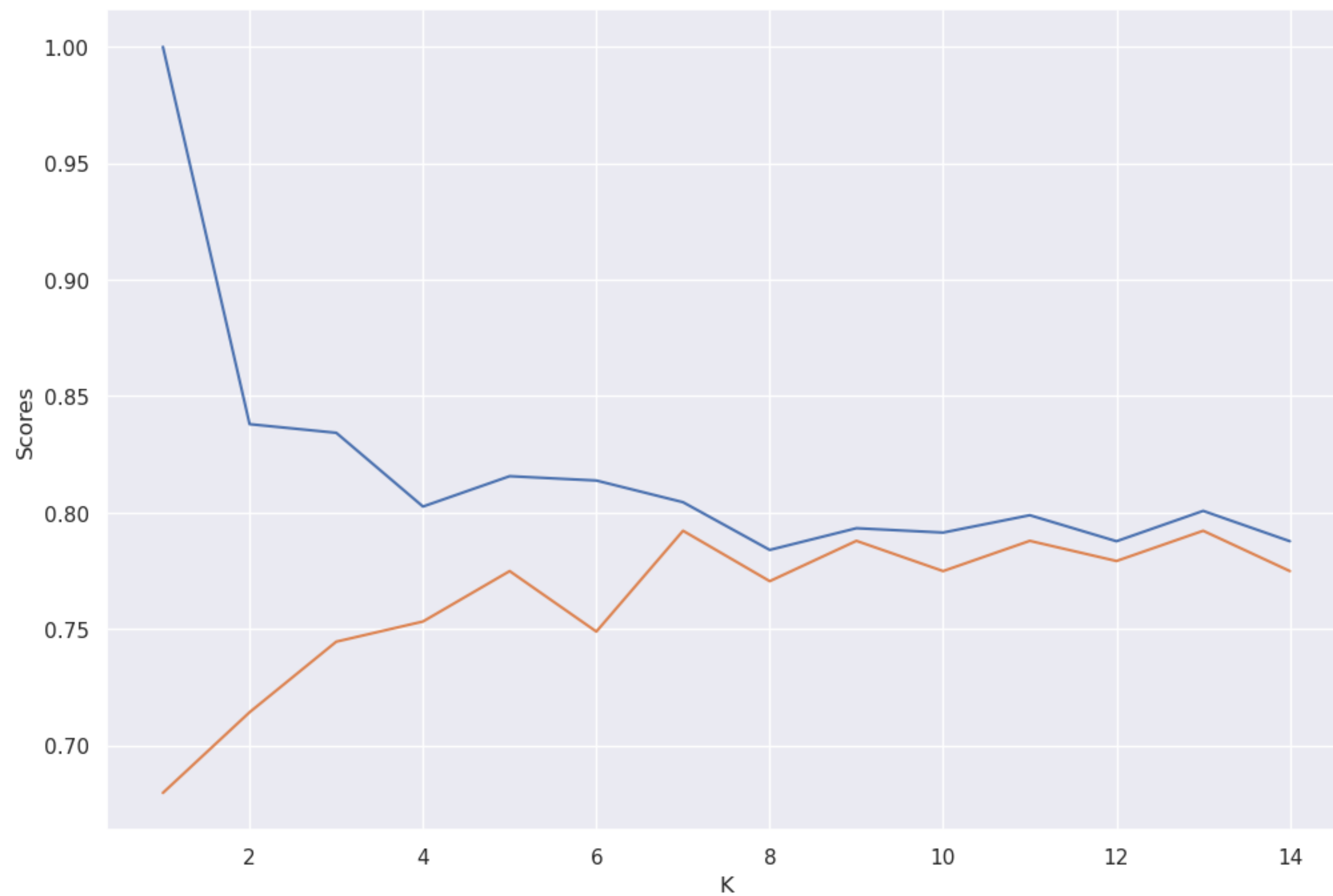| | Pregnancies | Glucose | BloodPressure | SkinThickness | Insulin | BMI | DiabetesPedigreeFunction | Age |
|---|---|---|---|---|---|---|---|---|
| 0 | 0.639947 | 0.865108 | -0.033518 | 0.670643 | -0.181541 | 0.166619 | 0.468492 | 1.425995 |
| 1 | -0.844885 | -1.206162 | -0.529859 | -0.012301 | -0.181541 | -0.852200 | -0.365061 | -0.190672 |
| 2 | 1.233880 | 2.015813 | -0.695306 | -0.012301 | -0.181541 | -1.332500 | 0.604397 | -0.105584 |
| 3 | -0.844885 | -1.074652 | -0.529859 | -0.695245 | -0.540642 | -0.633881 | -0.920763 | -1.041549 |
| 4 | -1.141852 | 0.503458 | -2.680669 | 0.670643 | 0.316566 | 1.549303 | 5.484909 | -0.020496 |

# Unsupervised Approach - KMeans

```python
kmeans_model = KMeans(init="random", n_clusters=2, n_init=10, max_iter=300, random_state=11)
y_pred = kmeans_model.fit_predict(diabetes)
```

```python
n_correct_predictions = 0
for i in range(diabetes.shape[0]):
    if diabetes["Outcome"][i] != y_pred[i]:
        n_correct_predictions += 1
print("Accuracy for kmeans:" +str(n_correct_predictions/diabetes.shape[0]))
```

```
Accuracy for kmeans:0.6536458333333334
```

# Supervised Approach - KNN

Confusion matrix

Knn(n_neighbors=7) ROC curve