

## Importing Libararies

```
In [1]: 1 import pandas as pd
        2 import numpy as np
        3 import seaborn as sns
        4 import matplotlib.pyplot as plt
        5 import warnings
        6 warnings.filterwarnings("ignore")
        7 from sklearn import metrics
```

## Retreiving Data

```
In [2]: 1 data=pd.read_csv(r"C:\DsTraining\five dataset for clening\diabetes.csv")
        2 data.head(10)
```

Out[2]:

	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
5	5	116	74	0	0	25.6	0.201	30	0
6	3	78	50	32	88	31.0	0.248	26	1
7	10	115	0	0	0	35.3	0.134	29	0
8	2	197	70	45	543	30.5	0.158	53	1
9	8	125	96	0	0	0.0	0.232	54	1

## Data Cleaning

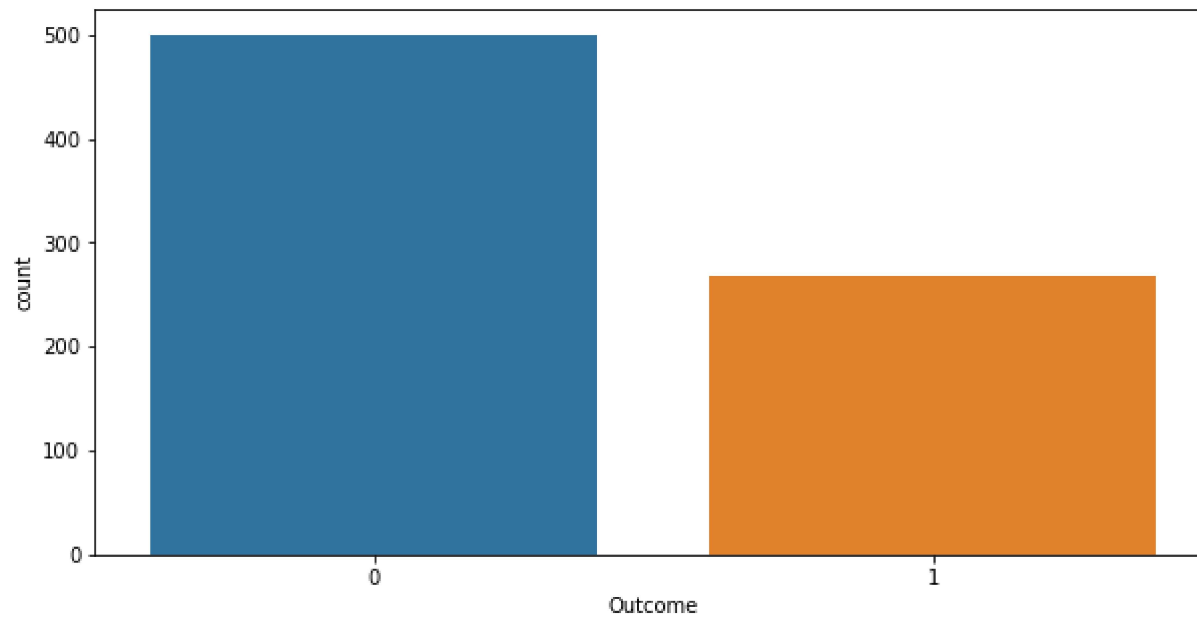
```
In [3]: 1 #find missing values and fill using appropriate central tendencies.(no missing values find in this data set
        2 data.isnull().sum()
```

```
Out[3]: Pregnancies      0
         Glucose          0
         BloodPressure    0
         SkinThickness    0
         Insulin          0
         BMI              0
         DiabetesPedigreeFunction  0
         Age              0
         Outcome          0
         dtype: int64
```

## EDA

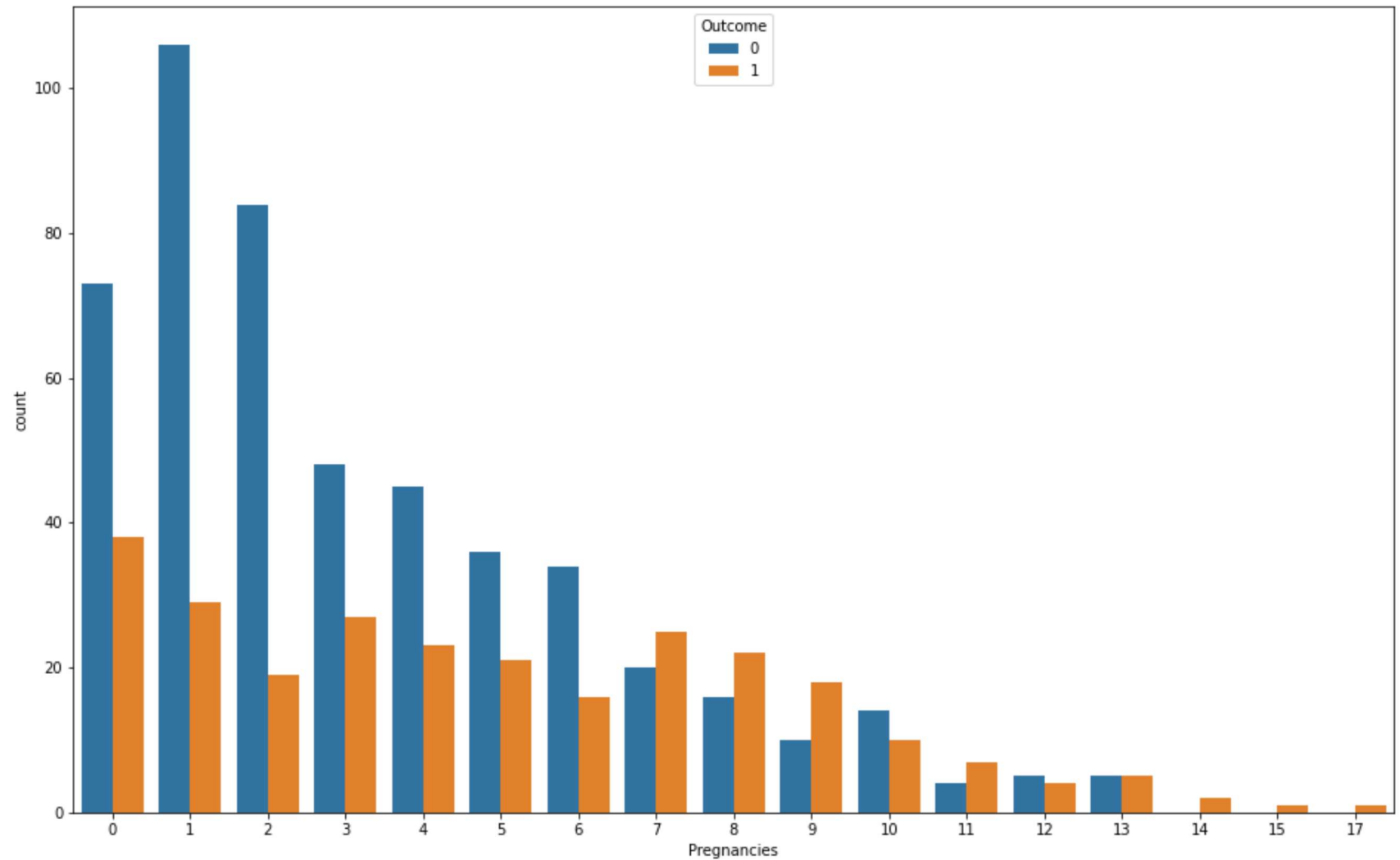
```
In [4]: 1 plt.figure(figsize=(10,5))
        2 sns.countplot(data["Outcome"])
```

```
Out[4]: <AxesSubplot:xlabel='Outcome', ylabel='count'>
```



```
In [5]: 1 plt.figure(figsize=(16,10))
        2 sns.countplot(data["Pregnancies"],hue="Outcome",data=data)
```

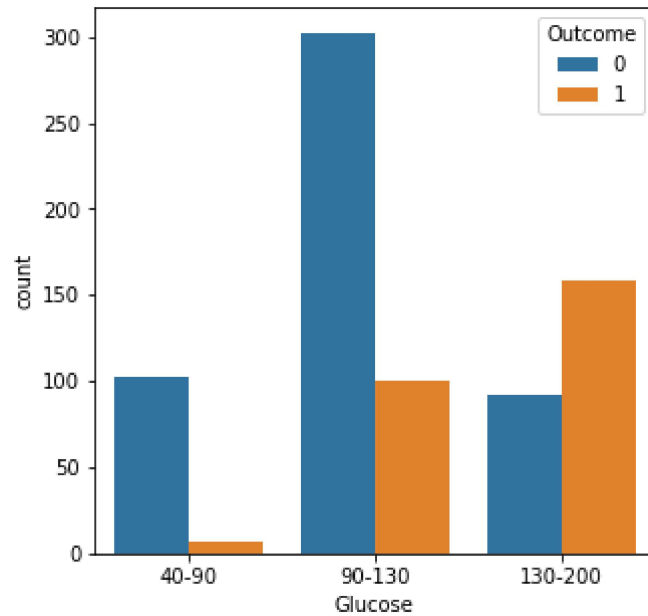
Out[5]: <AxesSubplot:xlabel='Pregnancies', ylabel='count'>



```
In [6]: 1 glucose_bins=pd.cut(data["Glucose"],bins=[40,90,130,200],labels=["40-90","90-130","130-200"])
```

```
In [7]: 1 plt.figure(figsize=(5,5))
        2 sns.countplot(glucose_bins,hue=data["Outcome"])
```

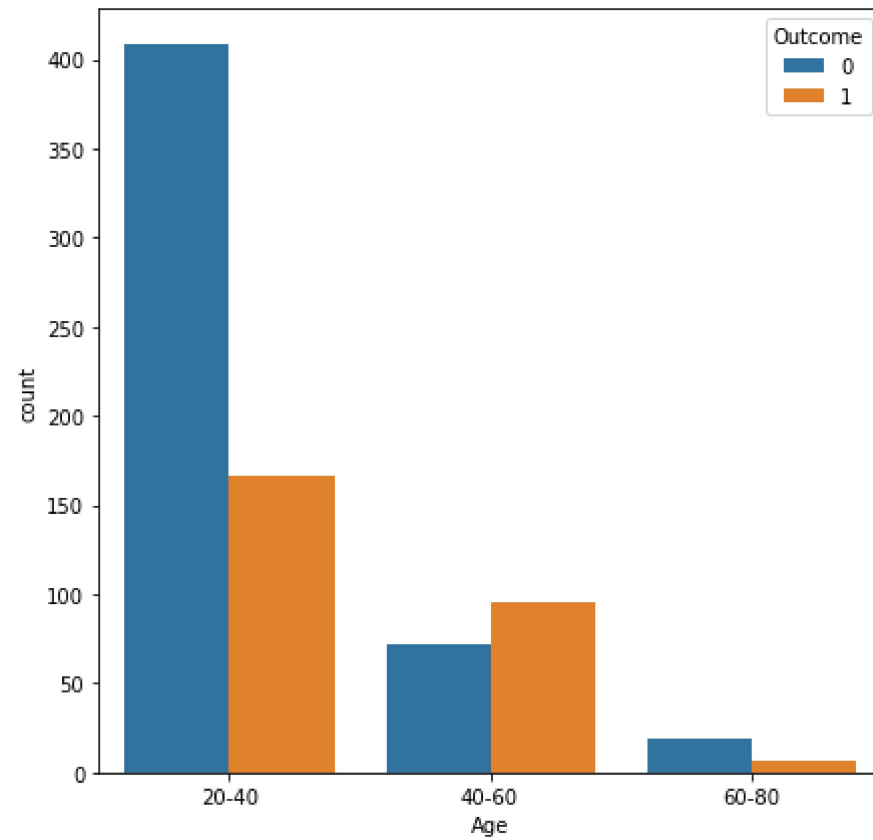
Out[7]: <AxesSubplot:xlabel='Glucose', ylabel='count'>



```
In [8]: 1 age_bins=pd.cut(data["Age"],bins=[20,40,60,80],labels=["20-40","40-60","60-80"])
```

```
In [9]: 1 plt.figure(figsize=(7,7))
        2 sns.countplot(age_bins,hue="Outcome",data=data)
```

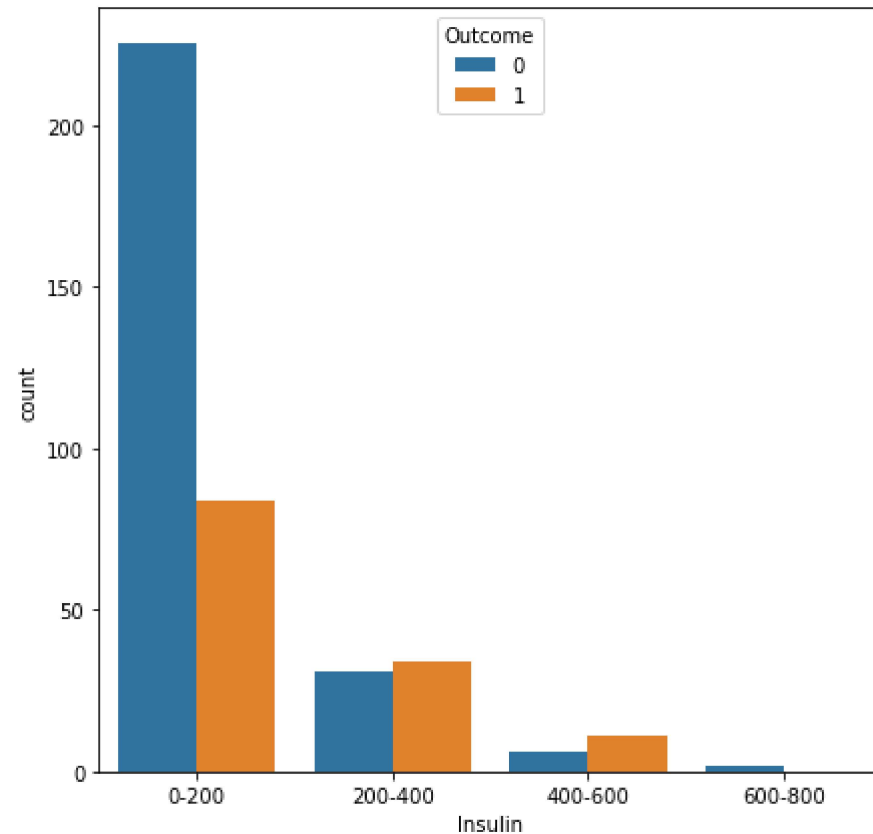
Out[9]: <AxesSubplot:xlabel='Age', ylabel='count'>



```
In [10]: 1 insulin_bins=pd.cut(data["Insulin"],bins=[0,200,400,600,800],labels=["0-200","200-400","400-600","600-800"])
```

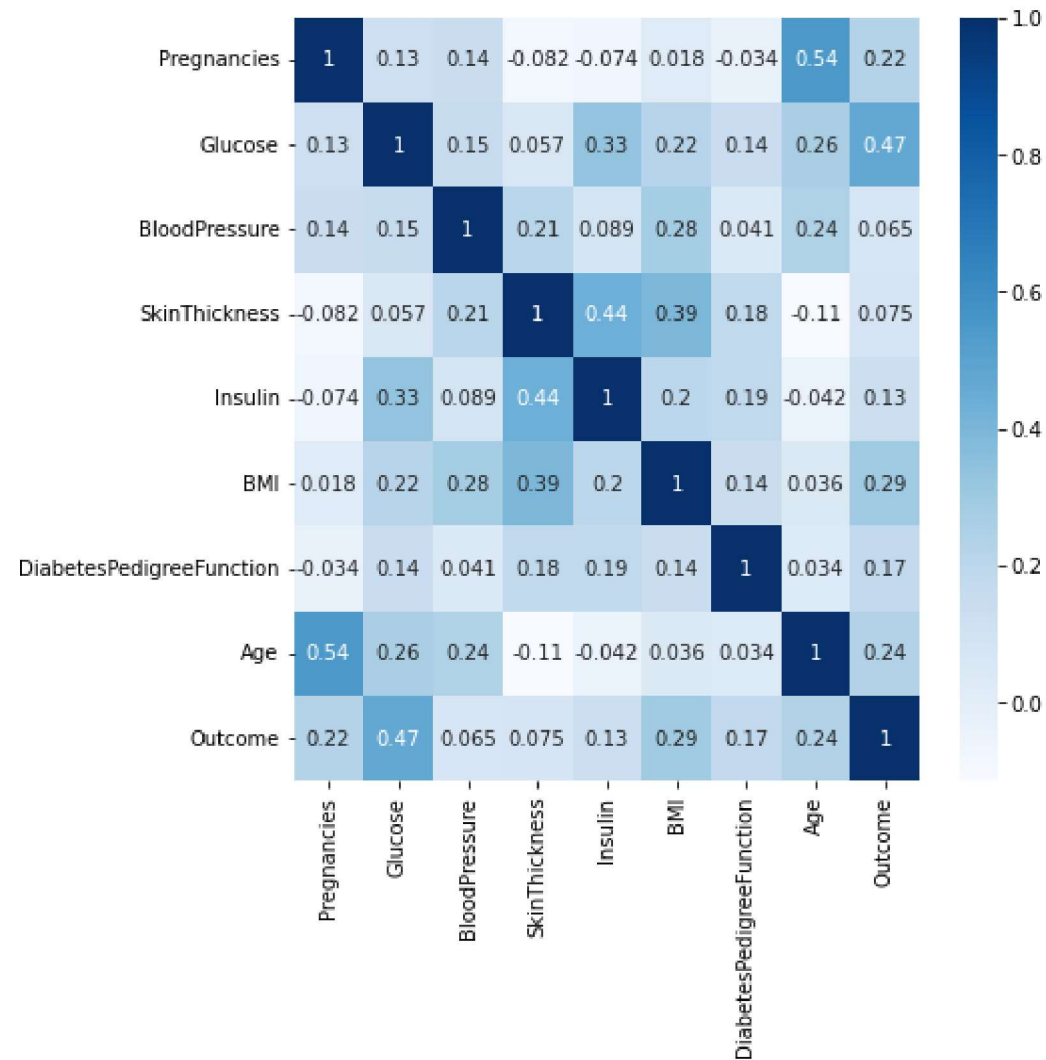
```
In [11]: 1 plt.figure(figsize=(7,7))  
2 sns.countplot(insulin_bins,hue=data["Outcome"])
```

```
Out[11]: <AxesSubplot:xlabel='Insulin', ylabel='count'>
```



```
In [16]: 1 plt.figure(figsize=(7,7))
          2 sns.heatmap(data.corr(),annot=True,cmap="Blues")
```

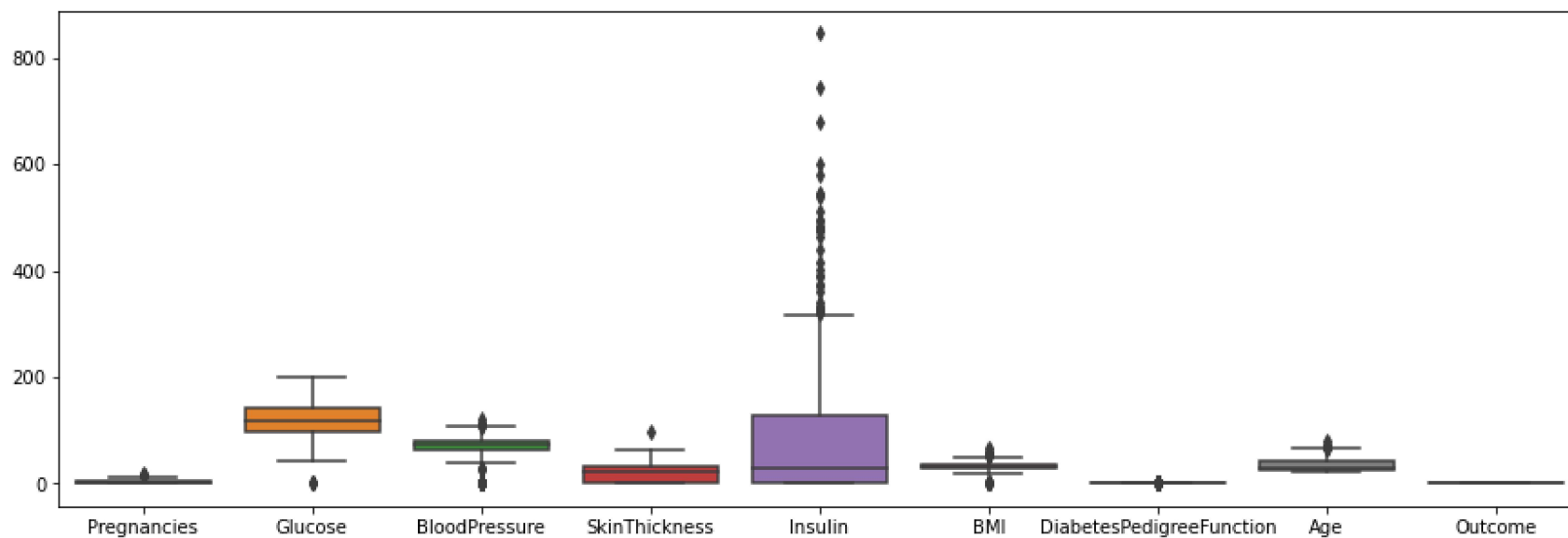
Out[16]: <AxesSubplot:>



## Outlier Detection

```
In [154]: 1 plt.figure(figsize=(15,5))  
          2 sns.boxplot(data=data)
```

Out[154]: <AxesSubplot:>



## Train-Test splitting.

```
In [155]: 1 x=data.iloc[:, :-1].values  
          2 y=data.iloc[:, -1].values
```

```
In [156]: 1 from sklearn.model_selection import train_test_split  
          2 x_train,x_test,y_train,y_test=train_test_split(x,y,random_state=42,test_size=0.20)
```



```
In [157]: 1 from sklearn.preprocessing import StandardScaler
2 sc=StandardScaler()
3 x_train=sc.fit_transform(x_train)
4 x_test=sc.transform(x_test)
```

## Logistic model

```
In [158]: 1 from sklearn.linear_model import LogisticRegression
2 classifier=LogisticRegression(random_state=0)
3 classifier.fit(x_train,y_train)
4
```

```
Out[158]: LogisticRegression(random_state=0)
```

```
In [159]: 1 y_pre=classifier.predict(x_test)
2 y_pre
```

```
Out[159]: array([0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 0, 1, 0, 0, 0, 0, 0, 0, 1, 1, 0, 0,
1, 0, 1, 1, 0, 0, 0, 0, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 0, 1, 1, 0,
0, 1, 1, 0, 0, 1, 0, 1, 1, 0, 0, 0, 1, 0, 0, 1, 1, 0, 0, 0, 0, 1,
0, 1, 0, 1, 1, 0, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 0, 1, 1, 0,
0, 0, 0, 0, 0, 1, 1, 1, 0, 0, 1, 0, 1, 0, 1, 0, 1, 0, 0, 1, 0, 1,
0, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 1, 1,
0, 0, 1, 0, 0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0],
dtype=int64)
```

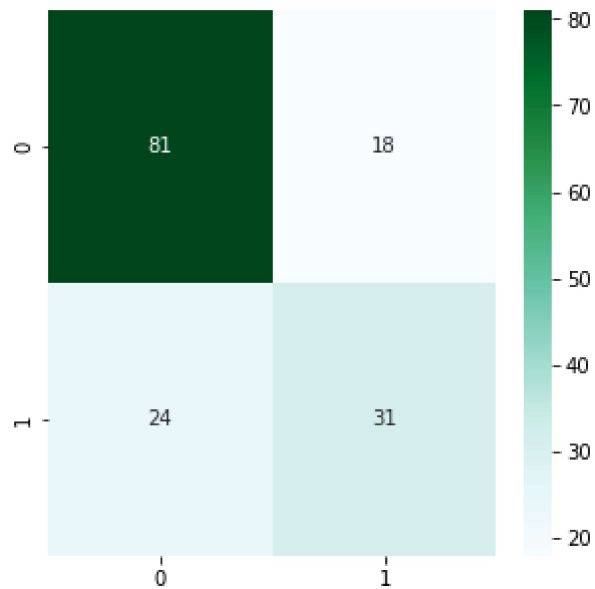
```
In [160]: 1 from sklearn.metrics import confusion_matrix,accuracy_score
2 cm=confusion_matrix(y_test,y_pred)
3 print(cm)
4 print("accuracy score is")
5 accuracy_score(y_test,y_pre)*100
6
```

```
[[81 18]
 [24 31]]
accuracy score is
```

```
Out[160]: 75.32467532467533
```

```
In [161]: 1 plt.figure(figsize=(5,5))
          2 sns.heatmap(cm,annot=True,cmap="BuGn")
```

Out[161]: <AxesSubplot:>



```
In [162]: 1 from sklearn.metrics import recall_score,precision_score,roc_auc_score
          2 rs1=recall_score(y_test,y_pre)*100
          3 print("Recall score=",rs)
          4 ps1=precision_score(y_test,y_pre)*100
          5 print("precision score=",ps)
          6 ras1=roc_auc_score(y_test,y_pre)*100
          7 print("ROC score=",ras1)
```

Recall score= 56.36363636363636  
precision score= 63.26530612244898  
ROC score= 73.53535353535354

## Random Forest

```
In [163]: 1 from sklearn.ensemble import RandomForestClassifier
          2 rfclassifier=RandomForestClassifier(n_estimators= 10, criterion="entropy")
          3 rfclassifier.fit(x_train,y_train)
```

```
Out[163]: RandomForestClassifier(criterion='entropy', n_estimators=10)
```

```
In [164]: 1 y_pred=rfclassifier.predict(x_test)
```

```
In [165]: 1 from sklearn.metrics import accuracy_score,recall_score,precision_score,confusion_matrix,roc_auc_score
          2 cm=confusion_matrix(y_test,y_pred)
          3 rs=recall_score(y_test,y_pred)*100
          4 ps=precision_score(y_test,y_pred)*100
          5 accs=accuracy_score(y_test,y_pred)*100
          6 ras=roc_auc_score(y_test,y_pred)*100
          7 print("confusion_matrix=\n",cm)
          8 print("recall_score=",rs)
          9 print("precision score=",ps)
         10 print("accuracy_score=",accs)
         11 print("ROC Score=",ras)
```

```
confusion_matrix=
```

```
[[85 14]
```

```
 [23 32]]
```

```
recall_score= 58.18181818181818
```

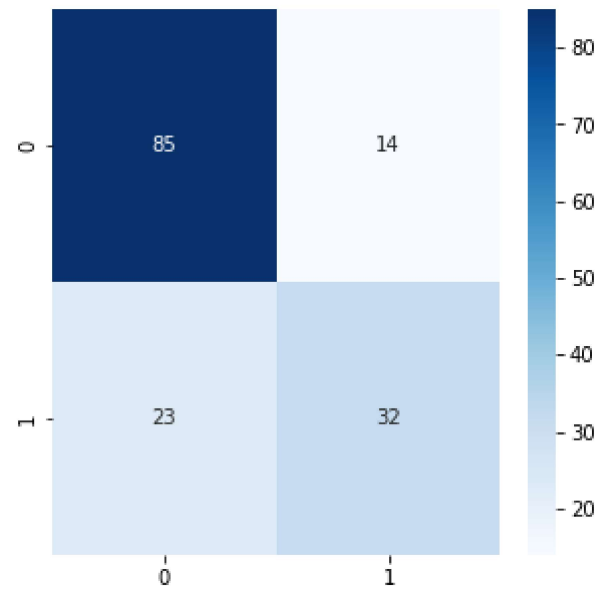
```
precision score= 69.56521739130434
```

```
accuracy_score= 75.97402597402598
```

```
ROC Score= 72.020202020202
```

```
In [166]: 1 plt.figure(figsize=(5,5))
          2 sns.heatmap(cm,cmap="Blues",annot=True)
```

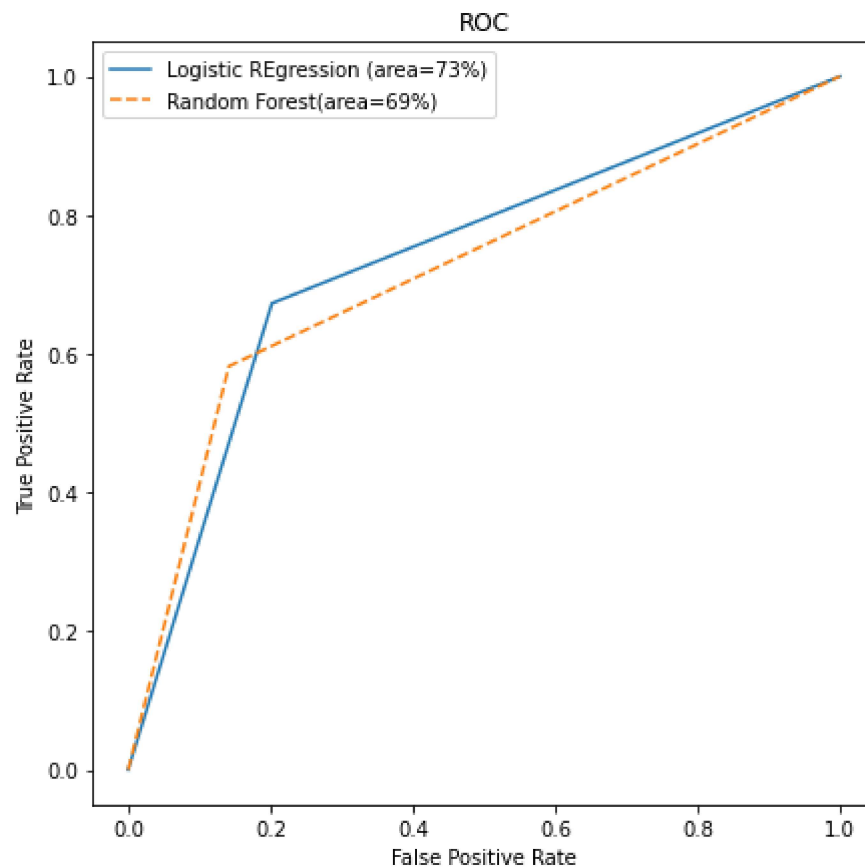
Out[166]: <AxesSubplot:>



## ROC Curve

```
In [167]: 1 fpr, tpr, _ = metrics.roc_curve(y_test, y_pre)
          2 fprd, tprd, _ = metrics.roc_curve(y_test, y_pred)
```

```
In [168]: 1 plt.figure(figsize=(7,7))
          2 plt.plot(fpr,tpr,label='Logistic REgression (area=73%)')
          3 plt.plot(fprd,tprd,"--",label="Random Forest(area=69%)")
          4 plt.ylabel('True Positive Rate')
          5 plt.xlabel('False Positive Rate')
          6 plt.title("ROC")
          7 plt.legend()
          8 plt.show()
```



```
In [ ]:
```

```
1
```

In [ ]:

1