

```
In [1]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
```

```
In [10]: df=pd.read_csv('diabetes.csv')
```

```
In [11]: df.head()
```

Out[11]:

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	Pedigree	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

```
In [12]: df.tail()
```

Out[12]:

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	Pedigree	Age	Outcome
763	10	101	76	48	180	32.9	0.171	63	0
764	2	122	70	27	0	36.8	0.340	27	0
765	5	121	72	23	112	26.2	0.245	30	0
766	1	126	60	0	0	30.1	0.349	47	1
767	1	93	70	31	0	30.4	0.315	23	0

```
In [13]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 768 entries, 0 to 767
Data columns (total 9 columns):
#   Column             Non-Null Count  Dtype  
---  -
0   Pregnancies        768 non-null   int64  
1   Glucose             768 non-null   int64  
2   BloodPressure       768 non-null   int64  
3   SkinThickness       768 non-null   int64  
4   Insulin             768 non-null   int64  
5   BMI                 768 non-null   float64 
6   Pedigree            768 non-null   float64 
7   Age                 768 non-null   int64  
8   Outcome             768 non-null   int64  
dtypes: float64(2), int64(7)
memory usage: 54.1 KB
```

```
In [14]: df.isnull().sum()
```

```
Out[14]: Pregnancies      0
Glucose      0
BloodPressure 0
SkinThickness 0
Insulin      0
BMI          0
Pedigree     0
Age         0
Outcome     0
dtype: int64
```


```
In [15]: df.shape
```

```
Out[15]: (768, 9)
```

```
In [16]: df.describe()
```

```
Out[16]:
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	Pedigree	
count	768.000000	768.000000	768.000000	768.000000	768.000000	768.000000	768.000000	76
mean	3.845052	120.894531	69.105469	20.536458	79.799479	31.992578	0.471876	3
std	3.369578	31.972618	19.355807	15.952218	115.244002	7.884160	0.331329	1
min	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.078000	2
25%	1.000000	99.000000	62.000000	0.000000	0.000000	27.300000	0.243750	2
50%	3.000000	117.000000	72.000000	23.000000	30.500000	32.000000	0.372500	2
75%	6.000000	140.250000	80.000000	32.000000	127.250000	36.600000	0.626250	4
max	17.000000	199.000000	122.000000	99.000000	846.000000	67.100000	2.420000	8



```
In [17]: df.head()
```

```
Out[17]:
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	Pedigree	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

```
In [18]: df.tail()
```

```
Out[18]:
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	Pedigree	Age	Outcome
763	10	101	76	48	180	32.9	0.171	63	0
764	2	122	70	27	0	36.8	0.340	27	0
765	5	121	72	23	112	26.2	0.245	30	0
766	1	126	60	0	0	30.1	0.349	47	1
767	1	93	70	31	0	30.4	0.315	23	0

In [19]: df.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 768 entries, 0 to 767
Data columns (total 9 columns):
#   Column                Non-Null Count  Dtype
---  -
0   Pregnancies            768 non-null   int64
1   Glucose                768 non-null   int64
2   BloodPressure          768 non-null   int64
3   SkinThickness          768 non-null   int64
4   Insulin                768 non-null   int64
5   BMI                   768 non-null   float64
6   Pedigree              768 non-null   float64
7   Age                   768 non-null   int64
8   Outcome               768 non-null   int64
dtypes: float64(2), int64(7)
memory usage: 54.1 KB
```

In [20]: df.isnull().sum()

```
Out[20]: Pregnancies    0
Glucose              0
BloodPressure        0
SkinThickness        0
Insulin              0
BMI                  0
Pedigree             0
Age                  0
Outcome              0
dtype: int64
```

In [21]: df.shape

Out[21]: (768, 9)

In [22]: df.describe()


```
Out[22]:
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	Pedigree	
count	768.000000	768.000000	768.000000	768.000000	768.000000	768.000000	768.000000	76
mean	3.845052	120.894531	69.105469	20.536458	79.799479	31.992578	0.471876	3
std	3.369578	31.972618	19.355807	15.952218	115.244002	7.884160	0.331329	1
min	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.078000	2
25%	1.000000	99.000000	62.000000	0.000000	0.000000	27.300000	0.243750	2
50%	3.000000	117.000000	72.000000	23.000000	30.500000	32.000000	0.372500	2
75%	6.000000	140.250000	80.000000	32.000000	127.250000	36.600000	0.626250	4
max	17.000000	199.000000	122.000000	99.000000	846.000000	67.100000	2.420000	8

In [23]: `df.corr()`

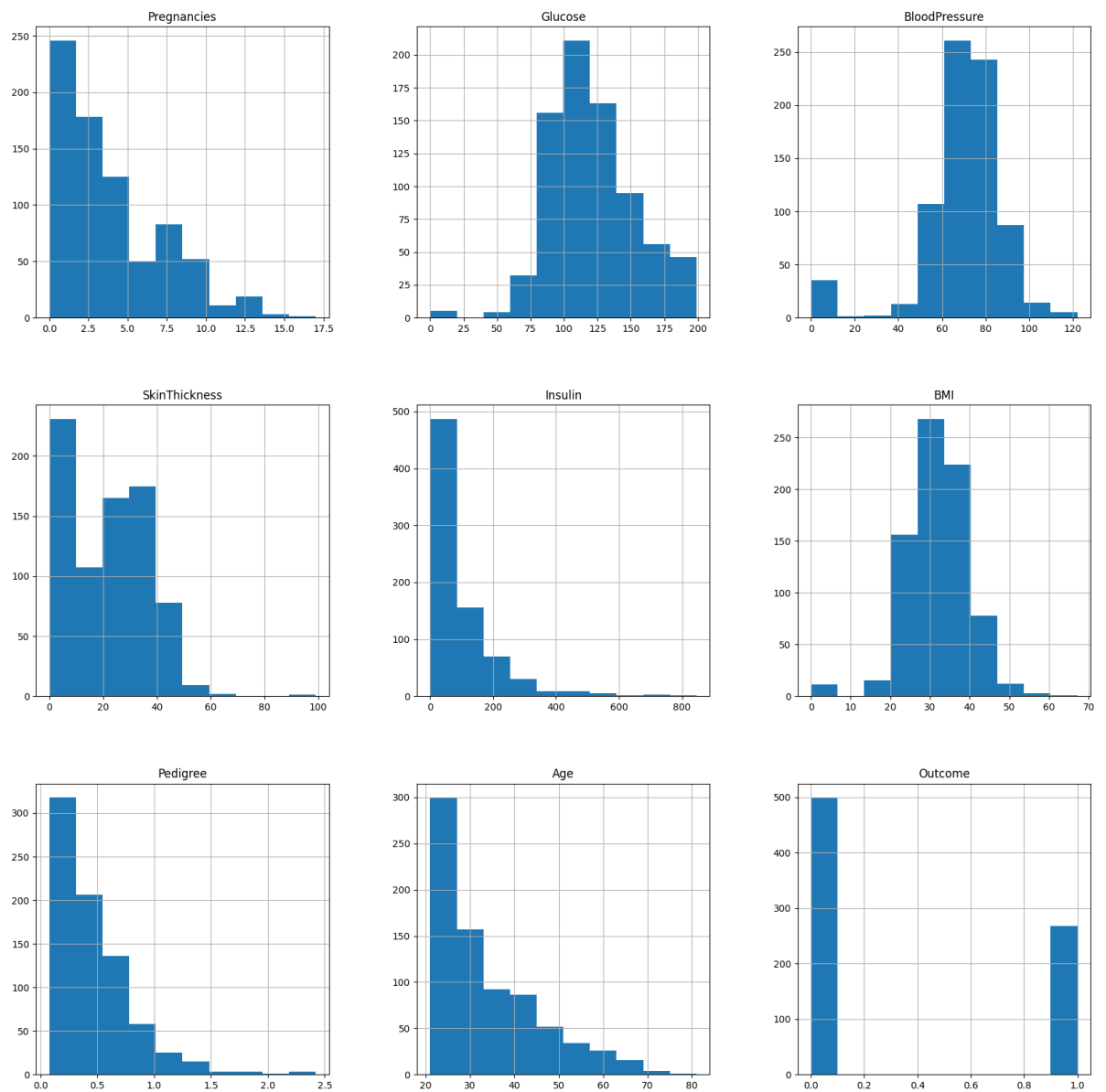
Out[23]:

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	Pedigree
Pregnancies	1.000000	0.129459	0.141282	-0.081672	-0.073535	0.017683	-0.033523
Glucose	0.129459	1.000000	0.152590	0.057328	0.331357	0.221071	0.137337
BloodPressure	0.141282	0.152590	1.000000	0.207371	0.088933	0.281805	0.041265
SkinThickness	-0.081672	0.057328	0.207371	1.000000	0.436783	0.392573	0.183928
Insulin	-0.073535	0.331357	0.088933	0.436783	1.000000	0.197859	0.185071
BMI	0.017683	0.221071	0.281805	0.392573	0.197859	1.000000	0.140647
Pedigree	-0.033523	0.137337	0.041265	0.183928	0.185071	0.140647	1.000000
Age	0.544341	0.263514	0.239528	-0.113970	-0.042163	0.036242	0.033561
Outcome	0.221898	0.466581	0.065068	0.074752	0.130548	0.292695	0.173844

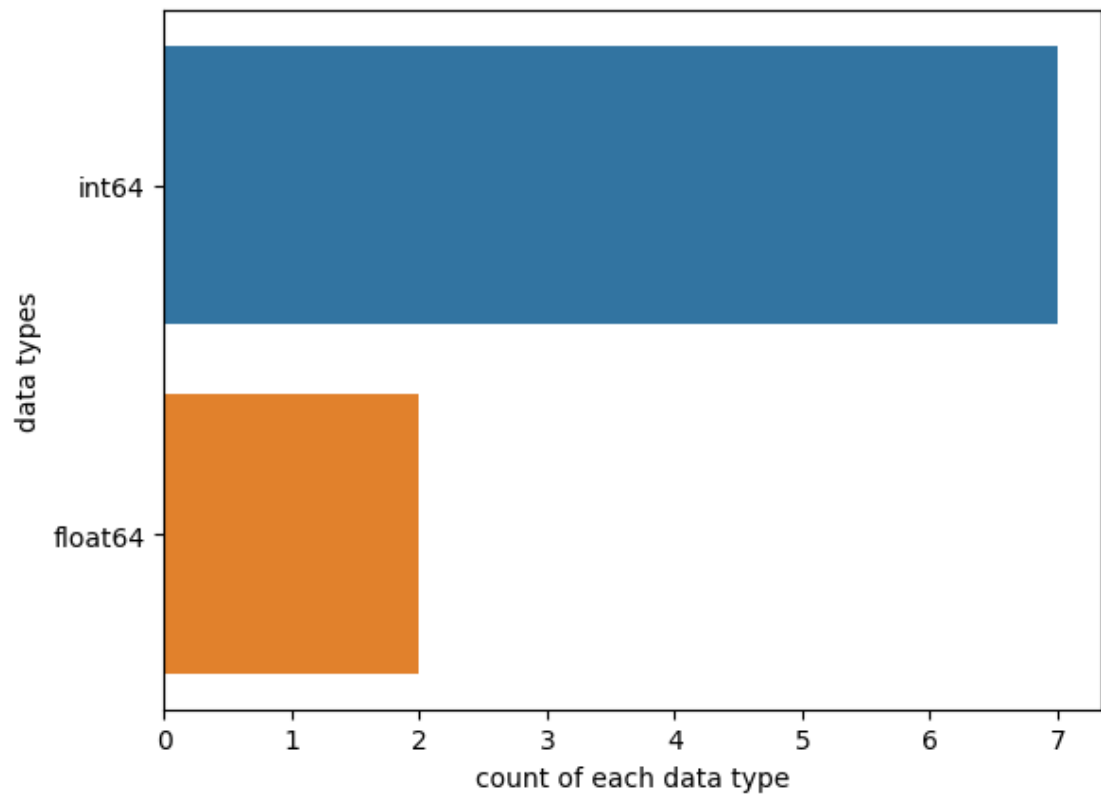


In [24]: `df['Glucose'].fillna(df['Glucose'].mean(), inplace = True)`
`df['BloodPressure'].fillna(df['BloodPressure'].mean(), inplace = True)`
`df['SkinThickness'].fillna(df['SkinThickness'].median(), inplace = True)`
`df['Insulin'].fillna(df['Insulin'].median(), inplace = True)`
`df['BMI'].fillna(df['BMI'].median(), inplace = True)`

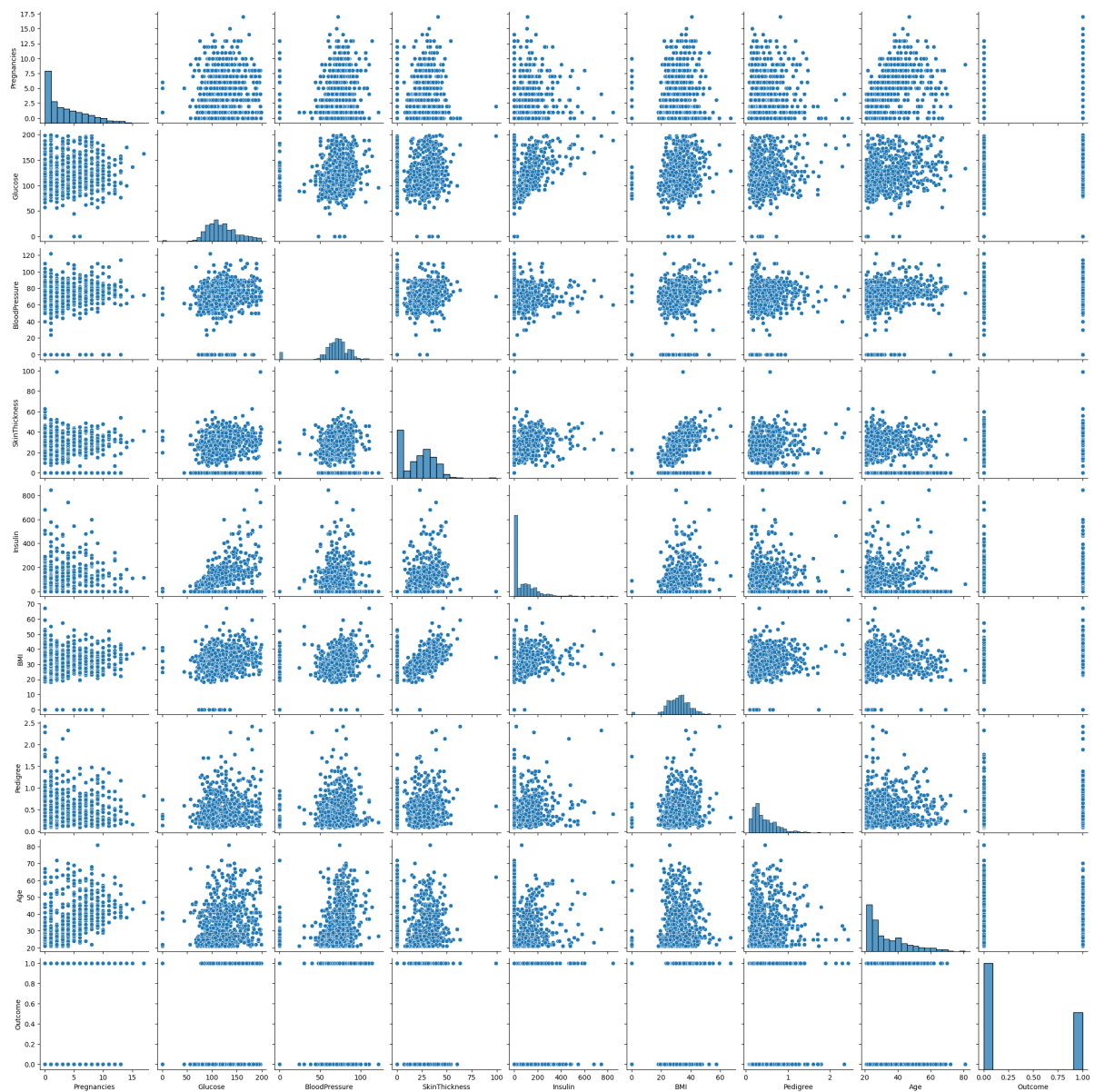
```
In [25]: p = df.hist(figsize = (20,20))
```



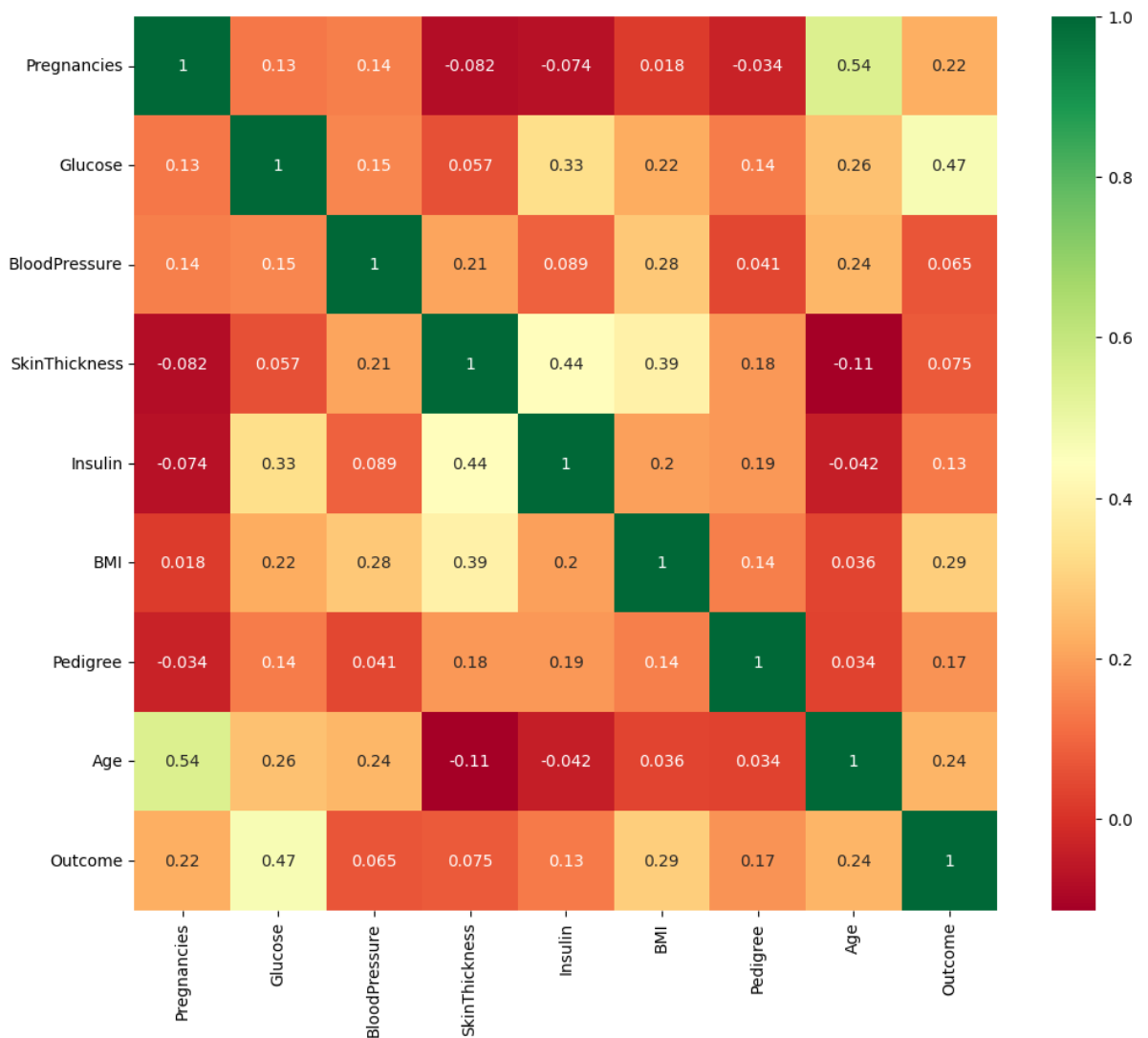
```
In [26]: sns.countplot(y=df.dtypes ,data=df)
plt.xlabel("count of each data type")
plt.ylabel("data types")
plt.show()
```



```
In [27]: p=sns.pairplot(df)
```



```
In [28]: plt.figure(figsize=(12,10)) # on this line I just set the size of figure to 12 by 10
p=sns.heatmap(df.corr(), annot=True,cmap = 'RdYlGn')
```



```
In [29]: from sklearn.preprocessing import StandardScaler
sc_X = StandardScaler()
X = pd.DataFrame(sc_X.fit_transform(df.drop(["Outcome"],axis = 1)),
columns=['Pregnancies', 'Glucose', 'BloodPressure', 'SkinThickness', 'Insulin',
'BMI', 'DiabetesPedigreeFunction', 'Age'])
```

```
In [30]: y = df.Outcome
```

```
In [31]: X.head()
```

```
Out[31]:
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
0	0.639947	0.848324	0.149641	0.907270	-0.692891	0.204013	0.468496	33	0
1	-0.844885	-1.123396	-0.160546	0.530902	-0.692891	-0.684422	-0.365000	23	1
2	1.233880	1.943724	-0.263941	-1.288212	-0.692891	-1.103255	0.604396	23	0
3	-0.844885	-0.998208	-0.160546	0.154533	0.123302	-0.494043	-0.920760	23	1
4	-1.141852	0.504055	-1.504687	0.907270	0.765836	1.409746	5.484900	23	0


```
In [32]: from sklearn.model_selection import train_test_split

# Split the data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=1/3, random_sta
```

```
In [33]: from sklearn.neighbors import KNeighborsClassifier
test_scores = []
train_scores = []
for i in range(1,15):
    knn = KNeighborsClassifier(i)
    knn.fit(X_train,y_train)

    train_scores.append(knn.score(X_train,y_train))
    test_scores.append(knn.score(X_test,y_test))
```

```
In [36]: max_train_score = max(train_scores)
train_scores_ind = [i for i, v in enumerate(train_scores) if v == max_train_score]
k_values_for_max_score = list(map(lambda x: x + 1, train_scores_ind))

print('Max train score {:.2f}% and k = {}'.format(max_train_score * 100, k_values_f
```

Max train score 100.00% and k = [1]

```
In [38]: max_test_score = max(test_scores)
test_scores_ind = [i for i, v in enumerate(test_scores) if v == max_test_score]
k_values_for_max_test_score = list(map(lambda x: x + 1, test_scores_ind))

print('Max test score {:.2f}% and k = {}'.format(max_test_score * 100, k_values_for
```

Max test score 76.95% and k = [11]

```
In [39]: knn = KNeighborsClassifier(11)
knn.fit(X_train,y_train)
knn.score(X_test,y_test)
```

Out[39]: 0.76953125

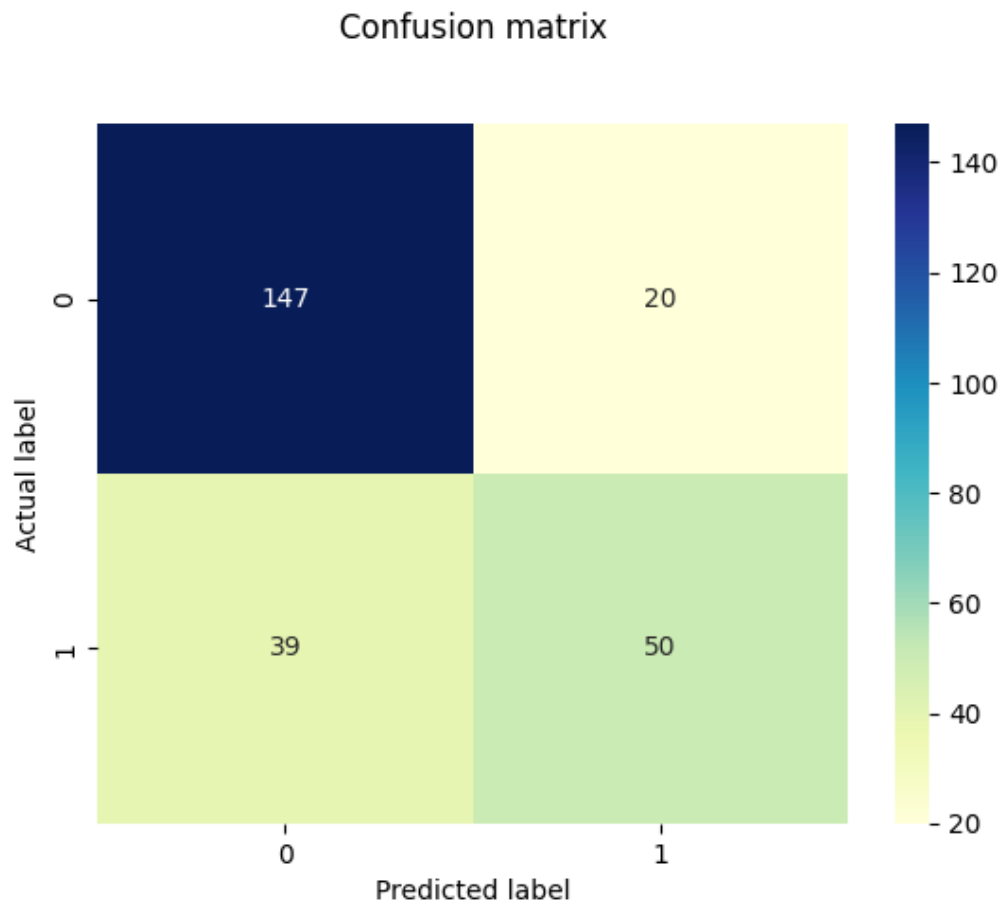
```
In [40]: from sklearn.metrics import confusion_matrix
#let us get the predictions using the classifier we had fit above
y_pred = knn.predict(X_test)
confusion_matrix(y_test,y_pred)
pd.crosstab(y_test, y_pred, rownames=['True'], colnames=['Predicted'], margins=True
```

Out[40]:

Predicted	0	1	All
True			
0	147	20	167
1	39	50	89
All	186	70	256

```
In [41]: y_pred = knn.predict(X_test)
from sklearn import metrics
cnf_matrix = metrics.confusion_matrix(y_test, y_pred)
p = sns.heatmap(pd.DataFrame(cnf_matrix), annot=True, cmap="YlGnBu" ,fmt='g')
plt.title('Confusion matrix', y=1.1)
plt.ylabel('Actual label')
plt.xlabel('Predicted label')
```

Out[41]: Text(0.5, 23.52222222222222, 'Predicted label')



```
In [42]: from sklearn.metrics import classification_report
print(classification_report(y_test,y_pred))
```

	precision	recall	f1-score	support
0	0.79	0.88	0.83	167
1	0.71	0.56	0.63	89
accuracy			0.77	256
macro avg	0.75	0.72	0.73	256
weighted avg	0.76	0.77	0.76	256

In []: