

```
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
import warnings
warnings.filterwarnings('ignore')

In [2]: df=pd.read_csv(r'C:\Users\sharma17\Downloads\data\classification dataset\diabetes.csv')

In [3]: df.head()

Out[3]:
   pregnant  glucose  bp  skin  insulin  bmi  predgree  age  target
0          0         126   72    35      0  33.6  0.627  50      0
1          1         85   66   29      0  26.6  0.351  31      0
2          8        183   84    0      0  23.3  0.672  32      1
3          1         89   66   23     94  26.1  0.167  21      0
4          0        137   40   35    168  43.1  2.286  33      1

In [4]: sns.heatmap(df.isnull())
plt.show()

In [5]: df.describe()

Out[5]:
   pregnant  glucose  bp  skin  insulin  bmi  predgree  age  target
count  768.000000  768.000000  768.000000  768.000000  768.000000  768.000000  768.000000  768.000000  768.000000
mean     3.845052  120.984531  69.106469  20.536456  79.799479  31.993576  0.471876  33.242885  0.349956
std     3.369578  31.872618  19.359807  15.952218  115.244002  7.884160  0.331329  11.760232  0.476951
min     0.000000  0.000000  0.000000  0.000000  0.000000  0.000000  0.078000  21.000000  0.000000
25%     1.000000  99.000000  62.000000  0.000000  0.000000  27.300000  0.247500  24.000000  0.000000
50%     3.000000  117.000000  72.000000  23.000000  30.000000  32.000000  0.372500  29.000000  0.000000
75%     6.000000  140.250000  80.000000  32.000000  127.250000  36.600000  0.626250  41.000000  1.000000
max    17.000000  199.000000  122.000000  99.000000  846.000000  67.100000  2.420000  81.000000  1.000000

In [6]: df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 768 entries, 0 to 767
Data columns (total 9 columns):
 #  Column  Non-Null Count  Dtype
---  --
 0  pregnant  768 non-null      int64
 1  glucose  768 non-null      int64
 2  bp       768 non-null      int64
 3  skin     768 non-null      int64
 4  insulin  768 non-null      float64
 5  bmi      768 non-null      float64
 6  predgree 768 non-null      float64
 7  age      768 non-null      int64
 8  target   768 non-null      int64
dtypes: float64(2), int64(7)
memory usage: 54.1 KB

In [7]: sns.pairplot(df,hue='target')

Out[7]:
<seaborn.axisgrid.PairGrid at 8x19f50c99798>

In [8]: df.target.value_counts()

Out[8]:
0    589
1    268
Name: target, dtype: int64

In [9]: plt.figure(figsize=(10,12))
sns.heatmap(df.corr(),annot=True)
plt.show()

In [10]: from scipy.stats import skew

In [11]: for col in df:
    print(col, '---', skew(df[col]))
    print(skew(df[col]))
    sns.distplot(df[col])
    plt.show()

In [12]: plt.figure(figsize=(15,12))
sns.boxplot(data=df)
plt.show()

In [13]: x=df.iloc[:,0:1].values
y=df.iloc[:,1].values

In [14]: from sklearn.linear_model import LogisticRegression
from sklearn.tree import DecisionTreeClassifier
from sklearn.svm import SVC
from sklearn.neighbors import KNeighborsClassifier
from sklearn.ensemble import RandomForestClassifier, AdaBoostClassifier, GradientBoostingClassifier
from xgboost import XGBClassifier
from sklearn.model_selection import train_test_split, StratifiedKFold, cross_val_score
from sklearn.metrics import accuracy_score, confusion_matrix, classification_report

In [15]: models=[(Logistic Regression,LogisticRegression()),
                 (Decision Tree,DecisionTreeClassifier()),
                 (Support Vector ,SVC()),
                 (Karest Neighbors ,KNeighborsClassifier()),
                 (Random Forest , RandomForestClassifier()),
                 (AdaBoostClassifier, AdaBoostClassifier()),
                 (Gradient Boosting,GradientBoostingClassifier()),
                 (Xtream Gradient Boosting,XGBClassifier())
                ]

In [16]: xtrain,xtest,ytrain,ytest=train_test_split(x,y,random_state=0,test_size=0.25,stratify=y)

In [17]: accuracy=[]

for name,model in models:
    print(name,' :-- ')
    print()
    model.fit(xtrain,ytrain)
    print()
    ypred=model.predict(xtest)
    print(confusion_matrix(ytest,ypred))
    print()
    print(classification_report(ytest,ypred))
    accuracy.append(accuracy_score(ytest,ypred))

Logistic Regression :--

[[118 15]
 [ 27 40]]

      precision  recall  f1-score  support
0      0.88      0.88      0.84     125
1      0.72      0.66      0.68      67
accuracy      0.77      0.74      0.75     192
macro avg      0.77      0.74      0.75     192
weighted avg      0.78      0.76      0.75     192

Decision Tree :--

[[184 21]
 [ 27 40]]

      precision  recall  f1-score  support
0      0.79      0.83      0.81     125
1      0.66      0.66      0.62      67
accuracy      0.72      0.71      0.72     192
macro avg      0.72      0.71      0.72     192
weighted avg      0.75      0.75      0.75     192

Support Vector :--

[[113 12]
 [ 33 34]]

      precision  recall  f1-score  support
0      0.77      0.98      0.83     125
1      1.00      0.74      0.83      67
accuracy      0.76      0.71      0.72     192
macro avg      0.76      0.71      0.72     192
weighted avg      0.76      0.77      0.75     192

Nearest Neighbors :--

[[184 21]
 [ 27 41]]

      precision  recall  f1-score  support
0      0.88      0.83      0.82     125
1      0.66      0.61      0.64      67
accuracy      0.73      0.72      0.75     192
macro avg      0.75      0.76      0.75     192
weighted avg      0.75      0.76      0.75     192

Random Forest :--

[[118 15]
 [ 28 39]]

      precision  recall  f1-score  support
0      0.88      0.88      0.84     125
1      1.00      0.58      0.63      67
accuracy      0.76      0.71      0.72     192
macro avg      0.77      0.74      0.74     192
weighted avg      0.77      0.78      0.77     192

AdaBoostClassifier :--

[[188 17]
 [ 28 39]]

      precision  recall  f1-score  support
0      0.79      0.86      0.83     125
1      0.78      0.58      0.63      67
accuracy      0.75      0.72      0.75     192
macro avg      0.75      0.76      0.76     192
weighted avg      0.76      0.77      0.76     192

Gradient Boosting :--

[[188 17]
 [ 22 45]]

      precision  recall  f1-score  support
0      0.83      0.86      0.85     125
1      1.00      0.67      0.78      67
accuracy      0.78      0.77      0.77     192
macro avg      0.79      0.76      0.78     192
weighted avg      0.78      0.79      0.79     192

Xtreme Gradient Boosting :--

[[136 17]
 [ 22 45]]
WARNING: C:\Users\Administrator\workspace\yghost-win64-release-1.5.1\src\learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'er
ror' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.

[[186 19]
 [ 22 45]]

      precision  recall  f1-score  support
0      0.83      0.85      0.84     125
1      0.78      0.67      0.69      67
accuracy      0.77      0.76      0.76     192
macro avg      0.77      0.76      0.76     192
weighted avg      0.78      0.79      0.79     192

In [18]: print('the mean accuracy of the model :-- ',np.array(accuracy).mean()*100)

the mean accuracy of the model :-- 77.233616666667

In [19]: xgb=GradientBoostingClassifier(ccp.alpha=0.088)
xgb.fit(xtrain,ytrain)
ypred=xgb.predict(xtest)
print('training score :-- ',xgb.score(xtrain,ytrain)*100)
print('testing score :-- ',xgb.score(xtest,ytest)*100)

training score :-- 74.8241363846346
testing score :-- 73.8562991583206

In [20]: skf=StratifiedKFold(n_splits=5)
skf.get_n_splits(x,y)

Out[21]: 5

In [22]: accuracy=[]

for train,test in skf.split(x,y):
    xtrain,xtest=ytrain[ytest]
    ytrain,ytest=ytrain[ytest]

    xgb=GradientBoostingClassifier(ccp.alpha=0.888)
    xgb.fit(xtrain,ytrain)
    ypred=xgb.predict(xtest)
    accuracy.append(accuracy_score(ytest,ypred))

In [23]: print('this is the minimum accuracy we can get from xtream boosting classifier :-- ',np.array(accuracy).mean()*100)

this is the minimum accuracy we can get from xtream boosting classifier :-- 74.34937614882

Because the dataset is imbalance we are getting less accuracy

In [24]: df.target.value_counts()

Out[24]:
0    589
1    268
Name: target, dtype: int64

In [25]: df.corr()

Out[25]:
   pregnant  glucose  bp  skin  insulin  bmi  predgree  age  target
pregnant  1.000000  0.150459  0.141282  0.061872  0.073525  0.017683  0.032523  0.544341  0.321098
glucose    0.149959  1.000000  0.152590  0.067328  0.319197  0.137337  0.263514  0.466581  0.650681
bp         0.141282  0.152590  1.000000  0.207371  0.088933  0.281805  0.041265  0.295528  0.466581
skin       0.061872  0.067328  0.207371  1.000000  0.436763  0.192573  0.189828  0.113970  0.074752
insulin    0.073525  0.319197  0.088933  0.436763  1.000000  0.197859  0.185071  0.042163  0.130548
bmi        0.017683  0.221071  0.281805  0.192573  0.197859  1.000000  0.140647  0.036242  0.282895
predgree   0.032523  0.137337  0.041265  0.189828  0.185071  0.140647  1.000000  0.033961  0.173844
age        0.544341  0.466581  0.295528  0.113970  0.042163  0.036242  0.033961  1.000000  0.238594
target     0.321098  0.466581  0.466581  0.074752  0.130548  0.282895  0.173844  0.238594  1.000000

In [ ]:
```