

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

In [2]: df=pd.read_csv('C:\Users\shane17\Downloads\1\classification dataset\HeartDisease.csv')

In [3]: df.head()

Out[3]:
   age  sex  BP  cholesterol  heart_disease
0    70    1    130    222             1
1    67    0    115    164             0
2    57    1    124    201             1
3    64    1    128    203             0
4    74    0    120    209             0

In [4]: df.info()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 270 entries, 0 to 269
Data columns (total 5 columns):
 #   Column        Non-Null Count  Dtype
---  --
 0   age           270 non-null      int64
 1   sex           270 non-null      int64
 2   BP            270 non-null      int64
 3   cholesterol   270 non-null      int64
 4   heart_disease 270 non-null      int64
dtypes: int64(5)
memory usage: 10.7 KB

In [5]: df.describe()

Out[5]:
      age      sex      BP      cholesterol      heart_disease
count 270.000000 270.000000 270.000000 270.000000 270.000000
mean  54.433333   0.777778  131.344444  249.692593   0.444444
std    9.109067   0.408196  17.861008   51.085237   0.497627
min     29.000000   0.000000   94.000000  126.000000   0.000000
25%    48.000000   0.000000  120.000000  213.000000   0.000000
50%    55.000000   1.000000  130.000000  245.000000   0.000000
75%    61.000000   1.000000  140.000000  280.000000   1.000000
max     77.000000   1.000000  200.000000  564.000000   1.000000

In [6]: sns.pairplot(df)

Out[6]: <seaborn.axisgrid.PairGrid at 0x203903d0e9>


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

Out[7]: <seaborn.axisgrid.PairGrid at 0x20392041139>


In [8]: plt.figure(figsize=(8,8))
sns.heatmap(df.corr(),annot=True)
plt.show()

Out[8]:


In [9]: plt.figure(figsize=(12,10))
sns.boxplot(data=df)
plt.xticks(rotation=90)
plt.show()

Out[9]:


In [10]: df.skew()

Out[10]:
age          -0.363635
sex          -0.760884
BP           0.722628
cholesterol  1.138272
heart_disease 0.224858
dtype: float64

In [11]: from scipy.stats import skew

In [12]: for col in df:
    print(col, ':-', skew(df[col]))
    print(skew(df[col]))
    sns.distplot(df[col])
    plt.show()

age :- -0.36270484746629945

sex :- -0.7608272549847587

BP :- 0.718597251657525

cholesterol :- 1.1771344865188054

heart_disease :- 0.22388679774997907


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

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

In [15]: xtrain,xtest,ytrain,ytest=train_test_split(x,y,random_state=0,test_size=0.25,stratify=y)

In [16]: models=[('Logistic Regression', LogisticRegression()),
              ('KNeighbors Classifier', KNeighborsClassifier()),
              ('Decision Tree', DecisionTreeClassifier()),
              ('Support vector Classifier', SVC()),
              ('Random Forest Classifier', RandomForestClassifier()),
              ('Ada Boost Classifier', AdaBoostClassifier()),
              ('Gradient Boosting Classifier', GradientBoostingClassifier()),
              ('Xtreme Boosting Classifier', XGBClassifier())
            ]

In [17]: accuracy=[]
for name,model in models:
    print(name)
    model.fit(xtrain,ytrain)
    ypred=model.predict(xtest)
    print('Confusion matrix:- ')
    print(confusion_matrix(ytest,ypred))
    print(classification_report(ytest,ypred))
    accuracy.append(accuracy_score(ytest,ypred))

Logistic Regression
Confusion matrix :-
[[20 13]
 [10 20]]
precision recall f1-score support
0 0.71 0.66 0.68 38
1 0.61 0.67 0.63 30
accuracy 0.66
macro avg 0.66
weighted avg 0.67
KNeighbors Classifier
Confusion matrix :-
[[27 11]
 [12 11]]
precision recall f1-score support
0 0.68 0.71 0.69 38
1 0.61 0.67 0.64 30
accuracy 0.64
macro avg 0.64
weighted avg 0.65
Decision Tree
Confusion matrix :-
[[20 12]
 [12 18]]
precision recall f1-score support
0 0.68 0.68 0.68 38
1 0.60 0.60 0.60 30
accuracy 0.64
macro avg 0.64
weighted avg 0.65
Support vector Classifier :-
Confusion matrix :-
[[33 9]
 [20 11]]
precision recall f1-score support
0 0.62 0.67 0.73 38
1 0.67 0.63 0.64 30
accuracy 0.63
macro avg 0.64
weighted avg 0.64
Random Forest Classifier :-
Confusion matrix :-
[[20 12]
 [10 20]]
precision recall f1-score support
0 0.72 0.68 0.70 38
1 0.62 0.67 0.65 30
accuracy 0.67
macro avg 0.67
weighted avg 0.68
Ada Boosting Classifier :-
Confusion matrix :-
[[22 16]
 [9 22]]
precision recall f1-score support
0 0.71 0.58 0.64 38
1 0.67 0.70 0.68 30
accuracy 0.64
macro avg 0.64
weighted avg 0.65
Gradient Boosting Classifier :-
Confusion matrix :-
[[24 14]
 [11 19]]
precision recall f1-score support
0 0.69 0.63 0.66 38
1 0.58 0.63 0.60 30
accuracy 0.63
macro avg 0.63
weighted avg 0.64
Xtreme Boosting Classifier :-
Confusion matrix :-
[[25 13]
 [10 20]]
precision recall f1-score support
0 0.71 0.66 0.68 38
1 0.63 0.67 0.65 30
accuracy 0.66
macro avg 0.66
weighted avg 0.67

In [18]: dtree=DecisionTreeClassifier()
dtree.fit(xtrain,ytrain)
ypred=dtree.predict(xtest)

In [19]: print(classification_report(ytest,ypred))

precision recall f1-score support
0 0.68 0.66 0.67 38
1 0.58 0.60 0.59 30
accuracy 0.63
macro avg 0.63
weighted avg 0.63

In [20]: print('training score :- ',dtree.score(xtrain,ytrain)*100)
print('testing score :- ',dtree.score(xtest,ytest)*100)

training score :- 189.0
testing score :- 63.23529411764706

In [21]: path=tree.cost_complexity_pruning_path(xtrain,ytrain)

Out[21]: ('ccp_alphas': array([0. , 0.00412941, 0.00448044, 0.00459903, 0.0049988 ,
0.0050905 , 0.00577658, 0.00594929, 0.00594929, 0.00615012,
0.00621465, 0.00660966, 0.00660966, 0.00660966, 0.00786454,
0.00825403, 0.00860966, 0.00877617, 0.00923348, 0.00960966,
0.01019977, 0.01155216, 0.01211384, 0.01309137, 0.01099524,
0.0097979 , 0.00422991, 0.00440997]),
'imprurities': array([0. , 0.00825983, 0.01785371, 0.03548783, 0.05520988,
0.0712241 , 0.00453074, 0.00459993, 0.12424712, 0.14209697,
0.16145084, 0.1688907 , 0.17465136, 0.18125262, 0.28544563,
0.21239646, 0.20555017, 0.14809261, 0.27239399, 0.26564665,
0.32868327, 0.33961473, 0.35127776, 0.36541913, 0.39940962,
0.42938041, 0.45462832, 0.45462832]))

In [22]: ccp_alphas,imprurities = path.ccp_alphas, path.imprurities

In [ ]:

In [ ]:

In [23]: cifs = []
for ccp_alpha in ccp_alphas:
    clf = DecisionTreeClassifier(random_state=0, ccp_alpha=ccp_alpha)
    clf.fit(xtrain, ytrain)
    cifs.append(clf)
    print(
        "Number of nodes in the last tree is: {} with ccp_alpha: {}".format(
            clf[-1].tree_.node_count, ccp_alpha[-1]
        )
    )

Number of nodes in the last tree is: 1 with ccp_alpha: 0.03944080892623343

In [24]: cifs = cifs[:-1]
ccp_alphas = ccp_alphas[:-1]

node_counts = [clf.tree_.node_count for clf in cifs]
depth = [clf.tree_.max_depth for clf in cifs]
fig, ax = plt.subplots(2, 3)
ax[0].plot(ccp_alphas, node_counts, marker="o", drawstyle="steps-post")
ax[0].set_xlabel("alpha")
ax[0].set_ylabel("Number of nodes")
ax[0].set_title("Number of nodes vs alpha")
ax[1].plot(ccp_alphas, depth, marker="o", drawstyle="steps-post")
ax[1].set_xlabel("alpha")
ax[1].set_ylabel("Depth of tree")
ax[1].set_title("Depth vs alpha")
fig.tight_layout()



In [25]: train_scores = [clf.score(xtrain, ytrain) for clf in cifs]
test_scores = [clf.score(xtest, ytest) for clf in cifs]

fig, ax = plt.subplots()
ax.set_xlabel("alpha")
ax.set_ylabel("accuracy")
ax.set_title("Accuracy vs alpha for training and testing sets")
ax.plot(ccp_alphas, train_scores, marker="o", label="train", drawstyle="steps-post")
ax.plot(ccp_alphas, test_scores, marker="o", label="test", drawstyle="steps-post")
ax.legend()
plt.show()



In [26]: randomf=RandomForestClassifier(ccp_alpha=0.0338)
randomf.fit(xtrain,ytrain)
ypred=randomf.predict(xtest)
print('training score :- ',randomf.score(xtrain,ytrain)*100)
print('testing score :- ',randomf.score(xtest,ytest)*100)

training score :- 78.297829782929
testing score :- 64.7858235294117

In [27]: print(confusion_matrix(ytest,ypred))

print(classification_report(ytest,ypred))

[[27 11]
 [13 17]]
precision recall f1-score support
0 0.68 0.71 0.69 38
1 0.61 0.61 0.61 30
accuracy 0.64
macro avg 0.64
weighted avg 0.65
```