## ⌄ Heart Disease Prediction

```python
# importing libraries
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.preprocessing import StandardScaler

import warnings
warnings.filterwarnings('ignore')

sns.set()
plt.style.use('ggplot')
%matplotlib inline
```

```python
#import dataset
heart_df = pd.read_csv('C:\Sunaina\Desktop\Dieseas Prediction app\dataset\heart.csv')
heart_df.head(10)
```

⇥

```python
# information about the dataset
heart_df.info()
```

⇥
```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 303 entries, 0 to 302
Data columns (total 14 columns):
 #   Column    Non-Null Count  Dtype
---  ------    --------------  -----
 0   age       303 non-null    int64
 1   sex       303 non-null    int64
 2   cp        303 non-null    int64
 3   trestbps  303 non-null    int64
 4   chol      303 non-null    int64
 5   fbs       303 non-null    int64
 6   restecg   303 non-null    int64
 7   thalach   303 non-null    int64
 8   exang     303 non-null    int64
 9   oldpeak   303 non-null    float64
 10  slope     303 non-null    int64
 11  ca        303 non-null    int64
 12  thal      303 non-null    int64
 13  target    303 non-null    int64
dtypes: float64(1), int64(13)
memory usage: 33.3 KB
```

```python
#description about dataset
heart_df.describe()
```

⇥

```python
heart_df.shape
```

⇥  (303, 14)

## ⌄ Checking null values

```python
heart_df.isnull().sum()
```

⇥
```
age        0
sex        0
cp         0
trestbps   0
chol       0
fbs        0
restecg    0
thalach    0
exang      0
oldpeak    0
slope      0
```

```
        ca          0
        thal        0
        target      0
        dtype: int64
```

```
heart_df.notnull().sum()
```

```
⇄   age         303
    sex         303
    cp          303
    trestbps    303
    chol        303
    fbs         303
    restecg     303
    thalach     303
    exang       303
    oldpeak     303
    slope       303
    ca          303
    thal        303
    target      303
    dtype: int64
```

```
heart_df.dtypes
```

```
⇄   age          int64
    sex          int64
    cp           int64
    trestbps     int64
    chol         int64
    fbs          int64
    restecg      int64
    thalach      int64
    exang        int64
    oldpeak    float64
    slope        int64
    ca           int64
    thal         int64
    target       int64
    dtype: object
```

## ⌄ Exploratory Data Analysis(EDA)

```
#Plotting the distribution plot.
plt.figure(figsize=(20,25))
plotnumber=1

for column in heart_df:
    if plotnumber<14:
        ax=plt.subplot(4,4,plotnumber)
        sns.distplot(heart_df[column])
        plt.xlabel(column,fontsize=20)
        plt.ylabel('Values',fontsize=20)
    plotnumber+=1
plt.show()
```

```
⇄
```

```
#Correlation matrix

plt.figure(figsize = (16, 8))

corr = heart_df.corr()
mask = np.triu(np.ones_like(corr, dtype = bool))
sns.heatmap(corr, mask = mask, annot = True, fmt = '.2g', linewidths = 1)
plt.show()
```

```
⇄
```

```
#checking the variance
heart_df.var()
```

```
⇄   age         82.484558
    sex          0.217166
    cp           1.065132
```

```
trestbps      307.586453
chol         2686.426748
fbs             0.126877
restecg         0.276528
thalach       524.646406
exang           0.220707
oldpeak         1.348095
slope           0.379735
ca              1.045724
thal            0.374883
target          0.248836
dtype: float64
```

We can see ,there is a huge variance.So,we should normalise it.

## ˅ Normalization

```python
heart_df['trestbps']=np.log(heart_df['trestbps'])
heart_df['chol']=np.log(heart_df['chol'])
heart_df['thalach']=np.log(heart_df['thalach'])
```

```python
np.var(heart_df[["trestbps",'chol','thalach']])
```

```
⮒   trestbps    0.016894
    chol        0.041401
    thalach     0.027054
    dtype: float64
```

```python
heart_df.isnull().sum()
```

```
⮒   age         0
    sex         0
    cp          0
    trestbps    0
    chol        0
    fbs         0
    restecg     0
    thalach     0
    exang       0
    oldpeak     0
    slope       0
    ca          0
    thal        0
    target      0
    dtype: int64
```

```python
x=heart_df.drop('target',axis=1)
y=heart_df['target']
```

```python
#spliting the dataset
from sklearn.model_selection import train_test_split
x_train,x_test,y_train,y_test = train_test_split(x, y, test_size=0.30, random_state=0)
```

```python
x.info()
```

```
⮒   <class 'pandas.core.frame.DataFrame'>
    RangeIndex: 303 entries, 0 to 302
    Data columns (total 13 columns):
     #   Column    Non-Null Count  Dtype
    ---  ------    --------------  -----
     0   age       303 non-null    int64
     1   sex       303 non-null    int64
     2   cp        303 non-null    int64
     3   trestbps  303 non-null    float64
     4   chol      303 non-null    float64
     5   fbs       303 non-null    int64
     6   restecg   303 non-null    int64
     7   thalach   303 non-null    float64
     8   exang     303 non-null    int64
     9   oldpeak   303 non-null    float64
     10  slope     303 non-null    int64
     11  ca        303 non-null    int64
     12  thal      303 non-null    int64
    dtypes: float64(4), int64(9)
    memory usage: 30.9 KB
```

## Logistic Regression

```
accuracies={}

from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score,confusion_matrix,classification_report
lr = LogisticRegression(penalty='l2')
lr.fit(x_train,y_train)

y_pred = lr.predict(x_test)

acc=accuracy_score(y_test,y_pred)
accuracies['LR']=acc*100
print("Training accuracy score of the model is:",accuracy_score(y_train, lr.predict(x_train))*100,"%")
print("Testing accuracy score of the model is:",accuracy_score(y_test,y_pred)*100,"%")
```

```
Training accuracy score of the model is: 85.37735849056604 %
Testing accuracy score of the model is: 80.21978021978022 %
```

```
print("Confusion matrix of the model",confusion_matrix(y_test,y_pred))

print("Classification Report",classification_report(y_test,y_pred))
```

```
Confusion matrix of the model [[32 12]
 [ 6 41]]
Classification Report               precision   recall  f1-score   support

           0       0.84      0.73      0.78        44
           1       0.77      0.87      0.82        47

    accuracy                           0.80        91
   macro avg       0.81      0.80      0.80        91
weighted avg       0.81      0.80      0.80        91
```

## KNearestNeighbors

```
from sklearn.neighbors import KNeighborsClassifier
knn=KNeighborsClassifier(n_neighbors=8)

knn.fit(x_train,y_train)

y_pred1 = knn.predict(x_test)

acc1=accuracy_score(y_test,y_pred1)
accuracies['KNN']=acc1*100

print("Training accuracy score of the model is:",accuracy_score(y_train, knn.predict(x_train))*100,"%")
print("Testing accuracy score of the model is:",accuracy_score(y_test,y_pred1)*100,"%")
```

```
Training accuracy score of the model is: 85.84905660377359 %
Testing accuracy score of the model is: 75.82417582417582 %
```

```
print("Confusion matrix of the model",confusion_matrix(y_test,y_pred1))

print("Classification Report",classification_report(y_test,y_pred1))
```

```
Confusion matrix of the model [[29 15]
 [ 7 40]]
Classification Report               precision   recall  f1-score   support

           0       0.81      0.66      0.72        44
           1       0.73      0.85      0.78        47

    accuracy                           0.76        91
   macro avg       0.77      0.76      0.75        91
weighted avg       0.77      0.76      0.76        91
```

## ∨ SVM

```
from sklearn.svm import SVC

svc = SVC(probability=True)
svc.fit(x_train, y_train)

y_pred2 = svc.predict(x_test)

acc2=accuracy_score(y_test,y_pred2)
accuracies['SVM']=acc2*100

print("Training accuracy score of the model is:",accuracy_score(y_train, svc.predict(x_train))*100,"%")
print("Testing accuracy score of the model is:",accuracy_score(y_test,y_pred2)*100,"%")
```

```
Training accuracy score of the model is: 55.660377358490564 %
Testing accuracy score of the model is: 51.64835164835166 %
```

```
print("Confusion matrix of the model",confusion_matrix(y_test,y_pred2))

print("Classification Report",classification_report(y_test,y_pred2))
```

```
Confusion matrix of the model [[ 0 44]
 [ 0 47]]
Classification Report               precision    recall  f1-score   support

           0       0.00      0.00      0.00        44
           1       0.52      1.00      0.68        47

    accuracy                           0.52        91
   macro avg       0.26      0.50      0.34        91
weighted avg       0.27      0.52      0.35        91
```

## ∨ Decision Tree

```
from sklearn.tree import DecisionTreeClassifier

dtc = DecisionTreeClassifier()
dtc.fit(x_train, y_train)

y_pred3 = dtc.predict(x_test)

acc3=accuracy_score(y_test,y_pred3)
accuracies['DT']=acc3*100

print("Training accuracy score of the model is:",accuracy_score(y_train, dtc.predict(x_train))*100,"%")
print("Testing accuracy score of the model is:",accuracy_score(y_test,y_pred3)*100,"%")
```

```
Training accuracy score of the model is: 100.0 %
Testing accuracy score of the model is: 71.42857142857143 %
```

```
print("Confusion matrix of the model",confusion_matrix(y_test,y_pred3))

print("Classification Report",classification_report(y_test,y_pred3))
```

```
Confusion matrix of the model [[32 12]
 [14 33]]
Classification Report               precision    recall  f1-score   support

           0       0.70      0.73      0.71        44
           1       0.73      0.70      0.72        47

    accuracy                           0.71        91
   macro avg       0.71      0.71      0.71        91
weighted avg       0.72      0.71      0.71        91
```

```
from sklearn.model_selection import GridSearchCV

grid_params = {
    'criterion' : ['gini', 'entropy'],
```

```
        'max_depth' : range(2, 32, 1),
        'min_samples_leaf' : range(1, 10, 1),
        'min_samples_split' : range(2, 10, 1),
        'splitter' : ['best', 'random']
}

grid_search = GridSearchCV(dtc, grid_params, cv = 10, n_jobs = -1, verbose = 1)
grid_search.fit(x_train, y_train)
```

⇥

```
grid_search.best_score_
```

⇥    np.float64(0.8538961038961039)

```
grid_search.best_params_
```

⇥
```
{'criterion': 'gini',
 'max_depth': 15,
 'min_samples_leaf': 7,
 'min_samples_split': 9,
 'splitter': 'random'}
```

```
dtc2 = DecisionTreeClassifier(criterion= 'entropy', max_depth= 12, min_samples_leaf= 1, min_samples_split= 2, splitter= 'random')
dtc2.fit(x_train, y_train)
```

⇥

```
y_pred4 = dtc2.predict(x_test)
acc4=accuracy_score(y_test,y_pred4)
accuracies['DT2']=acc4*100

print("Training accuracy score of the model is:",accuracy_score(y_train, dtc2.predict(x_train))*100,"%")
print("Testing accuracy score of the model is:",accuracy_score(y_test,y_pred4)*100,"%")
```

⇥   Training accuracy score of the model is: 100.0 %
      Testing accuracy score of the model is: 71.42857142857143 %

```
print("Confusion matrix of the model",confusion_matrix(y_test,y_pred4))

print("Classification Report",classification_report(y_test,y_pred4))
```

⇥   Confusion matrix of the model [[33 11]
     [15 32]]
     Classification Report             precision    recall   f1-score    support

```
               0       0.69      0.75      0.72        44
               1       0.74      0.68      0.71        47

        accuracy                           0.71        91
       macro avg       0.72      0.72      0.71        91
    weighted avg       0.72      0.71      0.71        91
```

```
# update dictionary
accuracies['DT']=acc4*100
del accuracies['DT2']
```

## ⌄ Random Forest

```
from sklearn.ensemble import RandomForestClassifier
```

```
rfc = RandomForestClassifier(criterion = 'gini', max_depth = 7, max_features = 'sqrt', min_samples_leaf = 2, min_samples_split = 4, n_estima
rfc.fit(x_train, y_train)
```

```
y_pred5 = rfc.predict(x_test)
```

```
acc5=accuracy_score(y_test,y_pred5)
accuracies['RF']=acc5*100

print("Training accuracy score of the model is:",accuracy_score(y_train, rfc.predict(x_train))*100,"%")
print("Testing accuracy score of the model is:",accuracy_score(y_test,y_pred5)*100,"%")
```

```
Training accuracy score of the model is: 97.16981132075472 %
Testing accuracy score of the model is: 82.41758241758241 %
```

```
print("Confusion matrix of the model",confusion_matrix(y_test,y_pred5))
```

```
print("Classification Report",classification_report(y_test,y_pred5))
```

```
Confusion matrix of the model [[31 13]
 [ 3 44]]
Classification Report              precision    recall  f1-score   support

           0       0.91      0.70      0.79        44
           1       0.77      0.94      0.85        47

    accuracy                           0.82        91
   macro avg       0.84      0.82      0.82        91
weighted avg       0.84      0.82      0.82        91
```

## ∨ Gradient Boosting

```
from sklearn.ensemble import GradientBoostingClassifier
```

```
gbc = GradientBoostingClassifier()
```

```
gbc = GradientBoostingClassifier(learning_rate = 0.05, loss = 'log_loss', n_estimators = 180)
gbc.fit(x_train, y_train)
```

```
y_pred6 = gbc.predict(x_test)
```

```
acc6 = accuracy_score(y_test,y_pred6)
accuracies['GradientBoosting']=acc6*100
```

```
print("Training accuracy score of the model is:",accuracy_score(y_train, gbc.predict(x_train))*100,"%")
print("Testing accuracy score of the model is:",accuracy_score(y_test,y_pred6)*100,"%")
```

```
Training accuracy score of the model is: 100.0 %
Testing accuracy score of the model is: 79.12087912087912 %
```

## ∨ XGBoost

```
from xgboost import XGBClassifier
```

```
xgb = XGBClassifier(objective = 'binary:logistic', learning_rate = 0.01, max_depth = 5, n_estimators = 180)
```

```
xgb.fit(x_train, y_train)
```

```
y_pred7 = xgb.predict(x_test)
```

```
acc7=accuracy_score(y_test,y_pred7)
```

```
accuracies['XGBoost']=acc7*100
print("Training accuracy score of the model is:",accuracy_score(y_train, xgb.predict(x_train))*100,"%")
print("Testing accuracy score of the model is:",accuracy_score(y_test,y_pred7)*100,"%")
```

```
Training accuracy score of the model is: 96.22641509433963 %
Testing accuracy score of the model is: 80.21978021978022 %
```

```
print("Confusion matrix of the model",confusion_matrix(y_test,y_pred7))
```

```
print("Classification Report",classification_report(y_test,y_pred7))
```

```
Confusion matrix of the model [[32 12]
 [ 6 41]]
Classification Report              precision    recall  f1-score   support

           0       0.84      0.73      0.78        44
           1       0.77      0.87      0.82        47
```

```
        accuracy                         0.80        91
       macro avg        0.81      0.80    0.80        91
    weighted avg        0.81      0.80    0.80        91
```

```python
colors = ["purple", "green", "orange", "magenta","blue","black"]

# sns.set_style("whitegrid")
plt.figure(figsize=(16,8))
plt.yticks(np.arange(0,1200,10))
plt.ylabel("Accuracy %")
plt.xlabel("Algorithms")
sns.barplot(x=list(accuracies.keys()), y=list(accuracies.values()), palette=colors )
plt.show()
```

```python
models = pd.DataFrame({
    'Model': ['Logistic Regression', 'KNN', 'SVM',  'Decision Tree', 'Random Forest', 'Gradient Boosting', 'XgBoost'],
    'Score': [acc, acc1, acc2, acc4, acc5, acc6, acc7]
})

models.sort_values(by = 'Score', ascending = False)
```

```python
import pickle
model = rfc
pickle.dump(model, open("heart.pkl",'wb'))
```

```python
from sklearn import metrics
plt.figure(figsize=(8,5))
models = [
{
    'label': 'LR',
    'model': lr,
},
{
    'label': 'DT',
    'model': dtc2,
},
{
    'label': 'SVM',
    'model': svc,
},
{
    'label': 'KNN',
    'model': knn,
},
{
    'label': 'XGBoost',
    'model': xgb,
},
{
    'label': 'RF',
    'model': rfc,
},
{
    'label': 'GBDT',
    'model': gbc,
}
]
for m in models:
    model = m['model']
    model.fit(x_train, y_train)
    y_pred=model.predict(x_test)
    fpr1, tpr1, thresholds = metrics.roc_curve(y_test, model.predict_proba(x_test)[:,1])
    auc = metrics.roc_auc_score(y_test,model.predict(x_test))
    plt.plot(fpr1, tpr1, label='%s - ROC (area = %0.2f)' % (m['label'], auc))

plt.plot([0, 1], [0, 1],'r--')
plt.xlim([-0.01, 1.0])
plt.ylim([0.0, 1.05])
plt.xlabel('1 - Specificity (False Positive Rate)', fontsize=12)
plt.ylabel('Sensitivity (True Positive Rate)', fontsize=12)
```

```
plt.title('ROC - Heart Disease Prediction', fontsize=12)
plt.legend(loc="lower right", fontsize=12)
plt.savefig("roc_heart.jpeg", format='jpeg', dpi=400, bbox_inches='tight')
plt.show()
```

```
from sklearn import metrics
import numpy as np
import matplotlib.pyplot as plt
models = [
{
    'label': 'LR',
    'model': lr,
},
{
    'label': 'DT',
    'model': dtc2,
},
{
    'label': 'SVM',
    'model': svc,
},
{
    'label': 'KNN',
    'model': knn,
},
{
    'label': 'XGBoost',
    'model': xgb,
},
{
    'label': 'RF',
    'model': rfc,
},
{
    'label': 'GBDT',
    'model': gbc,
}
]

means_roc = []
means_accuracy = [100*round(acc,4), 100*round(acc4,4), 100*round(acc2,4), 100*round(acc1,4), 100*round(acc7,4),
                  100*round(acc5,4), 100*round(acc6,4)]

for m in models:
    model = m['model']
    model.fit(x_train, y_train)
    y_pred=model.predict(x_test)
    fpr1, tpr1, thresholds = metrics.roc_curve(y_test, model.predict_proba(x_test)[:,1])
    auc = metrics.roc_auc_score(y_test,model.predict(x_test))
    auc = 100*round(auc,4)
    means_roc.append(auc)

print(means_accuracy)
print(means_roc)

# data to plot
n_groups = 7
means_accuracy = tuple(means_accuracy)
means_roc = tuple(means_roc)

# create plot
fig, ax = plt.subplots(figsize=(8,5))
index = np.arange(n_groups)
bar_width = 0.35
opacity = 0.8

rects1 = plt.bar(index, means_accuracy, bar_width,
alpha=opacity,
color='mediumpurple',
label='Accuracy (%)')

rects2 = plt.bar(index + bar_width, means_roc, bar_width,
alpha=opacity,
color='rebeccapurple',
label='ROC (%)')
```

```
plt.xlim([-1, 8])
plt.ylim([70, 105])

plt.title('Performance Evaluation - Heart Disease Prediction', fontsize=12)
plt.xticks(index, ('  LR', '   DT', '  SVM', '  KNN', 'XGBoost' , '   RF', '   GBDT'), rotation=40, ha='center', fontsize=12)
plt.legend(loc="upper right", fontsize=10)
plt.savefig("PE_heart.jpeg", format='jpeg', dpi=400, bbox_inches='tight')
plt.show()
```