from sklearn.model\_selection import train\_test\_split as t from sklearn.tree import DecisionTreeClassifier as dt df=pd.read\_csv('heart.csv') df.head() trestbps chol fbs restecg thalach exang oldpeak slope ca Out[3]: thal target age sex cp 145 150 37 130 250 187 0 3.5 0 0 0 1 130 204 0 172 0 1.4 2 0 2 1 2 0 56 120 236 178 0 8.0 1 1 163 In [4]: df.corr() exang Out[4]: trestbps chol fbs restecg thalach oldpeak slope ca thal target age sex ср **age** 1.000000 -0.098447 -0.068653 0.210013 -0.225439 0.279351 0.213678 0.121308 -0.116211 -0.398522 0.096801 -0.168814 0.276326 -0.197912 0.045032 -0.058196 1.000000 -0.049353 -0.056769 -0.044020 0.141664 0.096093 -0.030711 0.118261 0.210041 -0.280937 **sex** -0.098447 -0.068653 -0.049353 -0.149230 0.119717 -0.161736 1.000000 0.047608 -0.076904 0.094444 0.044421 0.295762 -0.394280 -0.181053 -0.121475 1.000000 0.177531 -0.114103 trestbps 0.279351 -0.056769 0.047608 0.123174 -0.046698 0.067616 0.193216 0.101389 0.062210 -0.144931 0.213678 -0.197912 -0.076904 0.123174 1.000000 0.013294 -0.151040 -0.009940 0.067023 0.053952 -0.004038 0.070511 0.098803 0.177531 0.005747 0.121308 0.045032 0.094444 0.013294 1.000000 -0.084189 -0.008567 0.025665 -0.059894 0.137979 -0.032019 -0.028046 -0.084189 -0.058770 0.093045 -0.011981 -0.116211 -0.058196 0.044421 -0.114103 -0.151040 1.000000 0.044123 -0.070733 -0.072042 0.137230 restecg thalach -0.398522 0.295762 -0.046698 -0.044020 -0.009940 -0.008567 0.044123 1.000000 -0.378812 -0.344187 0.386784 -0.213177 -0.096439 0.421741 0.096801 0.141664 -0.394280 0.067616 0.067023 0.025665 -0.070733 -0.378812 1.000000 0.288223 -0.257748 0.115739 0.206754 **oldpeak** 0.210013 0.096093 -0.149230 0.193216 0.053952 0.005747 -0.058770 -0.344187 0.288223 1.000000 -0.577537 0.222682 0.210244 -0.430696 -0.168814 -0.030711 0.119717 -0.121475 -0.004038 -0.059894 0.093045 0.386784 -0.257748 -0.577537 1.000000 -0.080155 -0.104764 -0.213177 0.101389 0.276326 0.118261 -0.181053 0.070511 0.137979 -0.072042 0.115739 0.222682 -0.080155 1.000000 0.151832 -0.391724 0.068001 0.210041 -0.161736 0.062210 0.098803 -0.032019 -0.011981 -0.096439 0.206754 0.210244 -0.104764 0.151832 target -0.225439 -0.280937 0.433798 -0.144931 -0.085239 -0.028046 0.137230 0.421741 -0.436757 -0.430696 0.345877 -0.391724 -0.344029 1.000000 In [ ]: In [ ]: In [5]: x=df.iloc[:,:-1] y=df.iloc[:,-1] xtrain, xtest, ytrain, ytest=t(x, y, random\_state=0) In [38]: dtree=dt() dtree.fit(xtrain,ytrain) ypred=dtree.predict(xtest) print(ac(ytest,ypred)) 0.8026315789473685 In [44]: print('Training score') dtree.score(xtrain,ytrain) Training score Out[44]: 1.0 In [45]: print('Testing score') dtree.score(xtest,ytest) Testing score Out[45]: 0.8026315789473685 In [41]: from sklearn import tree In the below the plot we have plotted the Decision Tree We can the max depth is 10 In [43]: plt.figure(figsize=(8,40)) chart=tree.plot\_tree(dtree, filled=True, fontsize=10, rounded=True, proportion=True) for node in chart: arrow=node.arrow\_patch if (arrow is not None): arrow.set\_edgecolor('red') arrow.set\_linewidth(3) X[2] <= 0.5 gini = 0.497 samples = 100.0% value = [0.463, 0.537] samples = 51.5% value = [0.727, 0.273] value = [0.214, 0.786] X[11] <= 1.5 gini = 0.043 samples = 19.8% value = [0.022, 0.978] X[7] <= 150.5 gini = 0.444 X[12] <= 2 gini = 0.49 samples = 22 samples = 31.7% value = [0.333, 0.667] samp value = [0.46, X[0] <= 57.5 gini = 0.305 samples = 10.6% samples = 13.2 samples = 21.1% value = [0.625, 0.375]alue = [0.188, 0.812] X[9] <= 1.2 gini = 0.469 gini = 0.434samples = 3.5% value = [0.682, 0.318] value value = [0.625, 0.375]X[4] gini samples samples samp samp value = value = [0.84]gini = 0.0 gin samples = 0.4% value = [1.0, 0.0] samples = 0.4% value = [1.0, 0.0 X[] g san samples = 0.4%samples = 14.5%X[0] <= 63.0 gini = 0.5 samples = 0.9% value = [0.5, 0.5] X[0] <= 51.0 gini = 0.444 samples = 1.3% value = [0.667, 0.333] Below we have use Post pruning decision tree with cost complexity to overcome the situation of overfitting In [ ]: path=dtree.cost\_complexity\_pruning\_path(xtrain,ytrain) In [13]: path.ccp\_alphas Out[13]: array([0. , 0.00367107, 0.00430739, 0.00437336, 0.00528634,  $\hbox{0.00587372, 0.00612103, 0.00660793, 0.00704846, 0.0074551, } \\$ 0.00751217, 0.00763395, 0.00837004, 0.00991189, 0.0100487, 0.01057269, 0.01109713, 0.01618943, 0.02529886, 0.03395007, 0.05769323, 0.13176578]) In [14]: path.impurities Out[14]: array([0. , 0.00734214, 0.01595693, 0.02470365, 0.03527634, 0.04702377, 0.0531448 , 0.05975273, 0.06680119, 0.08171139, 0.09673573, 0.11963757, 0.14474771, 0.1546596 , 0.17475701,  $0.19590238,\ 0.20699951,\ 0.22318894,\ 0.27378666,\ 0.30773674,$ 0.36542997, 0.49719575]) alphas,impure=path.ccp\_alphas , path.impurities In [16]: clfs=[] for i in alphas: i=dt(random\_state=0,ccp\_alpha=i) i.fit(xtrain,ytrain) clfs.append(i) In [17]: print(f'Number of nodes in the last tree is {clfs[-1].tree\_.node\_count} {alphas[-1]}') Number of nodes in the last tree is 1 0.13176578489314994 training\_score=[clf.score(xtrain,ytrain) for clf in clfs] testing\_score=[clf.score(xtest,ytest) for clf in clfs] In [19]: fig, ax=plt.subplots() plt.grid() ax.set\_xlabel('alpha') ax.set\_ylabel('accuracy') ax.set\_title('accuracy vs alpha for training and testing') ax.plot(alphas, training\_score, marker='o', label='train', drawstyle='steps-post') ax.plot(alphas, testing\_score, marker='o', label='test', drawstyle='steps-post') ax.legend() plt.show() accuracy vs alpha for training and testing 1.0 0.6 0.02 0.04 0.06 0.08 0.10 From the above the graph the training we can see the training and testing score are near to each other at the 0.03 value of aplha In [48]: dtree=dt(ccp\_alpha=0.03) dtree.fit(xtrain,ytrain) ypred=dtree.predict(xtest) In [49]: dtree.score(xtrain,ytrain) 0.8281938325991189 In [50]: dtree.score(xtest,ytest) Out[50]: 0.8157894736842105 From the above the score we can see the it has overcomed the situation of overfitting to besting fitting In [51]: plt.figure(figsize=(8,15)) chart=tree.plot\_tree(dtree, filled=True, rounded=True, proportion=True) for node in chart: arrow=node.arrow\_patch if (arrow is not None): arrow.set\_edgecolor('red') arrow.set\_linewidth(3)  $X[2] \le 0.5$ gini = 0.497 samples = 100.0% value = [0.463, 0.537] X[11] <= 0.5 gini = 0.397 gini = 0.336samples = 51.5% samples = 48.5%value = [0.214, 0.786] value = [0.727, 0.273] X[12] <= 2.5 gini = 0.497 gini = 0.095samples = 26.4% samples = 22.0%value = [0.95, 0.05]value = [0.46, 0.54] gini = 0.358gini = 0.32samples = 8.8% samples = 13.2%value = [0.233, 0.767] value = [0.8, 0.2]By comparing the above the decision plot to above plot we can see the max depth has been reduced to 3 from 10

In [1]:

import pandas as pd
import numpy as np

import seaborn as sns

import warnings

import matplotlib.pyplot as plt

warnings.filterwarnings('ignore')

from sklearn.metrics import accuracy\_score as ac ,confusion\_matrix as cm ,classification\_report as cr