Decision Tree for Classification of Breast Cancer Wisconsin Diagnostic In [1]: import pandas as pd import numpy as np
<pre>import seaborn as sns import matplotlib.pyplot as plt from matplotlib.pyplot import figure from sklearn.model_selection import train_test_split from sklearn.metrics import accuracy_score,precision_score,recall_score from sklearn.tree import DecisionTreeClassifier from sklearn import tree</pre>
<pre>#conda install -c conda-forge pygraphviz import graphviz from sklearn import preprocessing colname=['ID','Diagnosis']+['feature_%s' %i for i in range(1,31)] #print(colname) df = pd.read_csv('wdbc.data', header=None, names=colname)</pre>
1. After importing the data, let's check the if data successfully import. In [2]: print(df.head(5)) ID Diagnosis feature 1 feature 2 feature 4 feature 5 \
842302 M 17.99 10.38 122.80 1001.0 0.1340 1 842517 M 20.57 17.77 132.90 1326.0 0.08474 2 84300903 M 19.69 21.25 130.00 1203.0 0.10960 3 84348301 M 11.42 20.38 77.58 386.1 0.14250 4 84358402 M 20.29 14.34 135.10 1297.0 0.10930 feature 6 feature 7 feature 8 feature 21 feature 22 feature 23 \ 0 0.27760 0.3001 0.14710 25.38 17.33 184.60 1 0.07864 0.0869 0.07017 24.99 23.41 158.80 2 0.15990 0.1974 0.12790 24.99 23.41 158.80 2 0.15990 0.1974 0.10520 14.91 26.50 98.87 4 0.13280 0.1980 0.10540 22.54 16.67 152.20 feature 24 feature 25 feature 26 feature 27 feature 28 feature 29 \ 0 2019.0 0.1622 0.6656 0.7119 0.2654 0.4601 1 1956.0 0.1238 0.1866 0.2416 0.1860 0.2750 2 1709.0 0.1444 0.4245 0.4504 0.2430 0.3613 3 567.7 0.2099 0.663 0.6699 0.5755 0.6638 4 1575.0 0.1374 0.250 0.4000 0.1625 0.2364 feature 20 0 0.08788 3 0.17300 4 0.07678
<pre>[5 rows x 32 columns] In [3]: df.isnull().any() Out[3]: ID False Diagnosis False feature_1 False</pre>
feature 2 False feature 3 False feature 4 False feature 5 False feature 6 False feature 7 False feature 7 False feature 8 False feature 10 False feature 11 False feature 12 False feature 12 False feature 13 False feature 14 False feature 14 False feature 15 False feature 16 False feature 17 False feature 18 False feature 18 False feature 18 False feature 18 False feature 19 False feature 20 False feature 21 False feature 21 False feature 20 False feature 21 False feature 22 False feature 23 False feature 24 False feature 24 False feature 25 False feature 26 False feature 27 False feature 27 False feature 28 False feature 28 False feature 28 False feature 28 False feature 29 False feature 30 False feature 30 False
1. Let's understand more about the data. We will start by getting to know the type of each column values. We see that the width and length column are represented using float64 and the name of the species uses object or string. In [4]: df.dtypes Out [4]: Diagraphic species are object.
Diagnosis object feature_1 float64 feature_3 float64 feature_4 float64 feature_5 float64 feature_6 float64 feature_7 float64 feature_8 float64 feature_9 float64 feature_10 float64 feature_11 float64 feature_11 float64 feature_12 float64 feature_13 float64 feature_14 float64 feature_15 float64 feature_16 float64 feature_17 float64 feature_18 float64 feature_19 float64 feature_19 float64 feature_19 float64 feature_19 float64 feature_10 float64 feature_10 float64 feature_11 float64 feature_12 float64 feature_13 float64 feature_14 float64 feature_15 float64 feature_19 float64 feature_19 float64 feature_20 float64 feature_21 float64 feature_21 float64 feature_22 float64 feature_23 float64 feature_24 float64 feature_25 float64 feature_26 float64 feature_27 float64 feature_28 float64 feature_29 float64 feature_30 float64 feature_40 float64 feature_50 float64 feature_50 float64 feature_60 float64 feature_70 float64 feat
Out [5]: 1. Let's look at a quick summary of the data.
Out [6]: ID feature_1 feature_2 feature_3 feature_4 feature_5 feature_6 feature_7 feature_8 feature_9 feature_9 feature_21 feature_22 feature_23 feature_24 feature_24 feature_25 feature_24 feature_25 feature_25 feature_26 feature_27 feature_28 feature_29 feature_21 feature_22 feature_23 feature_24 feature_26 feature_27 feature_28 feature_29 feature_29 feature_21 feature_22 feature_23 feature_24 feature_27 feature_28 feature_29 feature_29 feature_29 feature_29 feature_21 feature_22 feature_23 feature_24 feature_28 feature_29 feature_28 feature_29 featu
1. Everything checks out. First, I use Information Gain to construct a DT. But choose which features? In [7]: #distribute M and B dfm=df[df['Diagnosis']=='M'] dfb=df[df['Diagnosis']=='B']
Out [7]: ID feature_1 feature_2 feature_3 feature_4 feature_5 feature_6 feature_7 feature_8 feature_9 feature_9 feature_21 feature_22 feature_23 feature_24 feature_24 feature_24 feature_25 feature_26 feature_26 feature_26 feature_26 feature_27 feature_28 feature_29
0ut[8]: ID feature_1 feature_2 feature_3 feature_4 feature_5 feature_6 feature_7 feature_8 feature_9 feature_9 feature_21 feature_22 feature_23 feature_24 feature_25 count 3.5700000e+02 357.00000 35
<pre>In [9]: #initialize avg_m=np.zeros(30) avg_b=np.zeros(30) for i_feature in range(30): #obtain average value of feature for M and B avg_m[i_feature] = np.mean(dfm[lambda dfm: dfm.columns[i_feature+2]]) avg_b[i_feature] = np.mean(dfb[lambda dfb: dfb.columns[i_feature+2]]) # find top 10 features with average value of M and B differs most diff=abs(avg_m-avg_b) diff_sorted = sorted(enumerate(diff),key=lambda x:x[1]) diff_sorted_index=[m[0] for m in diff_sorted] print(diff_sorted_index[0:10]) [18, 9, 14, 19, 17, 11, 4, 15, 29, 16]</pre>
5.construct decision tree only from the 10 features and loop for 20 times In [10]: X = df[['feature_18', 'feature_9', 'feature_19', 'feature_17',
<pre>dtl_accuracy_test=pp.zeros(20) dtl_precision=pp.zeros(20) dtl_recall=np.zeros(20) dtl_accuracy_test=np.zeros(20) dtl_accuracy_test=np.zeros(20) dtl_accuracy_test=np.zeros(20) dtl_precision=np.zeros(20) dtl_precision=np.zeros(20) for i_process in range(20): (X_train_X_test, Y_train_X_test, Y_train_X_test) = train_test_split(X, Y, random_state=2,test_size=0.3,stratify=Y) clf l = DecisionTreeclassifier(criterion='entropy') clf l = DecisionTreeclassifier(criterion='entropy') clf l.fit(X_train, Y_train) Y_pre_train=clf_l.predict(X_train) Y_pre_test=clf_l.predict(X_train) Y_pre_test=clf_l.predict(X_train) dtl_accuracy_train[i_process]=accuracy_score(Y_test, Y_pre_test) dtl_accuracy_train[i_process]=accuracy_score(Y_test, Y_pre_test) dtl_precision[i_process]=accuracy_test, Y_pre_test, Decision[score(Y_test, Y_pre_test)] dtl_recall[i_process]=ccuracy_test, Decision[score(Y_test, Y_pre_test)] dtl_recall[i_process]=ccuracy_test, Decision[score(Y_test, Y_pre_test)] dtl_recall[i_process]=ccuracy_test, Decision[score(Y_test, Y_pre_test)] dtl_recall[i_process]=ccuracy_test, Decision[score(Y_test, Y_pre</pre>
<pre>dt2_accuracy_train[i_process] = accuracy_score(Y_train, Y_pre_train) dt2_accuracy_test[i_process] = accuracy_score(Y_test, Y_pre_test) dt2_precision[i_process] = precision_score(Y_test, Y_pre_test,pos_label = 'M') dt2_recall[i_process] = recall_score(Y_test, Y_pre_test,pos_label = 'M') dt2_recall[i_process] = last_recall</pre>
<pre>dt1_avg_accuracy_train=np.mean(dt1_accuracy_train) dt1_avg_accuracy_test=np.mean(dt1_accuracy_test) dt1_avg_precision=np.mean(dt1_precision) dt1_avg_recall=np.mean(dt1_recall) dt2_avg_accuracy_train=np.mean(dt2_accuracy_train) dt2_avg_accuracy_test=np.mean(dt2_accuracy_test) dt2_avg_precision=np.mean(dt2_precision) dt2_avg_recall=np.mean(dt2_precision)</pre>
<pre>print([dt1_accuracy_train, dt2_accuracy_train]) print([dt1_accuracy_test,dt2_accuracy_test]) print([dt1_avg_accuracy_train,dt1_avg_accuracy_test,dt1_avg_precision,dt1_avg_recall]) print([dt2_avg_accuracy_train,dt2_avg_accuracy_test,dt2_avg_precision,dt2_avg_recall]) [array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.,</pre>
0.90643275, 0.90643275, 0.9122807, 0.9122807, 0.9122807, 0.9122807, 0.9005848, 0.9005848, 0.9005848, 0.90643275, 0.90643275, 0.90643275, 0.9005848, 0.91812865, 0.91812865, 0.90643275, 0.90643275, 0.9122807])] [1.0, 0.9064327485380117, 0.8811484600840431, 0.8671875] [0.9974874371859297, 0.9076023391812866, 0.8784219174642562, 0.89765625] 6.construct the decision tree from all 30 features and replace depth by minimum leaf samples
<pre>In [11]: X = df[lambda df: df.columns[2:32]].values</pre>
<pre>min_sample=2 last_recall=1 while dtl_recall[i_process] > dt2_recall[i_process] or last_recall < dt2_recall[i_process]: last_recall=dt2_recall[i_process] min_sample+=1 clf_2 = DecisionTreeClassifier(criterion='entropy',min_samples_split=min_sample) clf_2.fit(X_train, Y_train) Y pre_train=clf_2.predict(X_train) Y_pre_test=clf_2.predict(X_test) dt2_accuracy_train[i_process]=accuracy_score(Y_train, Y_pre_train) dt2_accuracy_test[i_process]=accuracy_score(Y_test, Y_pre_test) dt2_precision[i_process]=precision_score(Y_test, Y_pre_test,pos_label='M') dt2_recall[i_process]=last_recall</pre>
<pre>dt1_avg_accuracy_train=np.mean(dt1_accuracy_train) dt1_avg_accuracy_test=np.mean(dt1_accuracy_test) dt1_avg_precision=np.mean(dt1_precision) dt1_avg_recall=np.mean(dt1_recall) dt2_avg_accuracy_train=np.mean(dt2_accuracy_train) dt2_avg_accuracy_test=np.mean(dt2_accuracy_test)</pre>
dt2_avg_precision=np.mean(dt2_precision) dt2_avg_recall=np.mean(dt2_recall) print([dt1_accuracy_train,dt2_accuracy_train]) print([dt1_accuracy_test,dt2_accuracy_test]) print([dt1_avg_accuracy_train,dt1_avg_accuracy_test,dt1_avg_precision,dt1_avg_recall]) print([dt2_avg_accuracy_train,dt2_avg_accuracy_test,dt2_avg_precision,dt2_avg_recall]) [array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.,
0.91812865, 0.92397661, 0.90643275, 0.91812865, 0.9122807 , 0.93567251, 0.9122807 , 0.91812865, 0.92982456, 0.9122807]), array([0.9122807 , 0.92982456, 0.93567251, 0.93567251, 0.93567251, 0.93567251, 0.92397661, 0.93567251, 0.91812865, 0.92982456, 0.88304094, 0.92397661, 0.92982456, 0.91812865, 0.92397661, 0.92982456, 0.88304094, 0.92397661, 0.9122807 , 0.92397661, 0.92982456])] [1.0, 0.92046783625731, 0.9053330039099803, 0.8796875] [0.9868090452261307, 0.9219298245614036, 0.9039124400043848, 0.9]