

✓ Lab Assignment 6 - Decision Trees

In this lab exercise, you will learn a popular machine learning algorithm, Decision Tree. You will use this classification algorithm to build a model from historical data of patients, and their respond to different medications. Then you use the trained decision tree to predict the class of a unknown patient, or to find a proper drug for a new patient.

Import the Following Libraries:

- **numpy (as np)**
- **pandas**
- **DecisionTreeClassifier** from **sklearn.tree**

```
1 import numpy as np
2 import pandas as pd
3 from sklearn.tree import DecisionTreeClassifier
```

✓ About dataset

Imagine that you are a medical researcher compiling data for a study. You have collected data about a set of patients, all of whom suffered from the same illness. During their course of treatment, each patient responded to one of 5 medications, Drug A, Drug B, Drug c, Drug x and y.

Part of your job is to build a model to find out which drug might be appropriate for a future patient with the same illness. The feature sets of this dataset are Age, Sex, Blood Pressure, and Cholesterol of patients, and the target is the drug that each patient responded to.

It is a sample of binary classifier, and you can use the training part of the dataset to build a decision tree, and then use it to predict the class of a unknown patient, or to prescribe it to a new patient.

now, read data using pandas dataframe:

```
1 my_data = pd.read_csv("/content/drug200.csv", delimiter=",")
2 my_data[0:5]
```

	Age	Sex	BP	Cholesterol	Na_to_K	Drug
0	23	F	HIGH	HIGH	25.355	drugY
1	47	M	LOW	HIGH	13.093	drugC
2	47	M	LOW	HIGH	10.114	drugC
3	28	F	NORMAL	HIGH	7.798	drugX
4	61	F	LOW	HIGH	18.043	drugY

✓ Practice

What is the size of data?

```
1 # write your code here
2 data_size = my_data.size
3 print("The size of this data is: ",data_size)
```

The size of this data is: 1200

✓ Pre-processing

Using **my_data** as the Drug.csv data read by pandas, declare the following variables:

- **X** as the **Feature Matrix** (data of my_data)
- **y** as the **response vector (target)**

Remove the column containing the target name since it doesn't contain numeric values.

```
1 X = my_data[['Age', 'Sex', 'BP', 'Cholesterol', 'Na_to_K']].values
2 X[0:5]
```

```
array([[23, 'F', 'HIGH', 'HIGH', 25.355],
       [47, 'M', 'LOW', 'HIGH', 13.093],
       [47, 'M', 'LOW', 'HIGH', 10.114],
       [28, 'F', 'NORMAL', 'HIGH', 7.798],
       [61, 'F', 'LOW', 'HIGH', 18.043]], dtype=object)
```

As you may figure out, some features in this dataset are categorical such as **Sex** or **BP**. Unfortunately, Sklearn Decision Trees do not handle categorical variables. But still we can convert these features to numerical values. **pandas.get_dummies()** Convert categorical variable into dummy/indicator variables.

```
1 from sklearn import preprocessing
2 le_sex = preprocessing.LabelEncoder()
3 le_sex.fit(['F', 'M'])
4 X[:,1] = le_sex.transform(X[:,1])
5
6
7 le_BP = preprocessing.LabelEncoder()
8 le_BP.fit(['LOW', 'NORMAL', 'HIGH'])
9 X[:,2] = le_BP.transform(X[:,2])
10
11
12 le_Chol = preprocessing.LabelEncoder()
13 le_Chol.fit(['NORMAL', 'HIGH'])
14 X[:,3] = le_Chol.transform(X[:,3])
15
16 X[0:5]
17
```

```
array([[23, 0, 0, 0, 25.355],
       [47, 1, 1, 0, 13.093],
       [47, 1, 1, 0, 10.114],
       [28, 0, 2, 0, 7.798],
       [61, 0, 1, 0, 18.043]], dtype=object)
```

Now we can fill the target variable.

```
1 y = my_data["Drug"]
2 y[0:5]
```

```
0    drugY
1    drugC
2    drugC
3    drugX
4    drugY
Name: Drug, dtype: object
```

Setting up the Decision Tree

We will be using **train/test split** on our **decision tree**. Let's import **train_test_split** from **sklearn.cross_validation**.

```
1 from sklearn.model_selection import train_test_split
```

Now **train_test_split** will return 4 different parameters. We will name them:

X_trainset, X_testset, y_trainset, y_testset

The **train_test_split** will need the parameters:

X, y, test_size=0.3, and random_state=3.

The **X** and **y** are the arrays required before the split, the **test_size** represents the ratio of the testing dataset, and the **random_state** ensures that we obtain the same splits.

```
1 X_trainset, X_testset, y_trainset, y_testset = train_test_split(X, y, test_size=0.3, random_state=3)
```

Practice

Print the shape of X_trainset and y_trainset. Ensure that the dimensions match

```

1 # your code
2 print("The X_trainset is: ",X_trainset.shape)
3 print ("The y_trainset is: ",y_trainset.shape)
4
5 if X_trainset.shape[0] == y_trainset.shape[0]:
6     print("Everything is ok, both values from X_trainset and y_trainset are equal")
7 else:
8     print("Error, X_trainset and y_trainset has different values")

```

```

↗ The X_trainset is: (140, 5)
  The y_trainset is: (140,)
  Everything is ok, both values from X_trainset and y_trainset are equal

```

Print the shape of **X_testset** and **y_testset**. Ensure that the dimensions match

```

1 # your code
2 print("The X_trainset is: ",X_testset.shape)
3 print ("The y_trainset is: ",y_testset.shape)
4
5 if X_testset.shape[0] == y_testset.shape[0]:
6     print("Everything is ok, both values from X_trainset and y_trainset are equal")
7 else:
8     print("Error, X_trainset and y_trainset has different values")
9
10

```

```

↗ The X_trainset is: (60, 5)
  The y_trainset is: (60,)
  Everything is ok, both values from X_trainset and y_trainset are equal

```

✓ Modeling

We will first create an instance of the **DecisionTreeClassifier** called **drugTree**.

Inside of the classifier, specify *criterion="entropy"* so we can see the information gain of each node.

```

1 drugTree = DecisionTreeClassifier(criterion="entropy", max_depth = 4)
2 drugTree # it shows the default parameters

```

```

↗
  ▼ DecisionTreeClassifier
  DecisionTreeClassifier(criterion='entropy', max_depth=4)

```

Next, we will fit the data with the training feature matrix **X_trainset** and training response vector **y_trainset**

```

1 drugTree.fit(X_trainset,y_trainset)

```

```

↗
  ▼ DecisionTreeClassifier
  DecisionTreeClassifier(criterion='entropy', max_depth=4)

```

✓ Prediction

Let's make some **predictions** on the testing dataset and store it into a variable called **predTree**.

```

1 predTree = drugTree.predict(X_testset)

```

You can print out **predTree** and **y_testset** if you want to visually compare the prediction to the actual values.

```

1 print (predTree [0:5])
2 print (y_testset [0:5])
3

```

```


↗ ['drugY' 'drugX' 'drugX' 'drugX' 'drugX']
40      drugY
51      drugX
139     drugX
197     drugX
170     drugX
Name: Drug, dtype: object

```

✓ Evaluation

Next, let's import **metrics** from sklearn and check the accuracy of our model.

```
1 from sklearn import metrics
2 import matplotlib.pyplot as plt
3 print("DecisionTrees's Accuracy: ", metrics.accuracy_score(y_testset, predTree))
4 xteste = np.array(predTree)
5 xteste[0]
6
7 #predTree[0:1]
```

 DecisionTrees's Accuracy: 0.9833333333333333
'drugY'


Accuracy classification score computes subset accuracy: the set of labels predicted for a sample must exactly match the corresponding set of labels in y_true.

In multilabel classification, the function returns the subset accuracy. If the entire set of predicted labels for a sample strictly match with the true set of labels, then the subset accuracy is 1.0; otherwise it is 0.0.

✓ Practice

Can you calculate the accuracy score without sklearn ?

```
1 # your code here
2
3 i = y_testset.size
4 x = 0
5 ok = 0
6 fail = 0
7 y_testset_comp = np.array(y_testset)
8 predTree_comp = np.array(predTree)
9
10 for x in range(0,i):
11     if y_testset_comp[x] == predTree_comp[x]:
12         ok = ok + 1
13     else:
14         fail = fail + 1
15
16 accuracy = ok/i
17 print(accuracy)
18
```

 0.9833333333333333

✓ Visualization

Lets visualize the tree

```
1 !pip install six
```

 Requirement already satisfied: six in /usr/local/lib/python3.10/dist-packages (1.16.0)

```
1 from six import StringIO
2 import pydotplus
3 import matplotlib.image as mpimg
4 from sklearn import tree
5 %matplotlib inline
```

```
1 dot_data = StringIO()
2 filename = "drugtree.png"
3 featureNames = my_data.columns[0:5]
4 targetNames = my_data["Drug"].unique().tolist()
5 out=tree.export_graphviz(drugTree,feature_names=featureNames, out_file=dot_data, class_names= np.unique(y_trainset), filled=True, spe
6 graph = pydotplus.graph_from_dot_data(dot_data.getvalue())
7 graph.write_png(filename)
8 img = mpimg.imread(filename)
9 plt.figure(figsize=(100, 200))
10 plt.imshow(img,interpolation='nearest')
```

 <matplotlib.image.AxesImage at 0x7f06502332e8>