

# 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 response 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.

## Table of contents

1. [About the dataset \(https://#about\\_dataset\)](https://#about_dataset)
2. [Downloading the Data \(https://#downloading\\_data\)](https://#downloading_data)
3. [Pre-processing \(https://#pre-processing\)](https://#pre-processing)
4. [Setting up the Decision Tree \(https://#setting\\_up\\_tree\)](https://#setting_up_tree)
5. [Modeling \(https://#modeling\)](https://#modeling)
6. [Prediction \(https://#prediction\)](https://#prediction)
7. [Evaluation \(https://#evaluation\)](https://#evaluation)
8. [Visualization \(https://#visualization\)](https://#visualization)

---

Import the Following Libraries:

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

```
In [1]: import numpy as np
import pandas as pd
from sklearn.tree import DecisionTreeClassifier
```

## About the 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

## Downloading the Data

To download the data, we will use pandas read method.

```
In [3]: my_data = pd.read_csv("drug200.csv", delimiter=",")
my_data[0:5]
```

Out[3]:

	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?

```
In [ ]: # write your code here
```

## 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.

```
In [4]: X = my_data[['Age', 'Sex', 'BP', 'Cholesterol', 'Na_to_K']].values
X[0:5]
```

```
Out[4]: 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.

```
In [5]: from sklearn import preprocessing
le_sex = preprocessing.LabelEncoder()
le_sex.fit(['F', 'M'])
X[:,1] = le_sex.transform(X[:,1])

le_BP = preprocessing.LabelEncoder()
le_BP.fit([ 'LOW', 'NORMAL', 'HIGH'])
X[:,2] = le_BP.transform(X[:,2])

le_Chol = preprocessing.LabelEncoder()
le_Chol.fit([ 'NORMAL', 'HIGH'])
X[:,3] = le_Chol.transform(X[:,3])

X[0:5]
```

```
Out [5]: 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.

```
In [6]: y = my_data["Drug"]
y[0:5]
```

```
Out [6]: 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**.

```
In [7]: 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.

```
In [8]: X_trainset, X_testset, y_trainset, y_testset = train_test_split(X, y,
```

## Practice

Print the shape of X\_trainset and y\_trainset. Ensure that the dimensions match

```
In [ ]: # your code
```

Print the shape of X\_testset and y\_testset. Ensure that the dimensions match

```
In [ ]: # your code
```

---

## 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.

```
In [9]: drugTree = DecisionTreeClassifier(criterion="entropy", max_depth = 4)
drugTree # it shows the default parameters
```

```
Out[9]: 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**

```
In [10]: drugTree.fit(X_trainset,y_trainset)
```

```
Out[10]: DecisionTreeClassifier(criterion='entropy', max_depth=4)
```

---

## Prediction

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

```
In [11]: 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.

```
In [12]: print (predTree [0:5])
print (y_testset [0:5])

['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.

```
In [13]: from sklearn import metrics
import matplotlib.pyplot as plt
print("DecisionTrees's Accuracy: ", metrics.accuracy_score(y_testset,

DecisionTrees's Accuracy:  0.9833333333333333
```

**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 ?

```
In [ ]: # your code here
```

---

## Visualization

Lets visualize the tree

```
In [ ]: # Notice: You might need to uncomment and install the pydotplus and graphviz packages
# !conda install -c conda-forge pydotplus -y
# !conda install -c conda-forge python-graphviz -y
```

```
In [14]: from sklearn.externals.six import StringIO
import pydotplus
import matplotlib.image as mpimg
from sklearn import tree
%matplotlib inline
```

/Users/sumitkumarshukla/opt/anaconda3/lib/python3.8/site-packages/sklearn/externals/six.py:28: FutureWarning: The module is deprecated in version 0.21 and will be removed in version 0.23 since we've dropped support for Python 2.7. Please rely on the official version of six (<https://pypi.org/project/six/>).

warnings.warn("The module is deprecated in version 0.21 and will be removed ")

```
In [18]: dot_data = StringIO()
filename = "drugtree.png"
featureNames = my_data.columns[0:5]
targetNames = my_data["Drug"].unique().tolist()
out=tree.export_graphviz(drugTree,feature_names=featureNames, out_file=None)
graph = pydotplus.graph_from_dot_data(dot_data.getvalue())
graph.write_png(filename)
img = mpimg.imread(filename)
plt.figure(figsize=(100, 200))
plt.imshow(img,interpolation='nearest')

from sklearn import tree
tree_rep = tree.export_text(drugTree)
print(tree_rep)
```

```
|--- feature_4 <= 14.62
|   |--- feature_2 <= 0.50
|   |   |--- feature_0 <= 50.50
|   |   |   |--- class: drugA
|   |   |--- feature_0 > 50.50
|   |   |   |--- class: drugB
|   |--- feature_2 > 0.50
|   |   |--- feature_3 <= 0.50
|   |   |   |--- feature_2 <= 1.50
|   |   |   |   |--- class: drugC
|   |   |   |--- feature_2 > 1.50
|   |   |   |   |--- class: drugX
|   |   |--- feature_3 > 0.50
|   |   |   |--- class: drugX
|--- feature_4 > 14.62
|   |--- class: drugY
```

```
In [19]: fig = plt.figure(figsize=(30,25))
vis = tree.plot_tree(drugTree,
                    filled=True)
```

