

Classification Trees in python

from start to finish

* In this lesson I used scikit-learn and cost complexity pruning to build a classification tree.
Helps in predicting whether a patient has heart disease or not

- We will :
- 1) Import data from a file
 - 2) Identify & deal with missing data
 - 3) Format the data for decision trees
 - 4) Build a preliminary classification tree
 - 5) Use, cost complexity pruning to improve the tree
 - 6) Build, Evaluate, Draw and interpret final tree

Classification trees are an exceptionally useful machine learning method when you need to know how the decision are being made. For example, if you have to justify the predictions to your boss, classification trees are a good method because each step in the decision making process is easy to understand.

① Import the modules

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from sklearn.tree import DecisionTreeClassifier
from sklearn.tree import plot_tree
from sklearn.model_selection import train_test_split
from sklearn.model_selection import cross_val_score
from sklearn.metrics import confusion_matrix
from sklearn.metrics import plot_confusion_matrix
```


2) Import the data: (In Pandas read returns data frame)

> df = read_csv ('processed-cleveland-data', header=None)

↳ The above command loads the data into the dataframe df.

lets look at the first 5 rows using the head() function

» df.head()

Rename Columns

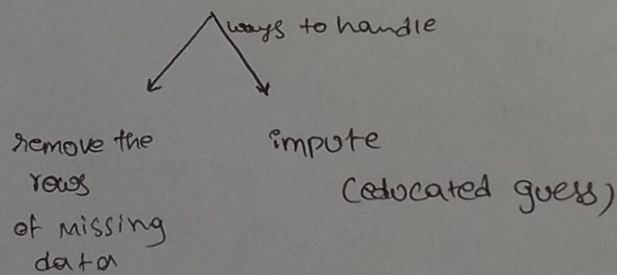
```
df.columns = ['age',  
              'some things',  
              ]
```

df.head()

keeps column names.

3) missing data (1)

like forgot to ask age



the option depends on the importance of it.

df.dtypes # gives the data types of them

Print unique values of a column)

Quiz [For practise]

10 Why do we use one-hot Encoding?

A) so that we sklearn's decision trees can use categorical data

20 sklearn decision trees can handle missing data

A) False

30 Why do we deal with missing data before splitting the full database into 'x', just the variables or features that will make classification, and 'y', the known classifications?

A) To ensure that each row in 'x' correctly corresponds to an item in 'y'.

40 What's the best way to find alpha for cost complexity pruning?

A) use cross validation

50 Why does pruning a tree improve accuracy?

A) Smaller trees ~~do not~~ ~~fix~~ do not overfit the training data

Project Name: Classification-Trees-in-Python-From-Start-to-Finish-Project

github link:

<https://github.com/Suryateja-Kolka/Classification-Trees-in-Python-From-Start-To-Finish-Project>