

ML 24 - Decision Tree Classifier With Pre-Pruning By Virat Tiwari

December 14, 2023

1 Decision Tree Pre-Pruning And Hyperparameter Tuning

```
[1]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import warnings
warnings.filterwarnings("ignore")
```

```
[2]: # We import " iris " dataset from sklearn

from sklearn.datasets import load_iris
```

```
[3]: # we store iris dataset in variable " dataset "

dataset=load_iris()
```

```
[4]: # print( dataset . DESCR ) function is used for understand the description of ↵
      ↪ dataset
      # . DESCR fucntion describe about the dataset

print(dataset.DESCR)
```

```
.. _iris_dataset:
```

Iris plants dataset

****Data Set Characteristics:****

```
:Number of Instances: 150 (50 in each of three classes)
:Number of Attributes: 4 numeric, predictive attributes and the class
:Attribute Information:
  - sepal length in cm
  - sepal width in cm
  - petal length in cm
  - petal width in cm
```

- class:
 - Iris-Setosa
 - Iris-Versicolour
 - Iris-Virginica

:Summary Statistics:

	Min	Max	Mean	SD	Class Correlation
sepal length:	4.3	7.9	5.84	0.83	0.7826
sepal width:	2.0	4.4	3.05	0.43	-0.4194
petal length:	1.0	6.9	3.76	1.76	0.9490 (high!)
petal width:	0.1	2.5	1.20	0.76	0.9565 (high!)

:Missing Attribute Values: None

:Class Distribution: 33.3% for each of 3 classes.

:Creator: R.A. Fisher

:Donor: Michael Marshall (MARSHALL%PLU@io.arc.nasa.gov)

:Date: July, 1988

The famous Iris database, first used by Sir R.A. Fisher. The dataset is taken from Fisher's paper. Note that it's the same as in R, but not as in the UCI Machine Learning Repository, which has two wrong data points.

This is perhaps the best known database to be found in the pattern recognition literature. Fisher's paper is a classic in the field and is referenced frequently to this day. (See Duda & Hart, for example.) The data set contains 3 classes of 50 instances each, where each class refers to a type of iris plant. One class is linearly separable from the other 2; the latter are NOT linearly separable from each other.

.. topic:: References

- Fisher, R.A. "The use of multiple measurements in taxonomic problems" Annual Eugenics, 7, Part II, 179-188 (1936); also in "Contributions to Mathematical Statistics" (John Wiley, NY, 1950).
- Duda, R.O., & Hart, P.E. (1973) Pattern Classification and Scene Analysis. (Q327.D83) John Wiley & Sons. ISBN 0-471-22361-1. See page 218.
- Dasarathy, B.V. (1980) "Nosing Around the Neighborhood: A New System Structure and Classification Rule for Recognition in Partially Exposed Environments". IEEE Transactions on Pattern Analysis and Machine Intelligence, Vol. PAMI-2, No. 1, 67-71.
- Gates, G.W. (1972) "The Reduced Nearest Neighbor Rule". IEEE Transactions on Information Theory, May 1972, 431-433.
- See also: 1988 MLC Proceedings, 54-64. Cheeseman et al's AUTOCLASS II conceptual clustering system finds 3 classes in the data.

- Many, many more ...

```
[5]: # Another way to import same dataset using seaborn
# we can also import the " iris " through the seaborn library

df=sns.load_dataset("iris")
```

```
[6]: # . head ( ) function gives the initial 5 datapoints from the entire dataset

df.head()
```

```
[6]:   sepal_length  sepal_width  petal_length  petal_width  species
0           5.1           3.5           1.4           0.2   setosa
1           4.9           3.0           1.4           0.2   setosa
2           4.7           3.2           1.3           0.2   setosa
3           4.6           3.1           1.5           0.2   setosa
4           5.0           3.6           1.4           0.2   setosa
```

```
[7]: # . target function is used for understanding the dataset more clearly in the
      ↪ form of 0's and 1's and so on

dataset.target
```

```
[7]: array([0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
        0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
        0, 0, 0, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
        1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
        1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 2, 2, 2, 2, 2, 2, 2, 2, 2,
        2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2,
        2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2])
```

```
[8]: # Independent and dependent features
# . iloc [ ] fucntion is used for defing the dataset that how much we want
      ↪ independent data or dependent data

x=df.iloc[:, :-1]
y=dataset.target
```

```
[9]: x
```

```
[9]:   sepal_length  sepal_width  petal_length  petal_width
0           5.1           3.5           1.4           0.2
1           4.9           3.0           1.4           0.2
2           4.7           3.2           1.3           0.2
3           4.6           3.1           1.5           0.2
4           5.0           3.6           1.4           0.2
..           ...           ...           ...           ...
```

```
[150 rows x 4 columns]
```

```
[11]: # Train Test Split
# import train_test_split function Train and Test the dataset
# train_test_split(x,y,test_size=0.33,random_state=42) - We pass x , y , test_
↳ size and random state for accuracy inside the train_test_split function
# This function require four parameters - x_train,x_test,y_train,y_test

from sklearn.model_selection import train_test_split
x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.
↳ 33,random_state=42)
```

```
[13]: from sklearn.model_selection import GridSearchCV
```

```
[15]: # TRAIN THE DATA

      clf.fit(x_train,y_train)
```

4

```

        'max_depth': [1, 2, 3, 4, 5],
        'max_features': ['auto', 'sqrt', 'log2'],
        'splitter': ['best', 'random']},
    scoring='accuracy')

```

```

[17]: # These are Pre-Prining Parametres
      # These parameters are created during the construction of tree

      clf.best_params_

```

```

[17]: {'criterion': 'log_loss',
      'max_depth': 5,
      'max_features': 'log2',
      'splitter': 'best'}

```

```

[18]: y_pred=clf.predict(x_test)

```

```

[19]: y_pred

```

```

[19]: array([1, 0, 2, 1, 1, 0, 1, 2, 1, 1, 2, 0, 0, 0, 0, 1, 2, 1, 1, 2, 0, 2,
           0, 2, 2, 2, 2, 2, 0, 0, 0, 0, 1, 0, 0, 2, 1, 0, 0, 0, 2, 1, 1, 0,
           0, 1, 1, 2, 1, 2])

```

```

[20]: # Now we import accuracy_score , classification_report for getting the accuracy_
      # or classification report of the dataset

      from sklearn.metrics import accuracy_score, classification_report

```

```

[21]: score=accuracy_score(y_pred,y_test)

```

```

[22]: score

```

```

[22]: 0.98

```

```

[23]: # Same thing we done here in case of " classification_report "

      print(classification_report(y_pred,y_test))

```

	precision	recall	f1-score	support
0	1.00	1.00	1.00	19
1	1.00	0.94	0.97	16
2	0.94	1.00	0.97	15
accuracy			0.98	50
macro avg	0.98	0.98	0.98	50

weighted avg	0.98	0.98	0.98	50
--------------	------	------	------	----

THANK YOU SO MUCH !!

YOURS VIRAT TIWARI :)