

# Decision Tree Classifier Tutorial with Python

Hello friends,

In this kernel, I build a Decision Tree Classifier to predict the safety of the car. I build two models, one with criterion `gini index` and another one with criterion `entropy`. I implement Decision Tree Classification with Python and Scikit-Learn.

**As always, I hope you find this kernel useful and your **UPVOTES** would be highly appreciated.**

# Table of Contents

1. [Introduction to Decision Tree algorithm](#)
2. [Classification and Regression Trees](#)
3. [Decision Tree algorithm terminology](#)
4. [Decision Tree algorithm intuition](#)
5. [Attribute selection measures](#)
  - [5.1 Information gain](#)
  - [5.2 Gini index](#)
6. [Overfitting in Decision-Tree algorithm](#)
7. [Import libraries](#)
8. [Import dataset](#)
9. [Exploratory data analysis](#)
10. [Declare feature vector and target variable](#)
11. [Split data into separate training and test set](#)
12. [Feature engineering](#)
13. [Decision Tree classifier with criterion gini-index](#)
14. [Decision Tree classifier with criterion entropy](#)
15. [Confusion matrix](#)
16. [Classification report](#)
17. [Results and conclusion](#)
18. [References](#)

# 1. Introduction to Decision Tree algorithm

## Table of Contents

A Decision Tree algorithm is one of the most popular machine learning algorithms. It uses a tree like structure and their possible combinations to solve a particular problem. It belongs to the class of supervised learning algorithms where it can be used for both classification and regression purposes.

A decision tree is a structure that includes a root node, branches, and leaf nodes. Each internal node denotes a test on an attribute, each branch denotes the outcome of a test, and each leaf node holds a class label. The topmost node in the tree is the root node.

We make some assumptions while implementing the Decision-Tree algorithm. These are listed below:-

1. At the beginning, the whole training set is considered as the root.
2. Feature values need to be categorical. If the values are continuous then they are discretized prior to building the model.
3. Records are distributed recursively on the basis of attribute values.
4. Order to placing attributes as root or internal node of the tree is done by using some statistical approach.

I will describe Decision Tree terminology in later section.

# 2. Classification and Regression Trees (CART)

## Table of Contents

Nowadays, Decision Tree algorithm is known by its modern name **CART** which stands for **Classification and Regression Trees**. Classification and Regression Trees or **CART** is a term introduced by Leo Breiman to refer to Decision Tree algorithms that can be used for classification and regression modeling problems.

The CART algorithm provides a foundation for other important algorithms like bagged decision trees, random forest and boosted decision trees. In this kernel, I will solve a classification problem. So, I will refer the algorithm also as Decision Tree Classification problem.

## 3. Decision Tree algorithm terminology

### Table of Contents

- In a Decision Tree algorithm, there is a tree like structure in which each internal node represents a test on an attribute, each branch represents the outcome of the test, and each leaf node represents a class label. The paths from the root node to leaf node represent classification rules.
- We can see that there is some terminology involved in Decision Tree algorithm. The terms involved in Decision Tree algorithm are as follows:-

### Root Node

- It represents the entire population or sample. This further gets divided into two or more homogeneous sets.

### Splitting

- It is a process of dividing a node into two or more sub-nodes.

### Decision Node

- When a sub-node splits into further sub-nodes, then it is called a decision node.

### Leaf/Terminal Node

- Nodes that do not split are called Leaf or Terminal nodes.

### Pruning

- When we remove sub-nodes of a decision node, this process is called pruning. It is the opposite process of splitting.

### Branch/Sub-Tree

- A sub-section of an entire tree is called a branch or sub-tree.

### Parent and Child Node

- A node, which is divided into sub-nodes is called the parent node of sub-nodes where sub-nodes are the children of a parent node.

The above terminology is represented clearly in the following diagram:-

## Decision-Tree terminology



# 4. Decision Tree algorithm intuition

## Table of Contents

The Decision-Tree algorithm is one of the most frequently and widely used supervised machine learning algorithms that can be used for both classification and regression tasks. The intuition behind the Decision-Tree algorithm is very simple to understand.

The Decision Tree algorithm intuition is as follows:-

1. For each attribute in the dataset, the Decision-Tree algorithm forms a node. The most important attribute is placed at the root node.
2. For evaluating the task in hand, we start at the root node and we work our way down the tree by following the corresponding node that meets our condition or decision.
3. This process continues until a leaf node is reached. It contains the prediction or the outcome of the Decision Tree.

# 5. Attribute selection measures

## Table of Contents

The primary challenge in the Decision Tree implementation is to identify the attributes which we consider as the root node and each level. This process is known as the **attributes selection**. There are different attributes selection measure to identify the attribute which can be considered as the root node at each level.

There are 2 popular attribute selection measures. They are as follows:-

- **Information gain**
- **Gini index**

While using **Information gain** as a criterion, we assume attributes to be categorical and for **Gini index** attributes are assumed to be continuous. These attribute selection measures are described below.

## 5.1 Information gain

### [Table of Contents](#)

By using information gain as a criterion, we try to estimate the information contained by each attribute. To understand the concept of Information Gain, we need to know another concept called **Entropy**.

## Entropy

Entropy measures the impurity in the given dataset. In Physics and Mathematics, entropy is referred to as the randomness or uncertainty of a random variable  $X$ . In information theory, it refers to the impurity in a group of examples. **Information gain** is the decrease in entropy. Information gain computes the difference between entropy before split and average entropy after split of the dataset based on given attribute values.

Entropy is represented by the following formula:-



Here,  $c$  is the number of classes and  $p_i$  is the probability associated with the  $i$ th class.

The ID3 (Iterative Dichotomiser) Decision Tree algorithm uses entropy to calculate information gain. So, by calculating decrease in **entropy measure** of each attribute we can calculate their information gain. The attribute with the highest information gain is chosen as the splitting attribute at the node.

## 5.2 Gini index

### [Table of Contents](#)

Another attribute selection measure that **CART (Categorical and Regression Trees)** uses is the **Gini index**. It uses the Gini method to create split points.

Gini index can be represented with the following diagram:-

## Gini index

$$Gini = 1 - \sum_{i=1}^C (p_i)^2$$

---

Here, again **c** is the number of classes and **p<sub>i</sub>** is the probability associated with the *i*th class.

Gini index says, if we randomly select two items from a population, they must be of the same class and probability for this is 1 if the population is pure.

It works with the categorical target variable “Success” or “Failure”. It performs only binary splits. The higher the value of Gini, higher the homogeneity. CART (Classification and Regression Tree) uses the Gini method to create binary splits.

Steps to Calculate Gini for a split

1. Calculate Gini for sub-nodes, using formula sum of the square of probability for success and failure ( $p^2 + q^2$ ).
2. Calculate Gini for split using weighted Gini score of each node of that split.

In case of a discrete-valued attribute, the subset that gives the minimum gini index for that chosen is selected as a splitting attribute. In the case of continuous-valued attributes, the strategy is to select each pair of adjacent values as a possible split-point and point with smaller gini index chosen as the splitting point. The attribute with minimum Gini index is chosen as the splitting attribute.

## 6. Overfitting in Decision Tree algorithm

### [Table of Contents](#)

Overfitting is a practical problem while building a Decision-Tree model. The problem of overfitting is considered when the algorithm continues to go deeper and deeper to reduce the training-set error but results with an increased test-set error. So, accuracy of prediction for our model goes down. It generally happens when we build many branches due to outliers and irregularities in data.

Two approaches which can be used to avoid overfitting are as follows:-

- Pre-Pruning
- Post-Pruning

### Pre-Pruning

In pre-pruning, we stop the tree construction a bit early. We prefer not to split a node if its goodness measure is below a threshold value. But it is difficult to choose an appropriate stopping point.

### Post-Pruning

In post-pruning, we go deeper and deeper in the tree to build a complete tree. If the tree shows the overfitting problem then pruning is done as a post-pruning step. We use the cross-validation data to check the effect of our pruning. Using cross-validation data, we test whether expanding a node will result in improve or not. If it shows an improvement, then we can continue by expanding that node. But if it shows a reduction in accuracy then it should not be expanded. So, the node should be converted to a leaf node.

## 7. Import libraries

### [Table of Contents](#)



In [1]:

```
# This Python 3 environment comes with many helpful analytics libraries installed
# It is defined by the kaggle/python docker image: https://github.com/kaggle/docker-python
# For example, here's several helpful packages to load in

import numpy as np # linear algebra
import pandas as pd # data processing, CSV file I/O (e.g. pd.read_csv)
import matplotlib.pyplot as plt # data visualization
import seaborn as sns # statistical data visualization
%matplotlib inline

# Input data files are available in the "../input/" directory.
# For example, running this (by clicking run or pressing Shift+Enter) will list all files under the input directory

import os
for dirname, _, filenames in os.walk('/kaggle/input'):
    for filename in filenames:
        print(os.path.join(dirname, filename))

# Any results you write to the current directory are saved as output.
```

/kaggle/input/car-evaluation-data-set/car\_evaluation.csv

In [2]:

```
import warnings

warnings.filterwarnings('ignore')
```

## 8. Import dataset

[Table of Contents](#)

In [3]:

```
data = '/kaggle/input/car-evaluation-data-set/car_evaluation.csv'  
  
df = pd.read_csv(data, header=None)
```

## 9. Exploratory data analysis

### [Table of Contents](#)

Now, I will explore the data to gain insights about the data.

In [4]:

```
# view dimensions of dataset  
  
df.shape
```

Out[4]:

```
(1728, 7)
```

We can see that there are 1728 instances and 7 variables in the data set.

### View top 5 rows of dataset

In [5]:

```
# preview the dataset
```

```
df.head()
```

Out[5]:

	0	1	2	3	4	5	6
0	vhigh	vhigh	2	2	small	low	unacc
1	vhigh	vhigh	2	2	small	med	unacc
2	vhigh	vhigh	2	2	small	high	unacc
3	vhigh	vhigh	2	2	med	low	unacc
4	vhigh	vhigh	2	2	med	med	unacc

## Rename column names

We can see that the dataset does not have proper column names. The columns are merely labelled as 0,1,2.... and so on. We should give proper names to the columns. I will do it as follows:-

In [6]:

```
col_names = ['buying', 'maint', 'doors', 'persons', 'lug_boot', 'safety', 'class']
```

```
df.columns = col_names
```

```
col_names
```

Out[6]:

```
['buying', 'maint', 'doors', 'persons', 'lug_boot', 'safety', 'class']
```

In [7]:

```
# let's again preview the dataset  
  
df.head()
```

Out[7]:

	buying	maint	doors	persons	lug_boot	safety	class
0	vhigh	vhigh	2	2	small	low	unacc
1	vhigh	vhigh	2	2	small	med	unacc
2	vhigh	vhigh	2	2	small	high	unacc
3	vhigh	vhigh	2	2	med	low	unacc
4	vhigh	vhigh	2	2	med	med	unacc

We can see that the column names are renamed. Now, the columns have meaningful names.

## View summary of dataset

In [8]:

```
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>  
RangeIndex: 1728 entries, 0 to 1727  
Data columns (total 7 columns):  
buying      1728 non-null object  
maint       1728 non-null object  
doors       1728 non-null object  
persons     1728 non-null object  
lug_boot    1728 non-null object  
safety      1728 non-null object  
class       1728 non-null object  
dtypes: object(7)  
memory usage: 94.6+ KB
```

## Frequency distribution of values in variables

Now, I will check the frequency counts of categorical variables.

In [9]:

```
col_names = ['buying', 'maint', 'doors', 'persons', 'lug_boot', 'safety', 'class']

for col in col_names:

    print(df[col].value_counts())
```

```
low      432
high     432
med      432
vhigh    432
Name: buying, dtype: int64
low      432
high     432
med      432
vhigh    432
Name: maint, dtype: int64
5more    432
2         432
3         432
4         432
Name: doors, dtype: int64
2         576
more      576
4         576
Name: persons, dtype: int64
med       576
small     576
big       576
Name: lug_boot, dtype: int64
low       576
high      576
med       576
Name: safety, dtype: int64
unacc    1210
acc       384
good      69
vgood     65
Name: class, dtype: int64
```

We can see that the `doors` and `persons` are categorical in nature. So, I will treat them as categorical variables.

## Summary of variables

- There are 7 variables in the dataset. All the variables are of categorical data type.
- These are given by `buying`, `maint`, `doors`, `persons`, `lug_boot`, `safety` and `class`.
- `class` is the target variable.

## Explore `class` variable

In [10]:

```
df['class'].value_counts()
```

Out[10]:

```
unacc    1210
acc       384
good       69
vgood     65
Name: class, dtype: int64
```

The `class` target variable is ordinal in nature.

## Missing values in variables

In [11]:

```
# check missing values in variables  
  
df.isnull().sum()
```

Out[11]:

```
buying      0  
maint      0  
doors      0  
persons    0  
lug_boot   0  
safety     0  
class      0  
dtype: int64
```

We can see that there are no missing values in the dataset. I have checked the frequency distribution of values previously. It also confirms that there are no missing values in the dataset.

## 10. Declare feature vector and target variable

[Table of Contents](#)

In [12]:

```
X = df.drop(['class'], axis=1)  
  
y = df['class']
```

## 11. Split data into separate training and test set

[Table of Contents](#)



In [13]:

```
# split X and y into training and testing sets

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)
```

In [14]:

```
# check the shape of X_train and X_test

X_train.shape, X_test.shape
```

Out[14]:

```
((1157, 6), (571, 6))
```

## 12. Feature Engineering

### [Table of Contents](#)

**Feature Engineering** is the process of transforming raw data into useful features that help us to understand our model better and increase its predictive power. I will carry out feature engineering on different types of variables.

First, I will check the data types of variables again.

In [15]:

```
# check data types in X_train
```

```
X_train.dtypes
```

Out[15]:

```
buying      object
maint       object
doors       object
persons     object
lug_boot    object
safety      object
dtype: object
```

## Encode categorical variables

Now, I will encode the categorical variables.

In [16]:

```
X_train.head()
```

Out[16]:

	buying	maint	doors	persons	lug_boot	safety
48	vhigh	vhigh	3	more	med	low
468	high	vhigh	3	4	small	low
155	vhigh	high	3	more	small	high
1721	low	low	5more	more	small	high
1208	med	low	2	more	small	high

We can see that all the variables are ordinal categorical data type.

In [17]:

```
# import category encoders

import category_encoders as ce
```

In [18]:

```
# encode variables with ordinal encoding

encoder = ce.OrdinalEncoder(cols=['buying', 'maint', 'doors', 'persons', 'lug_boot', 'safety'])

X_train = encoder.fit_transform(X_train)

X_test = encoder.transform(X_test)
```

In [19]:

```
X_train.head()
```

Out[19]:

	buying	maint	doors	persons	lug_boot	safety
48	1	1	1	1	1	1
468	2	1	1	2	2	1
155	1	2	1	1	2	2
1721	3	3	2	1	2	2
1208	4	3	3	1	2	2

In [20]:

```
X_test.head()
```

Out[20]:

	buying	maint	doors	persons	lug_boot	safety
599	2	2	4	3	1	2
1201	4	3	3	2	1	3
628	2	2	2	3	3	3
1498	3	2	2	2	1	3
1263	4	3	4	1	1	1

We now have training and test set ready for model building.

## 13. Decision Tree Classifier with criterion gini index

[Table of Contents](#)

In [21]:

```
# import DecisionTreeClassifier

from sklearn.tree import DecisionTreeClassifier
```

In [22]:

```
# instantiate the DecisionTreeClassifier model with criterion gini index

clf_gini = DecisionTreeClassifier(criterion='gini', max_depth=3, random_state=0)

# fit the model
clf_gini.fit(X_train, y_train)
```

Out[22]:

```
DecisionTreeClassifier(class_weight=None, criterion='gini', max_depth=3,
                        max_features=None, max_leaf_nodes=None,
                        min_impurity_decrease=0.0, min_impurity_split=None,
                        min_samples_leaf=1, min_samples_split=2,
                        min_weight_fraction_leaf=0.0, presort=False,
                        random_state=0, splitter='best')
```

## Predict the Test set results with criterion gini index

In [23]:

```
y_pred_gini = clf_gini.predict(X_test)
```

## Check accuracy score with criterion gini index

In [24]:

```
from sklearn.metrics import accuracy_score

print('Model accuracy score with criterion gini index: {0:0.4f}'.format(accuracy_score(y_test, y_pred_gini)))
```

Model accuracy score with criterion gini index: 0.8021

Here, **y\_test** are the true class labels and **y\_pred\_gini** are the predicted class labels in the test-set.

## Compare the train-set and test-set accuracy

Now, I will compare the train-set and test-set accuracy to check for overfitting.

In [25]:

```
y_pred_train_gini = clf_gini.predict(X_train)

y_pred_train_gini
```

Out[25]:

```
array(['unacc', 'unacc', 'unacc', ..., 'unacc', 'unacc', 'acc'],
      dtype=object)
```

In [26]:

```
print('Training-set accuracy score: {0:0.4f}'.format(accuracy_score(y_train, y_
pred_train_gini)))
```

Training-set accuracy score: 0.7865

## Check for overfitting and underfitting

In [27]:

```
# print the scores on training and test set

print('Training set score: {:.4f}'.format(clf_gini.score(X_train, y_train)))

print('Test set score: {:.4f}'.format(clf_gini.score(X_test, y_test)))
```

Training set score: 0.7865

Test set score: 0.8021

Here, the training-set accuracy score is 0.7865 while the test-set accuracy to be 0.8021. These two values are quite comparable. So, there is no sign of overfitting.

## Visualize decision-trees

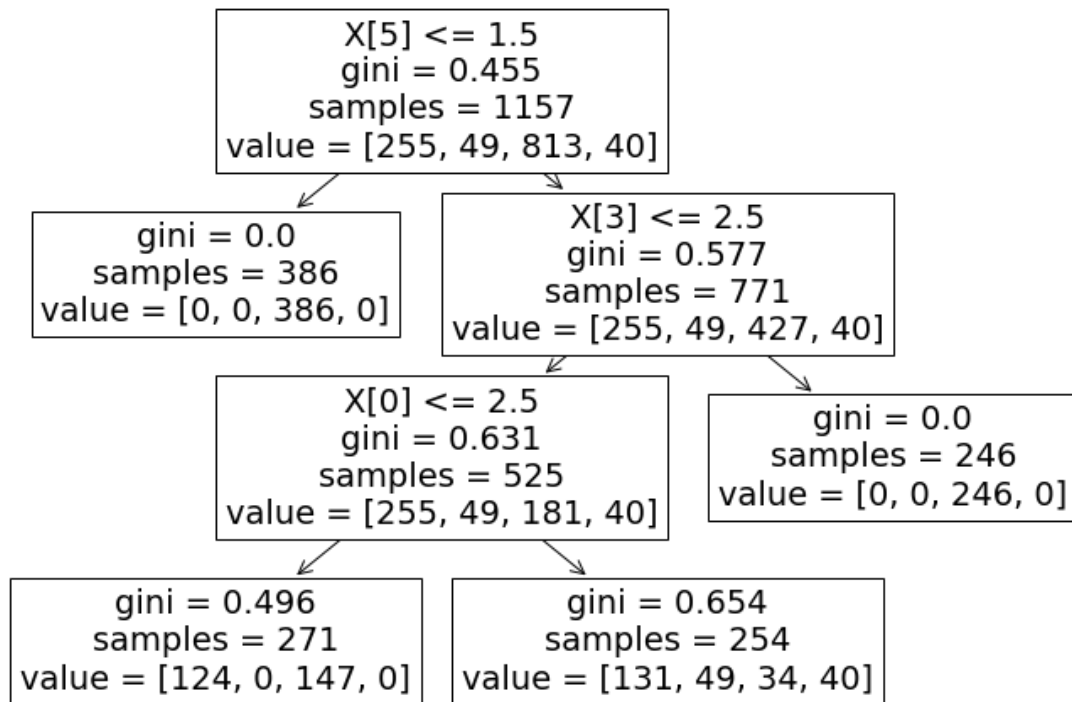
In [28]:

```
plt.figure(figsize=(12,8))  
  
from sklearn import tree  
  
tree.plot_tree(clf_gini.fit(X_train, y_train))
```



Out[28]:

```
[Text(267.84000000000003, 380.52, 'X[5] <= 1.5\ngini = 0.455\nsamples = 1157\nvalue = [255, 49, 813, 40]'),  
Text(133.92000000000002, 271.8, 'gini = 0.0\nsamples = 386\nvalue = [0, 0, 386, 0]'),  
Text(401.76000000000005, 271.8, 'X[3] <= 2.5\ngini = 0.577\nsamples = 771\nvalue = [255, 49, 427, 40]'),  
Text(267.84000000000003, 163.07999999999998, 'X[0] <= 2.5\ngini = 0.631\nsamples = 525\nvalue = [255, 49, 181, 40]'),  
Text(133.92000000000002, 54.360000000000014, 'gini = 0.496\nsamples = 271\nvalue = [124, 0, 147, 0]'),  
Text(401.76000000000005, 54.360000000000014, 'gini = 0.654\nsamples = 254\nvalue = [131, 49, 34, 40]'),  
Text(535.68000000000001, 163.07999999999998, 'gini = 0.0\nsamples = 246\nvalue = [0, 0, 246, 0]')]
```



Visualize decision-trees with graphviz

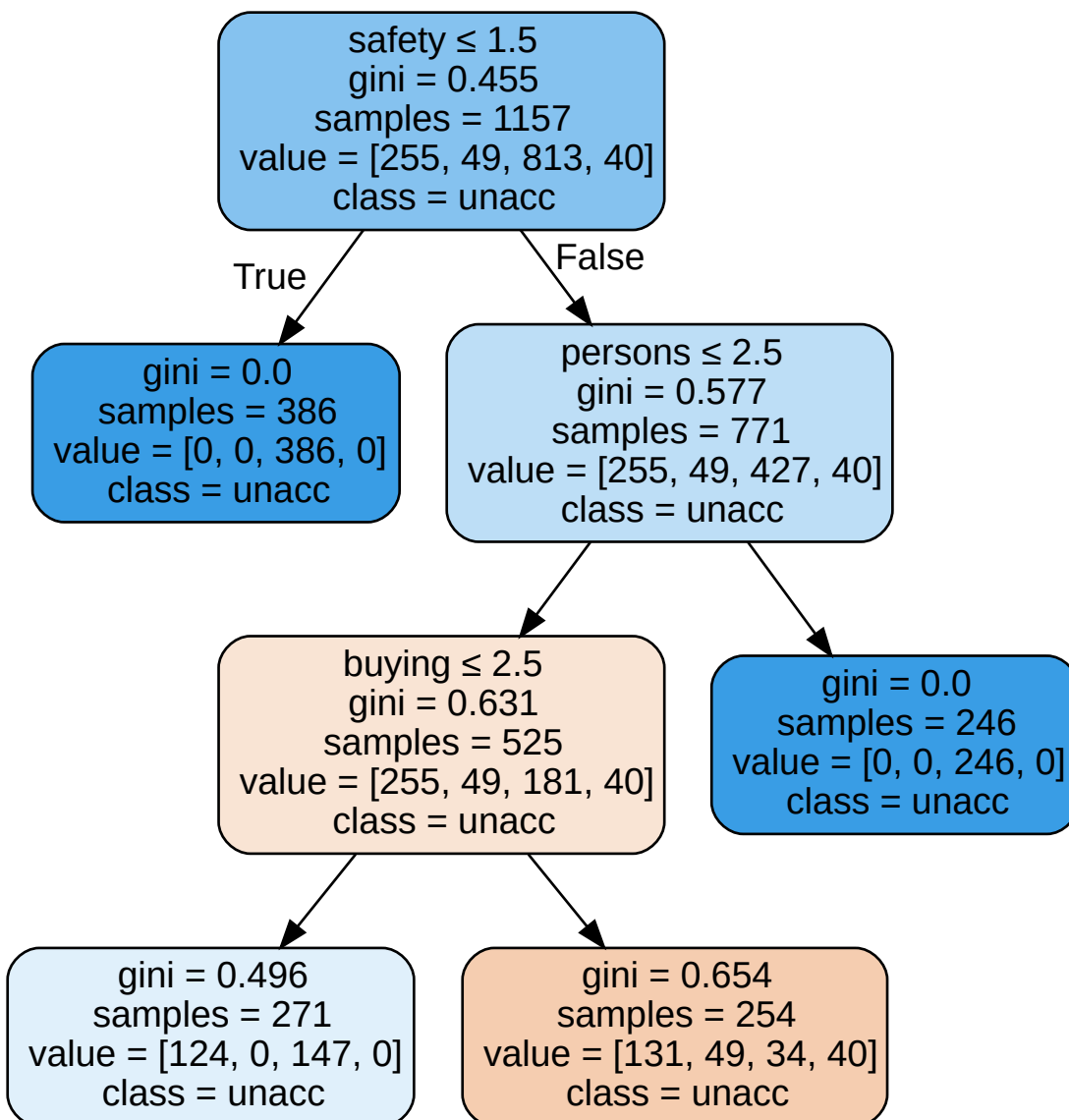
In [29]:

```
import graphviz
dot_data = tree.export_graphviz(clf_gini, out_file=None,
                                feature_names=X_train.columns,
                                class_names=y_train,
                                filled=True, rounded=True,
                                special_characters=True)

graph = graphviz.Source(dot_data)

graph
```

Out[29]:



## 14. Decision Tree Classifier with criterion entropy

### [Table of Contents](#)

In [30]:

```
# instantiate the DecisionTreeClassifier model with criterion entropy

clf_en = DecisionTreeClassifier(criterion='entropy', max_depth=3, random_state=0
)

# fit the model
clf_en.fit(X_train, y_train)
```

Out[30]:

```
DecisionTreeClassifier(class_weight=None, criterion='entropy', max_
depth=3,
                        max_features=None, max_leaf_nodes=None,
                        min_impurity_decrease=0.0, min_impurity_spli
t=None,
                        min_samples_leaf=1, min_samples_split=2,
                        min_weight_fraction_leaf=0.0, presort=False,
                        random_state=0, splitter='best')
```

### Predict the Test set results with criterion entropy

In [31]:

```
y_pred_en = clf_en.predict(X_test)
```

### Check accuracy score with criterion entropy

In [32]:

```
from sklearn.metrics import accuracy_score

print('Model accuracy score with criterion entropy: {0:0.4f}'.format(accuracy_score(y_test, y_pred_en)))
```

Model accuracy score with criterion entropy: 0.8021

## Compare the train-set and test-set accuracy

Now, I will compare the train-set and test-set accuracy to check for overfitting.

In [33]:

```
y_pred_train_en = clf_en.predict(X_train)

y_pred_train_en
```

Out[33]:

```
array(['unacc', 'unacc', 'unacc', ..., 'unacc', 'unacc', 'acc'],
      dtype=object)
```

In [34]:

```
print('Training-set accuracy score: {0:0.4f}'.format(accuracy_score(y_train, y_pred_train_en)))
```

Training-set accuracy score: 0.7865

## Check for overfitting and underfitting

In [35]:

```
# print the scores on training and test set

print('Training set score: {:.4f}'.format(clf_en.score(X_train, y_train)))

print('Test set score: {:.4f}'.format(clf_en.score(X_test, y_test)))
```

Training set score: 0.7865

Test set score: 0.8021

We can see that the training-set score and test-set score is same as above. The training-set accuracy score is 0.7865 while the test-set accuracy to be 0.8021. These two values are quite comparable. So, there is no sign of overfitting.

## Visualize decision-trees

In [36]:

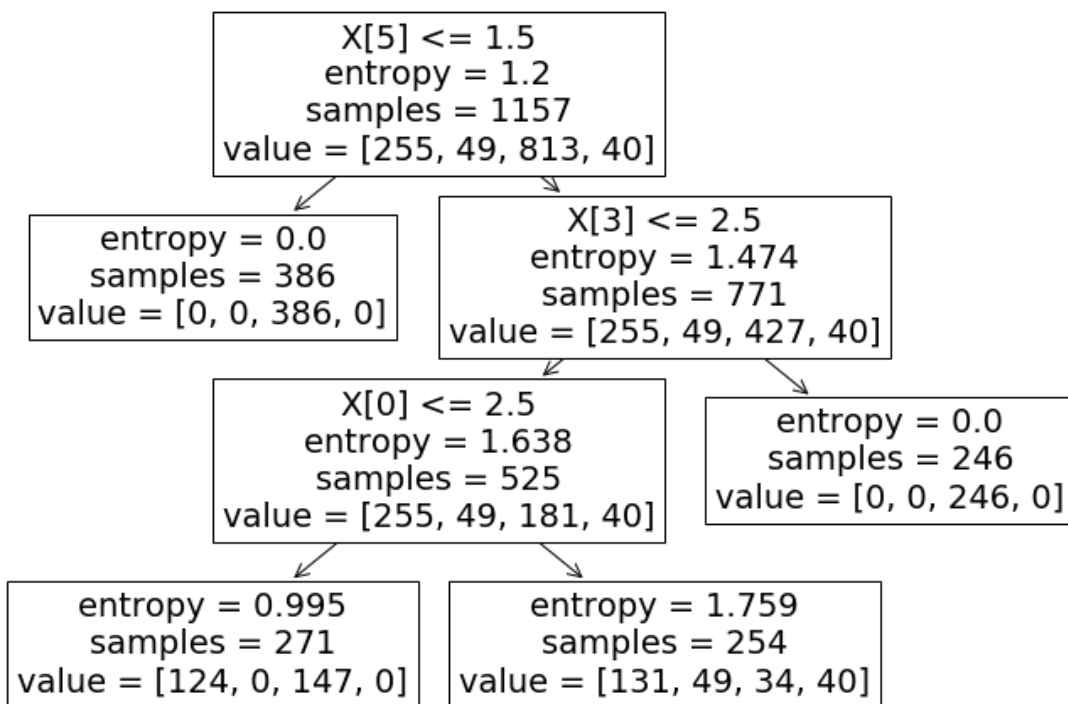
```
plt.figure(figsize=(12,8))

from sklearn import tree

tree.plot_tree(clf_en.fit(X_train, y_train))
```

Out[36]:

```
[Text(267.84000000000003, 380.52, 'X[5] <= 1.5\nentropy = 1.2\nsamples = 1157\nvalue = [255, 49, 813, 40]'),
 Text(133.92000000000002, 271.8, 'entropy = 0.0\nsamples = 386\nvalue = [0, 0, 386, 0]'),
 Text(401.76000000000005, 271.8, 'X[3] <= 2.5\nentropy = 1.474\nsamples = 771\nvalue = [255, 49, 427, 40]'),
 Text(267.84000000000003, 163.07999999999998, 'X[0] <= 2.5\nentropy = 1.638\nsamples = 525\nvalue = [255, 49, 181, 40]'),
 Text(133.92000000000002, 54.360000000000014, 'entropy = 0.995\nsamples = 271\nvalue = [124, 0, 147, 0]'),
 Text(401.76000000000005, 54.360000000000014, 'entropy = 1.759\nsamples = 254\nvalue = [131, 49, 34, 40]'),
 Text(535.68000000000001, 163.07999999999998, 'entropy = 0.0\nsamples = 246\nvalue = [0, 0, 246, 0]')]
```



## Visualize decision-trees with graphviz

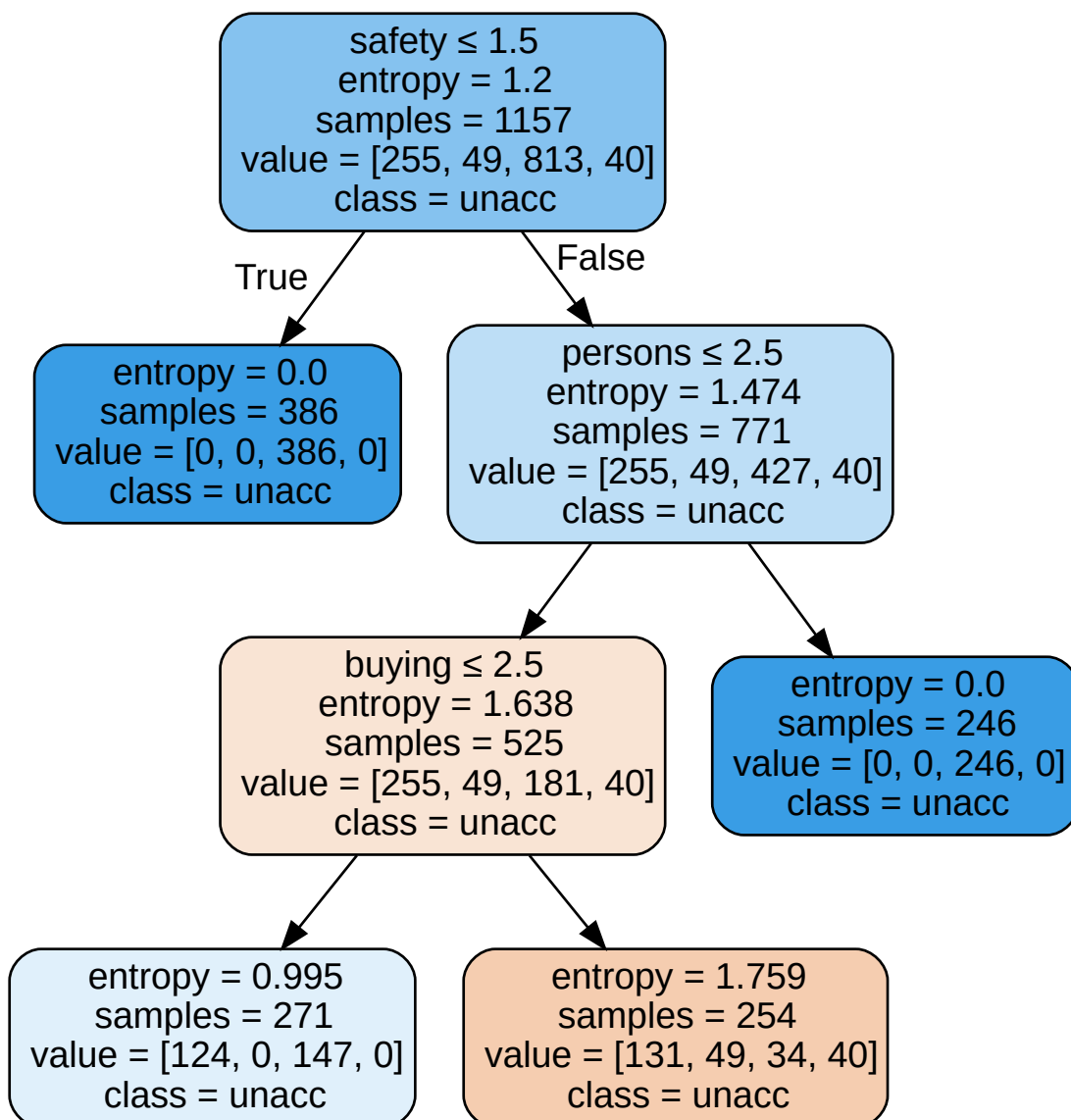
In [37]:

```
import graphviz
dot_data = tree.export_graphviz(clf_en, out_file=None,
                                feature_names=X_train.columns,
                                class_names=y_train,
                                filled=True, rounded=True,
                                special_characters=True)

graph = graphviz.Source(dot_data)

graph
```

Out[37]:





Now, based on the above analysis we can conclude that our classification model accuracy is very good. Our model is doing a very good job in terms of predicting the class labels.

But, it does not give the underlying distribution of values. Also, it does not tell anything about the type of errors our classifier is making.

We have another tool called `Confusion matrix` that comes to our rescue.

## 15. Confusion matrix

### [Table of Contents](#)

A confusion matrix is a tool for summarizing the performance of a classification algorithm. A confusion matrix will give us a clear picture of classification model performance and the types of errors produced by the model. It gives us a summary of correct and incorrect predictions broken down by each category. The summary is represented in a tabular form.

Four types of outcomes are possible while evaluating a classification model performance. These four outcomes are described below:-

**True Positives (TP)** – True Positives occur when we predict an observation belongs to a certain class and the observation actually belongs to that class.

**True Negatives (TN)** – True Negatives occur when we predict an observation does not belong to a certain class and the observation actually does not belong to that class.

**False Positives (FP)** – False Positives occur when we predict an observation belongs to a certain class but the observation actually does not belong to that class. This type of error is called **Type I error**.

**False Negatives (FN)** – False Negatives occur when we predict an observation does not belong to a certain class but the observation actually belongs to that class. This is a very serious error and it is called **Type II error**.

These four outcomes are summarized in a confusion matrix given below.

In [38]:

```
# Print the Confusion Matrix and slice it into four pieces

from sklearn.metrics import confusion_matrix

cm = confusion_matrix(y_test, y_pred_en)

print('Confusion matrix\n\n', cm)
```

Confusion matrix

```
[[ 73   0  56   0]
 [ 20   0   0   0]
 [ 12   0 385   0]
 [ 25   0   0   0]]
```

## 16. Classification Report

### [Table of Contents](#)

**Classification report** is another way to evaluate the classification model performance. It displays the **precision**, **recall**, **f1** and **support** scores for the model. I have described these terms in later.

We can print a classification report as follows:-

In [39]:

```
from sklearn.metrics import classification_report

print(classification_report(y_test, y_pred_en))
```

	precision	recall	f1-score	support
acc	0.56	0.57	0.56	129
good	0.00	0.00	0.00	20
unacc	0.87	0.97	0.92	397
vgood	0.00	0.00	0.00	25
accuracy			0.80	571
macro avg	0.36	0.38	0.37	571
weighted avg	0.73	0.80	0.77	571

## 17. Results and conclusion

### Table of Contents

1. In this project, I build a Decision-Tree Classifier model to predict the safety of the car. I build two models, one with criterion `gini index` and another one with criterion `entropy`. The model yields a very good performance as indicated by the model accuracy in both the cases which was found to be 0.8021.
2. In the model with criterion `gini index`, the training-set accuracy score is 0.7865 while the test-set accuracy to be 0.8021. These two values are quite comparable. So, there is no sign of overfitting.
3. Similarly, in the model with criterion `entropy`, the training-set accuracy score is 0.7865 while the test-set accuracy to be 0.8021. We get the same values as in the case with criterion `gini`. So, there is no sign of overfitting.
4. In both the cases, the training-set and test-set accuracy score is the same. It may happen because of small dataset.
5. The confusion matrix and classification report yields very good model performance.

# 18. References

## Table of Contents

The work done in this project is inspired from following books and websites:-

1. Hands on Machine Learning with Scikit-Learn and Tensorflow by Aurélien Geron
2. Introduction to Machine Learning with Python by Andreas C. Müller and Sarah Guido
3. [https://en.wikipedia.org/wiki/Decision\\_tree](https://en.wikipedia.org/wiki/Decision_tree) ([https://en.wikipedia.org/wiki/Decision\\_tree](https://en.wikipedia.org/wiki/Decision_tree))
4. [https://en.wikipedia.org/wiki/Information\\_gain\\_in\\_decision\\_trees](https://en.wikipedia.org/wiki/Information_gain_in_decision_trees)  
([https://en.wikipedia.org/wiki/Information\\_gain\\_in\\_decision\\_trees](https://en.wikipedia.org/wiki/Information_gain_in_decision_trees))
5. [https://en.wikipedia.org/wiki/Entropy\\_\(information\\_theory](https://en.wikipedia.org/wiki/Entropy_(information_theory))  
([https://en.wikipedia.org/wiki/Entropy\\_\(information\\_theory\)\)](https://en.wikipedia.org/wiki/Entropy_(information_theory))))
6. <https://www.datacamp.com/community/tutorials/decision-tree-classification-python>  
(<https://www.datacamp.com/community/tutorials/decision-tree-classification-python>)
7. <https://stackabuse.com/decision-trees-in-python-with-scikit-learn/>  
(<https://stackabuse.com/decision-trees-in-python-with-scikit-learn/>)
8. <https://acadgild.com/blog/decision-tree-python-code> (<https://acadgild.com/blog/decision-tree-python-code>)

So, now we will come to the end of this kernel.

I hope you find this kernel useful and enjoyable.

Your comments and feedback are most welcome.

Thank you

[Go to Top](#)