

DEVYANI BEOHAR

J068

EXP 9 Decision Tree with Cross Validation and GridSearchCV

In [1]:

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import os
print(os.getcwd())
%matplotlib inline
```

C:\Users\USER

In [2]:

```
df=pd.read_csv("C:\\Users\\USER\\Desktop\\Data Science\\SEM 5\\ML\\car_evaluation.csv",h
eader=None)
df.head()
```

Out[2]:

	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

In [3]:

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

Out[3]:

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

In [4]:

```
df.head()
```

Out[4]:

	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

In [5]:

```
for i in col_names:
```

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

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

In [6]:

```
df.shape
```

Out[6]:

```
(1728, 7)
```

In [7]:

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

In [8]:

```
from sklearn.model_selection import train_test_split
X_train,X_test,y_train,y_test=train_test_split(X,y,test_size=0.3,random_state=42)
```

In [10]:

```
from sklearn.preprocessing import OrdinalEncoder
enc=OrdinalEncoder()
X_train = enc.fit_transform(X_train)
X_test = enc.transform((X_test))
```

Gini index as criterion

In [11]:

```
from sklearn.tree import DecisionTreeClassifier
```

In [12]:

```
clf_gini = DecisionTreeClassifier(criterion='gini',max_depth=3,random_state=42)
```

```
clf_gini.fit(X_train,y_train)
```

Out[12]:

```
DecisionTreeClassifier(max_depth=3, random_state=42)
```

In [13]:

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

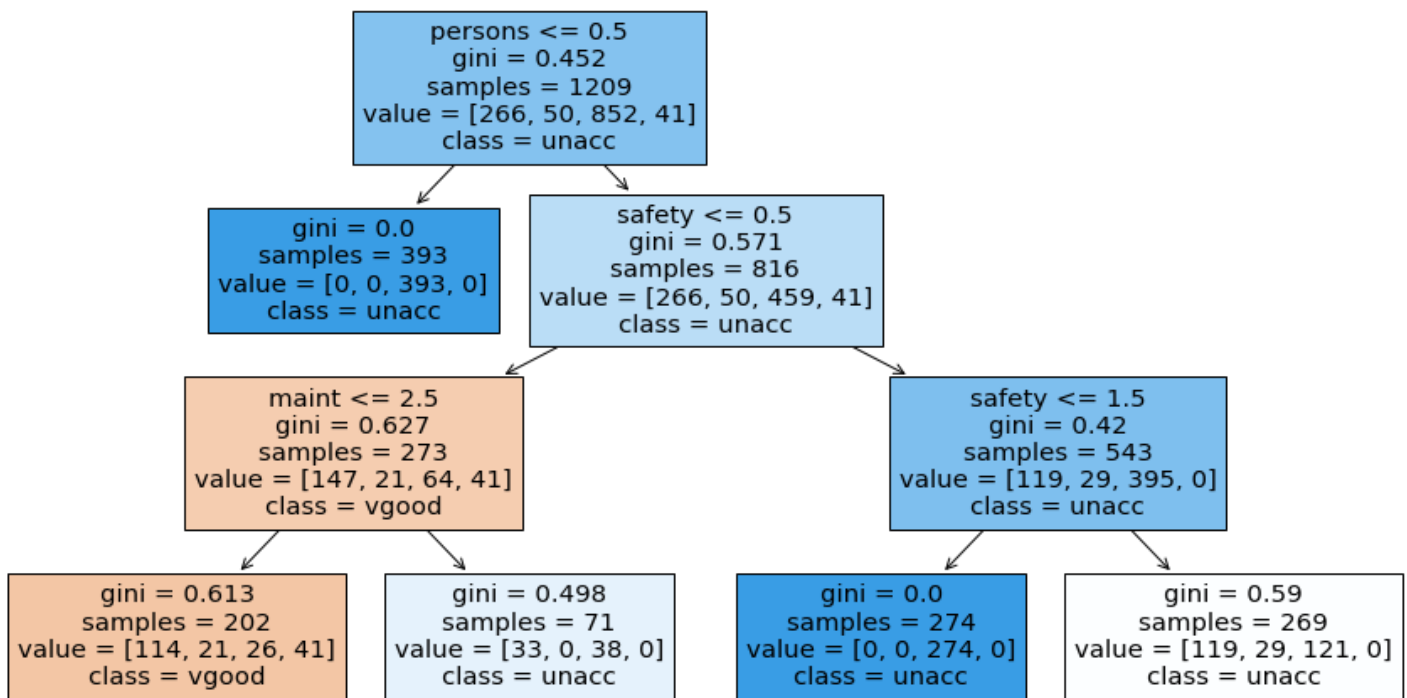
In [14]:

```
from sklearn.metrics import accuracy_score
print(f'Model with gini index gives an accuracy of: {accuracy_score(y_test, y_pred)}')
```

Model with gini index gives an accuracy of: 0.7572254335260116

In [15]:

```
from sklearn import tree
plt.figure(figsize=(15,8))
tree.plot_tree(clf_gini,feature_names=['buying','maint','doors','persons','lug_boot','safety'],
               class_names=list(set(y_train)),filled=True)
plt.show()
```



In [17]:

```
print(f'Training set score: {clf_gini.score(X_train,y_train)}')
print(f'Test set score: {clf_gini.score(X_test,y_test)}')
```

Training set score: 0.7775020678246485

Test set score: 0.7572254335260116

Entropy as criterion

In [18]:

```
clf_entropy=DecisionTreeClassifier(criterion='entropy',max_depth=3,random_state=42)
clf_entropy.fit(X_train,y_train)
```

Out[18]:

```
DecisionTreeClassifier(criterion='entropy', max_depth=3, random_state=42)
```

In [19]:

```
y_pred=clf_entropy.predict(X_test)
```

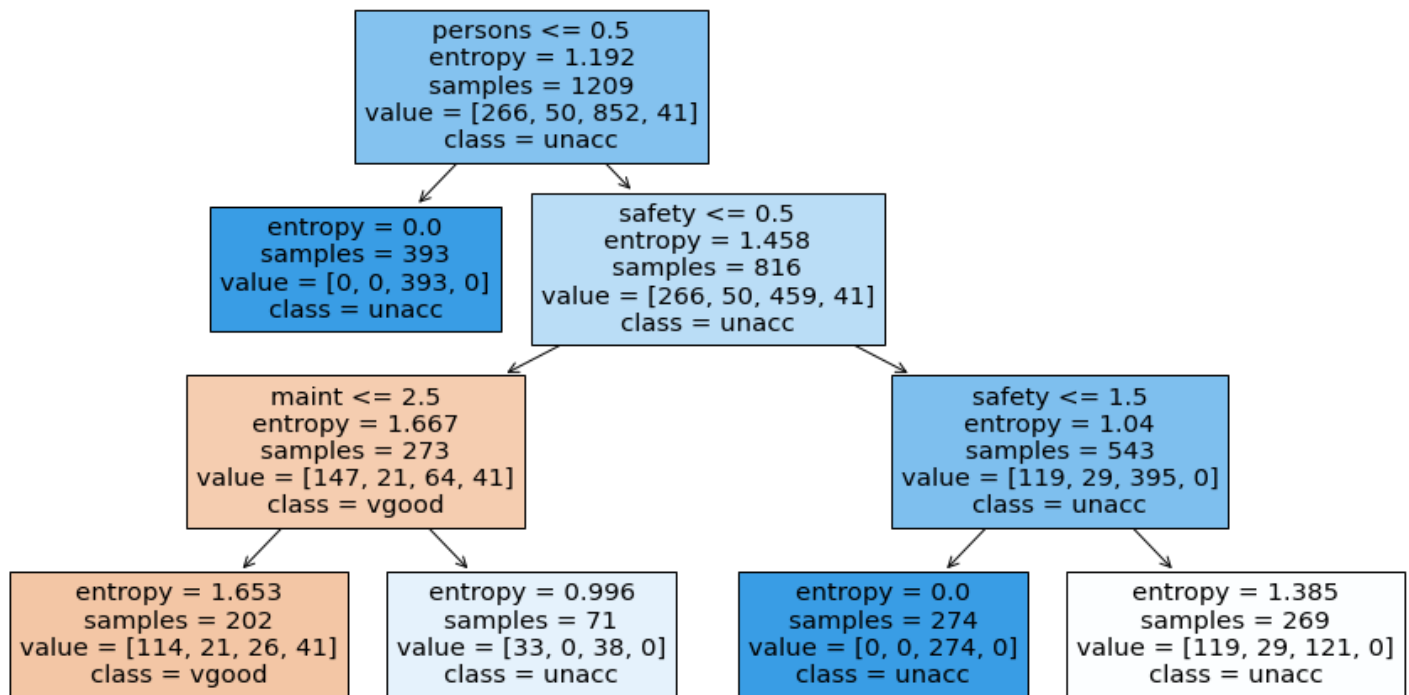
In [20]:

```
from sklearn.metrics import accuracy_score
print(f'Model with gini index gives an accuracy of: {accuracy_score(y_test, y_pred)}')
```

Model with gini index gives an accuracy of: 0.7572254335260116

In [21]:

```
plt.figure(figsize=(15,8))
tree.plot_tree(clf_entropy,feature_names=['buying','maint','doors','persons','lug_boot','safety'],
               class_names=list(set(y_train)),filled=True)
plt.show()
```



In [22]:

```
print(f'Training set score: {clf_entropy.score(X_train,y_train)}')
print(f'Test set score: {clf_entropy.score(X_test,y_test)}')
```

Training set score: 0.7775020678246485

Test set score: 0.7572254335260116

Grid Search CV

In [23]:

```
from sklearn.model_selection import GridSearchCV
```

In [24]:

```
params={'criterion':['gini','entropy'],'max_depth':list(range(3,7)),'min_samples_split':list(range(3,7)),'min_samples_leaf':list(range(3,7)),'max_leaf_nodes':list(range(3,12))}
```

In [25]:

```
decision_tree=DecisionTreeClassifier()
dt=GridSearchCV(decision_tree,params,cv=10,scoring='accuracy')
dt.fit(X_train,y_train)
```

```
dt.best_score_
```

```
Out[25]:
```

```
0.8461088154269971
```

```
In [26]:
```

```
dt.best_params_
```

```
Out[26]:
```

```
{'criterion': 'entropy',  
 'max_depth': 6,  
 'max_leaf_nodes': 11,  
 'min_samples_leaf': 3,  
 'min_samples_split': 3}
```

```
In [27]:
```

```
y_pred=dt.predict(X_test)
```

```
In [28]:
```

```
print(f'Model with Decision Tree gives an accuracy of: {accuracy_score(y_test, y_pred)}')
```

```
Model with Decision Tree gives an accuracy of: 0.861271676300578
```

After GridSearchCV

```
In [29]:
```

```
dt=DecisionTreeClassifier(criterion='entropy',max_depth=6,max_leaf_nodes=11,min_samples_1  
eaf=3,  
                          min_samples_split=3,splitter='best')
```

```
In [30]:
```

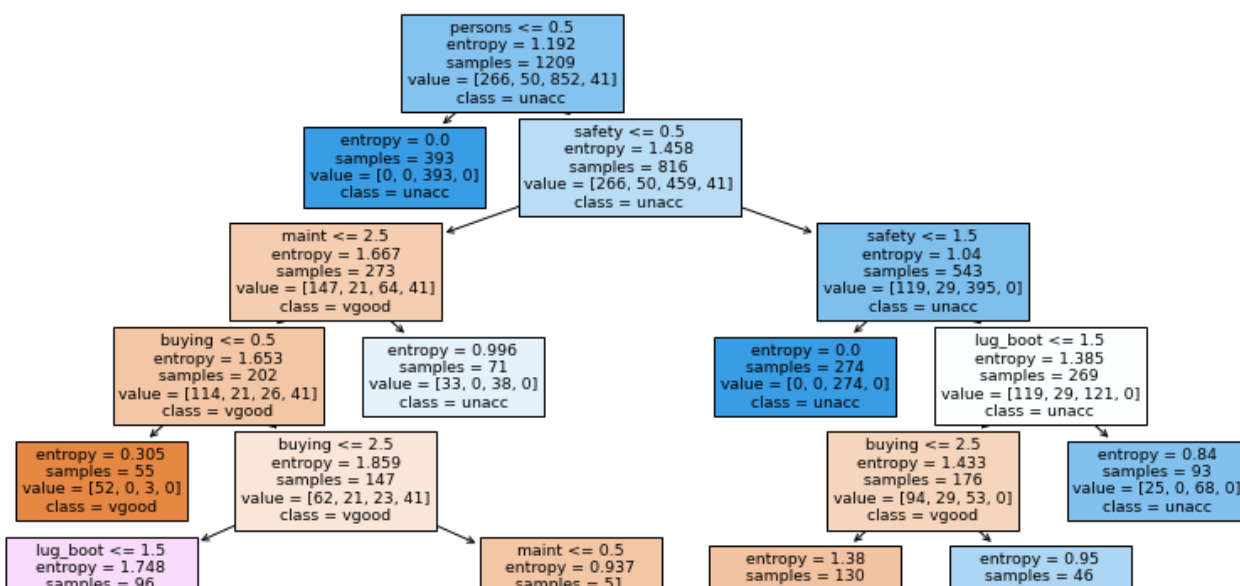
```
dt.fit(X_train,y_train)
```

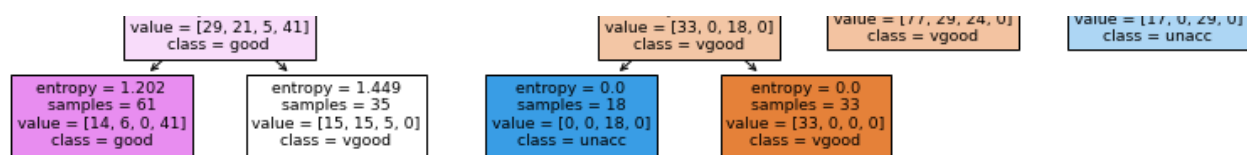
```
Out[30]:
```

```
DecisionTreeClassifier(criterion='entropy', max_depth=6, max_leaf_nodes=11,  
                      min_samples_leaf=3, min_samples_split=3)
```

```
In [31]:
```

```
plt.figure(figsize=(15,8))  
tree.plot_tree(dt,feature_names=['buying','maint','doors','persons','lug_boot','safety'],  
               class_names=list(set(y_train)),filled=True)  
plt.show()
```





Cross Validation

In [32]:

```
from sklearn.model_selection import cross_val_score
```

In [33]:

```
score=cross_val_score(dt,X_train,y_train,cv=10,scoring='accuracy')
score.mean()
```

Out[33]:

```
0.8461088154269971
```

In [34]:

```
print(f'Training set score: {dt.score(X_train,y_train)}')
print(f'Test set score: {dt.score(X_test,y_test)}')
```

```
Training set score: 0.858560794044665
Test set score: 0.861271676300578
```

In [35]:

```
from sklearn.metrics import confusion_matrix,classification_report
cm=confusion_matrix(y_test,y_pred)
```

In [36]:

```
print(cm)
```

```
[[ 78   0  32   8]
 [ 16   0   0   3]
 [ 13   0 345   0]
 [  0   0   0 24]]
```

In [37]:

```
print(classification_report(y_test,y_pred))
```

	precision	recall	f1-score	support
acc	0.73	0.66	0.69	118
good	0.00	0.00	0.00	19
unacc	0.92	0.96	0.94	358
vgood	0.69	1.00	0.81	24
accuracy			0.86	519
macro avg	0.58	0.66	0.61	519
weighted avg	0.83	0.86	0.84	519

C:\Users\USER\anaconda3\lib\site-packages\sklearn\metrics_classification.py:1221: UndefinedMetricWarning: Precision and F-score are ill-defined and being set to 0.0 in labels with no predicted samples. Use `zero_division` parameter to control this behavior.
_warn_prf(average, modifier, msg_start, len(result))