

[illegible]

wine.data

[illegible]

```
wine.data.shape
```

### 3.1 2.1把行列数据合并成表格形式

```
import pandas as pd
```

```
pd.concat([pd.DataFrame(wine.data),pd.DataFrame(wine.target)],axis=1)
```

178 rows × 14 columns

```
wine.feature_names
```

```
wine.target_names
```

### 3.2 2.2分割训练集和测试集

```
Xtrain, Xtest, Ytrain, Ytest = train_test_split(wine.data, wine.target, test_size=0.3)
```

In [12]:	Xtrain.shape
(124, 13)	
In [13]:	Ytrain
array([2, 0, 0, 0, 1, 2, 1, 2, 1, 0, 1, 1, 2, 1, 1, 0, 1, 2, 2, 0, 0, 2, 1, 0, 0, 2, 2, 1, 1, 0, 0, 0, 1, 1, 2, 0, 1, 2, 1, 0, 1, 1, 2, 2, 0, 0, 0, 0, 0, 0, 2, 1, 1, 2, 0, 2, 0, 1, 0, 0, 2, 1, 0, 0, 1, 0, 1, 1, 1, 1, 2, 0, 0, 0, 1, 1, 1, 2, 2, 2, 0, 2, 0, 1, 1, 2, 0, 0, 0, 2, 2, 0, 0, 1, 0, 0, 1, 1, 1, 2, 1, 0, 2, 1, 1, 2, 1, 1, 1, 2, 1, 2, 2, 2, 1, 2, 1, 0, 0, 1, 0, 1])	
4 3.训练模型	
4.1 3.1搭建模型	
In [14]:	<div><div>▼</div><div>clf = tree.DecisionTreeClassifier(criterion="entropy"  # 随机性  ,random_state=30  ,splitter="random"  # 剪枝  ,max_depth = 4 #最大层数为3 不含根  # ,min_samples_leaf=5 #叶子结点包含最少样本数  # ,min_samples_split=20 #中间节点包含最少样本数  )  clf = clf.fit(Xtrain, Ytrain) score = clf.score(Xtest, Ytest) #返回预测的准确度<i>accuracy</i>  score</div></div>
0.9074074074074074	
4.2 3.2画树	
In [15]:	<div><div>▼</div><div>feature_name = ['酒精','苹果酸','灰','灰的碱性','镁','总酚','类黄酮','非黄烷类酚类','花青素','颜色强度','色调','od280/od315稀  import graphviz dot_data = tree.export_graphviz(clf  # ,feature_names = feature_name  # ,class_names=["琴酒","雪莉","贝尔摩德"]  # ,filled=True #填充颜色，颜色深度和不纯度有关  # ,rounded=True #框的形状  )</div></div>
In [16]:	<div><div>graph = graphviz.Source(dot_data) graph</div></div>
<graphviz.files.Source at 0x1ab6e1bbc48>	
4.3 3.3测试各个特征值的权重	
In [17]:	clf.feature_importances_  array([0.14919183, 0.02028786, 0.02624354, 0.01995807, 0. , 0. , 0.49343547, 0. , 0.01523205, 0.01116594, 0. , 0.18185863, 0.08262662])
4.4 3.4给特征值添加名称	
In [18]:	[*zip(feature_name,clf.feature_importances_)]  [('酒精', 0.14919183271291966), ( '苹果酸', 0.020287855558148975), ( '灰', 0.026243535272012935), ( '灰的碱性', 0.01995807000571754), ( '镁', 0.0), ( '总酚', 0.0), ( '类黄酮', 0.49343546921712306), ( '非黄烷类酚类', 0.0), ( '花青素', 0.015232053722644873), ( '颜色强度', 0.011165938842047928), ( '色调', 0.0), ( 'od280/od315稀释葡萄酒', 0.1818586259206247), ( '脯氨酸', 0.0826266187487606)]

In [19]:

score = clf.score(Xtrain,Ytrain)  
score

0.9596774193548387

## 5 4.用matplotlib画图，调参数

In [20]:

import matplotlib.pyplot as plt  
test = []  
for i in range(10):  
 clf = tree.DecisionTreeClassifier(max\_depth=i+1  
 ,criterion="entropy"  
 ,random\_state=30  
 ,splitter="random"  
 )  
  
 clf = clf.fit(Xtrain, Ytrain)  
 score = clf.score(Xtest, Ytest)  
 test.append(score)  
plt.plot(range(1,11),test,color="red",label="max\_depth")  
plt.legend()  
plt.show()

Matplotlib is building the font cache; this may take a moment.



max_depth	score
1	0.58
2	0.78
3	0.91
4	0.91
5	0.91
6	0.95
7	0.95
8	0.95
9	0.95
10	0.95

In [21]:

*#apply*返回每个测试样本所在的叶子节点的索引  
clf.apply(Xtest)

array([ 7, 9, 23, 9, 31, 14, 23, 20, 7, 7, 20, 20, 20, 4, 31, 7, 7,  
 7, 23, 7, 7, 31, 23, 7, 14, 23, 23, 31, 20, 31, 31, 31, 7, 23,  
 23, 31, 9, 7, 7, 7, 31, 14, 20, 20, 31, 7, 28, 14, 31, 28, 11,  
 28, 31, 23], dtype=int64)

In [22]:

*#predict*返回每个测试样本的分类/回归结果  
clf.predict(Xtest)

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

In [ ]: