

```
In [18]: import sklearn
sklearn.__version__
```

```
Out[18]: '0.23.2'
```

```
In [19]: %matplotlib inline
```

Post pruning decision trees with cost complexity pruning

..currentmodule:: sklearn.tree

The :class: DecisionTreeClassifier provides parameters such as min_samples_leaf and max_depth to prevent a tree from overfitting. Cost complexity pruning provides another option to control the size of a tree. In :class: DecisionTreeClassifier , this pruning technique is parameterized by the cost complexity parameter, ccp_alpha . Greater values of ccp_alpha increase the number of nodes pruned. Here we only show the effect of ccp_alpha on regularizing the trees and how to choose a ccp_alpha based on validation scores.

See also minimal_cost_complexity_pruning for details on pruning.

```
In [20]: print(__doc__)
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
from sklearn.datasets import load_breast_cancer
from sklearn.tree import DecisionTreeClassifier

Automatically created module for IPython interactive environment
```

```
In [21]: X, y = load_breast_cancer(return_X_y=True)
X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=0)

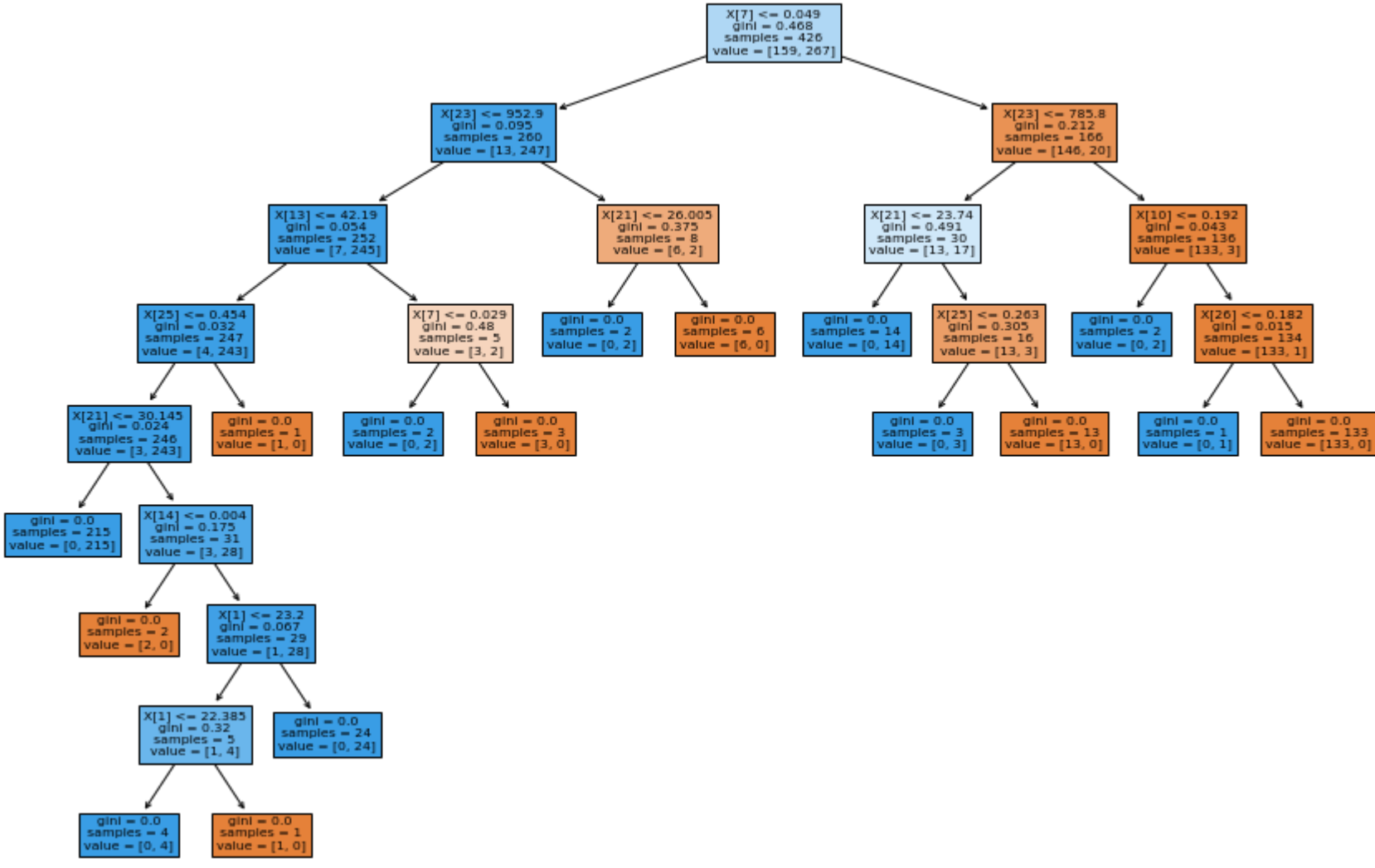
clf = DecisionTreeClassifier(random_state=0)
clf.fit(X_train,y_train)
```

```
Out[21]: DecisionTreeClassifier(random_state=0)
```

```
In [22]: pred=clf.predict(X_test)
from sklearn.metrics import accuracy_score
accuracy_score(y_test, pred)
```

```
Out[22]: 0.8811188811188811
```

```
In [23]: from sklearn import tree
plt.figure(figsize=(15,10))
tree.plot_tree(clf,filled=True)
plt.show()
```



```
In [24]: path = clf.cost_complexity_pruning_path(X_train, y_train)
ccp_alphas, impurities = path.ccp_alphas, path.impurities
```

```
In [25]: ccp_alphas
```

```
Out[25]: array([0.          , 0.00226647, 0.00464743, 0.0046598 , 0.0056338 ,
        0.00704225, 0.00784194, 0.00911402, 0.01144366, 0.018988 ,
        0.02314163, 0.03422475, 0.32729844])
```

```
In [26]: clfs = []
for ccp_alpha in ccp_alphas:
    clf = DecisionTreeClassifier(random_state=0, ccp_alpha=ccp_alpha)
    clf.fit(X_train, y_train)
    clfs.append(clf)
print("Number of nodes in the last tree is: {} with ccp_alpha: {}".format(
    clfs[-1].tree_.node_count, ccp_alphas[-1]))
```

Number of nodes in the last tree is: 1 with ccp_alpha: 0.3272984419327777

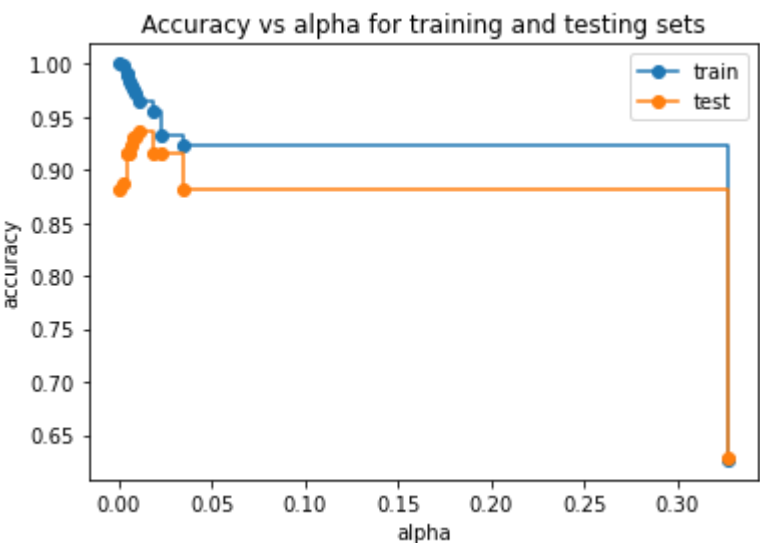
For the remainder of this example, we remove the last element in clfs and ccp_alphas , because it is the trivial tree with only one node. Here we show that the number of nodes and tree depth decreases as alpha increases.

Accuracy vs alpha for training and testing sets

When ccp_alpha is set to zero and keeping the other default parameters of :class: DecisionTreeClassifier , the tree overfits, leading to a 100% training accuracy and 88% testing accuracy. As alpha increases, more of the tree is pruned, thus creating a decision tree that generalizes better. In this example, setting ccp_alpha=0.015 maximizes the testing accuracy.

```
In [27]: train_scores = [clf.score(X_train, y_train) for clf in clfs]
test_scores = [clf.score(X_test, y_test) for clf in clfs]

fig, ax = plt.subplots()
ax.set_xlabel("alpha")
ax.set_ylabel("accuracy")
ax.set_title("Accuracy vs alpha for training and testing sets")
ax.plot(ccp_alphas, train_scores, marker='o', label="train",
        drawstyle="steps-post")
ax.plot(ccp_alphas, test_scores, marker='o', label="test",
        drawstyle="steps-post")
ax.legend()
plt.show()
```



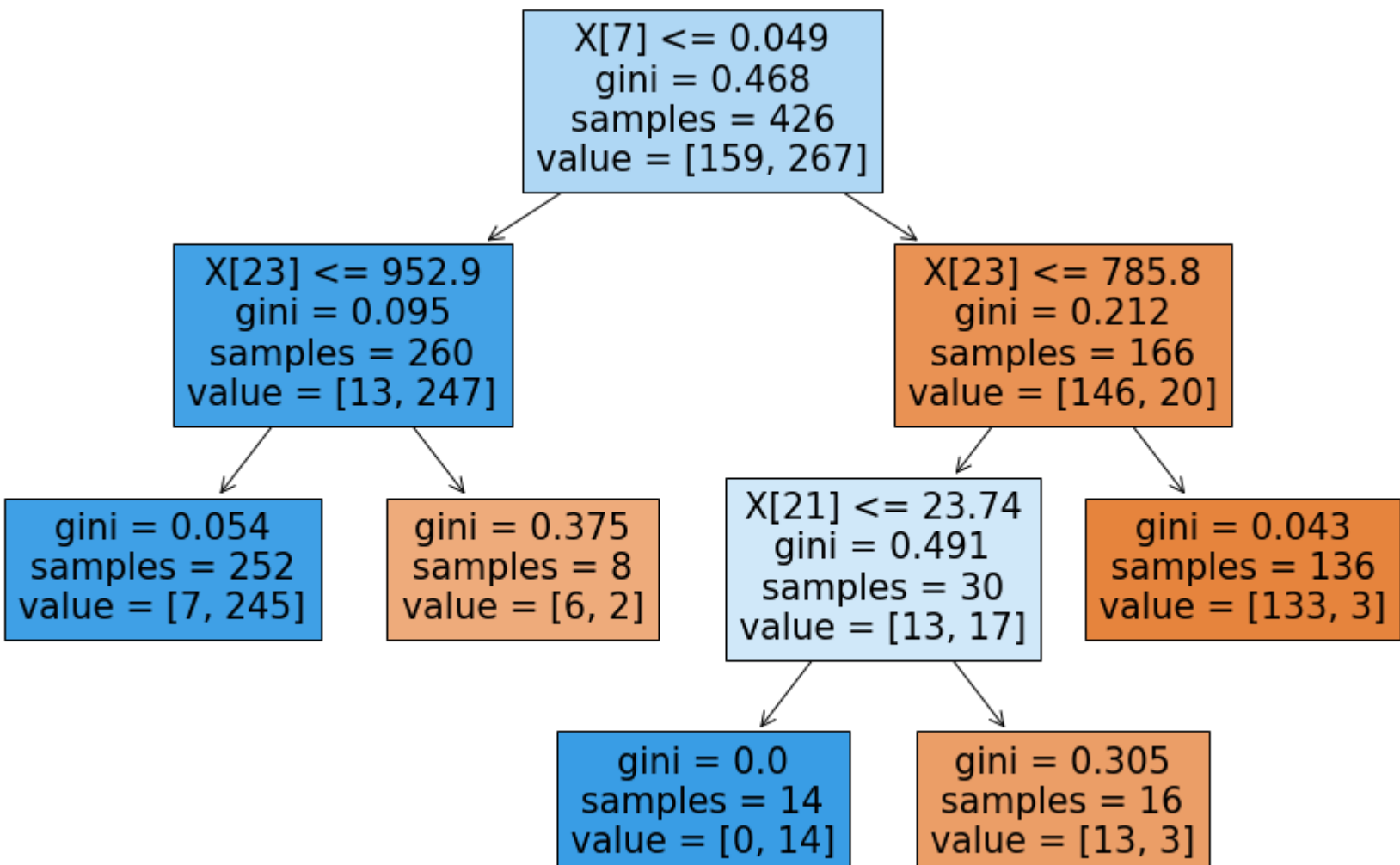
```
In [28]: clf = DecisionTreeClassifier(random_state=0, ccp_alpha=0.012)
clf.fit(X_train,y_train)
```

```
Out[28]: DecisionTreeClassifier(ccp_alpha=0.012, random_state=0)
```

```
In [29]: pred=clf.predict(X_test)
from sklearn.metrics import accuracy_score
accuracy_score(y_test, pred)
```

```
Out[29]: 0.9370629370629371
```

```
In [31]: from sklearn import tree
plt.figure(figsize=(15,10))
tree.plot_tree(clf,filled=True)
plt.show()
```



```
In [ ]:
```

```
In [ ]:
```