# semi-supervised

November 21, 2024

```
[147]: from sklearn import semi_supervised
       from sklearn.ensemble import RandomForestClassifier
       from sklearn.model_selection import train_test_split
       from sklearn.metrics import accuracy_score
       from sklearn.datasets import make_classification
       from sklearn import datasets
       from sklearn.linear_model import LogisticRegression
       from sklearn.tree import export_graphviz
       import pydot
       import numpy as np
       import matplotlib.pyplot as plt
       from sklearn.tree import plot_tree
```

```
[148]: x,y =␣
       ↪make_classification(n_samples=1000,n_features=20,n_classes=2,random_state=42)
```

```
[149]: x_train_val,x_test,y_train_val,y_test = train_test_split(x,y,test_size=0.
       ↪2,random_state=42)
       x_labled,x_unlabeled, y_labled,y_unlabled =␣
       ↪train_test_split(x_train_val,y_train_val,test_size=0.9,random_state=42)
```

```
[150]: print("Total no of samples:",len(x))
       print("Total no of labled samples:",len(x_labled))
       print("Total no of unlabled samples:",len(x_unlabeled))
       print("Total no of test samples:",len(x_test))
```

```
Total no of samples: 1000
Total no of labled samples: 80
Total no of unlabled samples: 720
Total no of test samples: 200
```

```
[151]: model = RandomForestClassifier(n_estimators=100,random_state=42)

       max_iteration=10
       confidence_threshold =0.6
       for iteration in range(max_iteration):
           model.fit(x_labled,y_labled)
```

```python
    probas = model.predict_proba(x_unlabeled)
    confident_indices = np.where(np.max(probas,axis=1)>=confidence_threshold)[0]
    if len(confident_indices) == 0:
        print("No confident Prediction found stop  training")
        break
    x_labled = np.vstack((x_labled,x_unlabeled[confident_indices]))
    y_labled = np.hstack((y_labled,np.argmax(probas[confident_indices],axis=1)))

    x_unlabeled = np.delete(x_unlabeled,confident_indices,axis=0)

    print(f"Iteration:{iteration+1}")
    print(f"Added {len(confident_indices)} confident prediction to the labled␣
 ↪dataset")
    print(f"Remaining unlabled smaples: {len(x_unlabeled)}")
    print(f"Number of training samples after iteration {iteration+1}:
 ↪{len(x_labled)}")
    if len(x_unlabeled)==0:
        break

# final retraining on originally labelled data
model.fit(x_labled,y_labled)

y_pred = model.predict(x_test)
test_acc = accuracy_score(y_test,y_pred)
print(f"Test accuracy after final retraining on test data: {test_acc:.4f}")

y_test_pred = model.predict(x_labled)
test_acc = accuracy_score(y_labled,y_test_pred)
print(f"Test accuracy after final retraining on training data: {test_acc:.4f}")

tree = model.estimators_[0]
plt.figure(figsize=(12, 8))
plot_tree(tree, filled=True)
plt.show()
```

```
Iteration:1
Added 653 confident prediction to the labled dataset
Remaining unlabled smaples: 67
Number of training samples after iteration 1:733
Iteration:2
Added 54 confident prediction to the labled dataset
Remaining unlabled smaples: 13
Number of training samples after iteration 2:787
Iteration:3
Added 4 confident prediction to the labled dataset
Remaining unlabled smaples: 9
Number of training samples after iteration 3:791
```

```
Iteration:4
Added 1 confident prediction to the labled dataset
Remaining unlabled smaples: 8
Number of training samples after iteration 4:792
No confident Prediction found stop  training
Test accuracy after final retraining on test data: 0.8650
Test accuracy after final retraining on training data: 1.0000
```
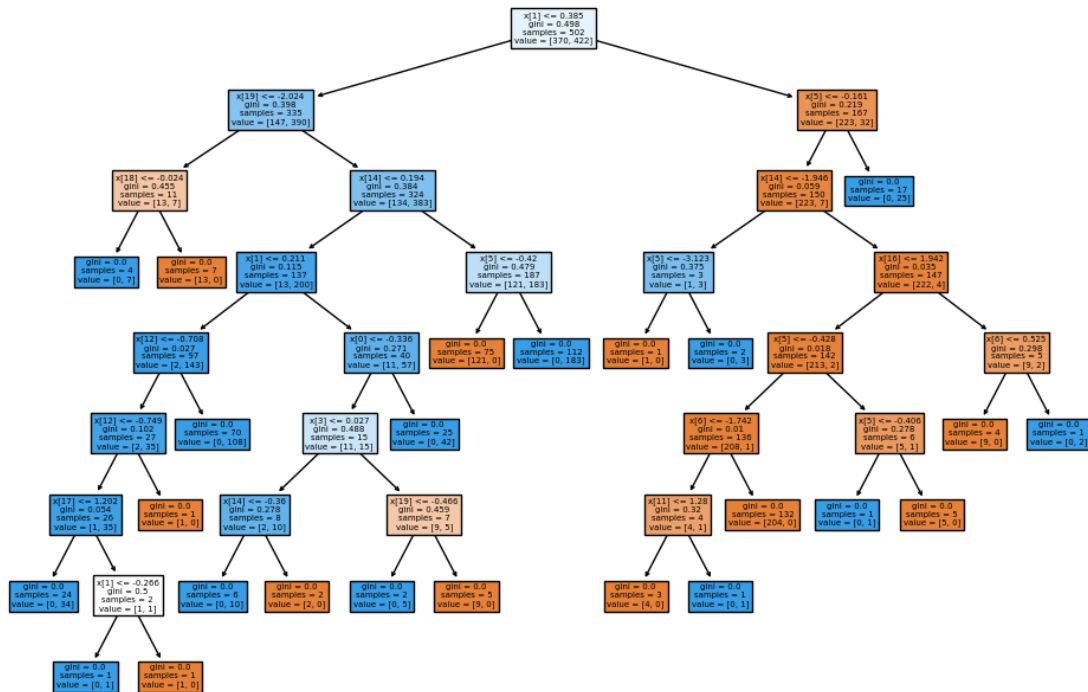


```python
[152]: iris = datasets.load_iris()
       x = iris.data
       y = iris.target
```

```python
[153]: x_train_val,x_test,y_train_val,y_test = train_test_split(x,y,test_size=0.
       ↪2,random_state=42)
       x_labled,x_unlabeled, y_labled,y_unlabeled =␣
       ↪train_test_split(x_train_val,y_train_val,test_size=0.9,random_state=42)
```

```python
[154]: print("Total no of samples:",len(x))
       print("Total no of labled samples:",len(x_labled))
       print("Total no of unlabled samples:",len(x_unlabeled))
       print("Total no of test samples:",len(x_test))
```

```
Total no of samples: 150
Total no of labled samples: 12
```

```
Total no of unlabled samples: 108
Total no of test samples: 30
```

[155]:
```python
model = RandomForestClassifier(n_estimators=100,random_state=42)
max_iteration=10
confidence_threshold =0.7
for iteration in range(max_iteration):
    model.fit(x_labled,y_labled)
    probas = model.predict_proba(x_unlabeled)
    confident_indices = np.where(np.max(probas,axis=1)>=confidence_threshold)[0]
    if len(confident_indices) == 0:
        print("No confident Prediction found stop  training")
        break

    x_labled = np.vstack((x_labled,x_unlabeled[confident_indices]))
    y_labled = np.hstack((y_labled,np.argmax(probas[confident_indices],axis=1)))

    x_unlabeled = np.delete(x_unlabeled,confident_indices,axis=0)

    print(f"Iteration:{iteration+1}")
    print(f"Added {len(confident_indices)} confident prediction to the labled␣
  ↪dataset")
    print(f"Remaining unlabled smaples: {len(x_unlabeled)}")
    print(f"Number of training samples after iteration {iteration+1}:␣
  ↪{len(x_labled)}")
    if len(x_unlabeled)==0:
        break

# final retraining on originally labelled data
model.fit(x_labled,y_labled)

y_pred = model.predict(x_test)
test_acc = accuracy_score(y_test,y_pred)
print(f"Test accuracy after final retraining on test data: {test_acc:.4f}")

y_test_pred = model.predict(x_labled)
test_acc = accuracy_score(y_labled,y_test_pred)
print(f"Test accuracy after final retraining on training data: {test_acc:.4f}")


tree = model.estimators_[0]
plt.figure(figsize=(12, 8))
plot_tree(tree, feature_names=iris.feature_names, class_names=iris.target_names.
  ↪tolist(), filled=True)
plt.show()
```
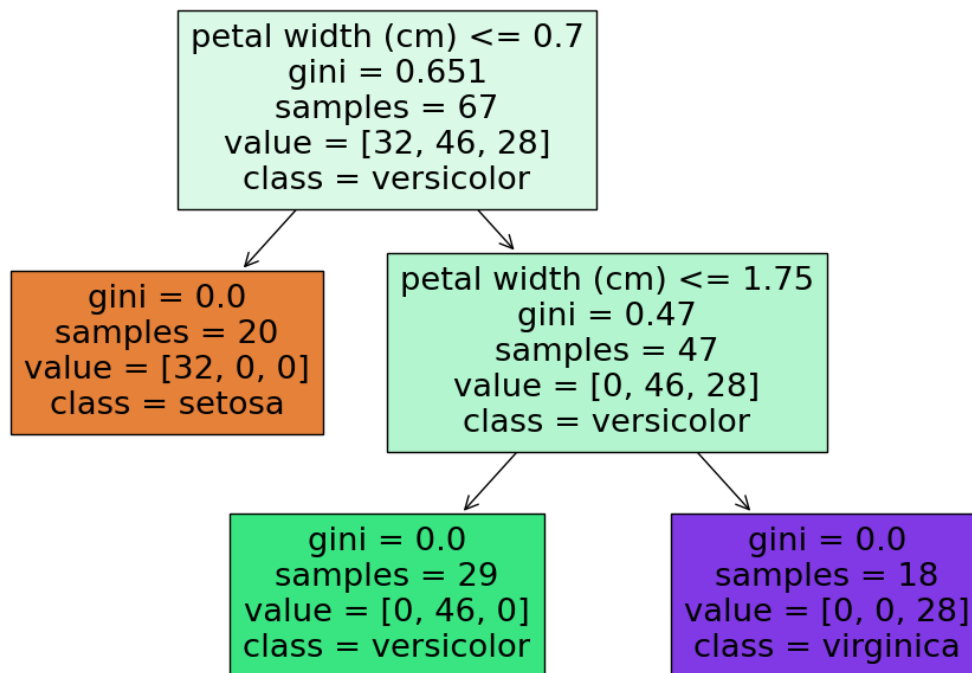
```
Iteration:1
```

```
Added 92 confident prediction to the labled dataset
Remaining unlabled smaples: 16
Number of training samples after iteration 1:104
Iteration:2
Added 2 confident prediction to the labled dataset
Remaining unlabled smaples: 14
Number of training samples after iteration 2:106
No confident Prediction found stop  training
Test accuracy after final retraining on test data: 0.9000
Test accuracy after final retraining on training data: 1.0000
```

```
                    petal width (cm) <= 0.7
                         gini = 0.651
                         samples = 67
                      value = [32, 46, 28]
                       class = versicolor


        gini = 0.0                petal width (cm) <= 1.75
      samples = 20                     gini = 0.47
    value = [32, 0, 0]                samples = 47
      class = setosa              value = [0, 46, 28]
                                   class = versicolor


                          gini = 0.0              gini = 0.0
                        samples = 29            samples = 18
                      value = [0, 46, 0]      value = [0, 0, 28]
                       class = versicolor      class = virginica
```

[ ]: 

[ ]: