

Chapter 7 – Ensemble Learning and Random Forests

This notebook contains all the sample code and solutions to the exercises in chapter 7.

[Open in Colab](#)

[Open in Kaggle](#)

Setup

This project requires Python 3.7 or above:

```
In [1]: import sys  
  
assert sys.version_info >= (3, 7)
```

It also requires Scikit-Learn ≥ 1.0.1:

```
In [2]: from packaging import version  
import sklearn  
  
assert version.parse(sklearn.__version__) >= version.parse("1.0.1")
```

As we did in previous chapters, let's define the default font sizes to make the figures prettier:

```
In [3]: import matplotlib.pyplot as plt  
  
plt.rc('font', size=14)  
plt.rc('axes', labelsize=14, titlesize=14)  
plt.rc('legend', fontsize=14)  
plt.rc('xtick', labelsize=10)  
plt.rc('ytick', labelsize=10)
```

And let's create the `images/ensembles` folder (if it doesn't already exist), and define the `save_fig()` function which is used through this notebook to save the figures in high-res for the book:

```
In [4]: from pathlib import Path  
  
IMAGES_PATH = Path() / "images" / "ensembles"  
IMAGES_PATH.mkdir(parents=True, exist_ok=True)  
  
def save_fig(fig_id, tight_layout=True, fig_extension="png", resolution=300):  
    path = IMAGES_PATH / f"{fig_id}.{fig_extension}"  
    if tight_layout:  
        plt.tight_layout()  
    plt.savefig(path, format=fig_extension, dpi=resolution)
```

Voting Classifiers

Let's build a voting classifier:

```
In [5]: from sklearn.datasets import make_moons  
from sklearn.ensemble import RandomForestClassifier, VotingClassifier  
from sklearn.linear_model import LogisticRegression  
from sklearn.model_selection import train_test_split  
from sklearn.svm import SVC  
  
X, y = make_moons(n_samples=500, noise=0.30, random_state=42)  
X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=42)  
  
voting_clf = VotingClassifier(  
    estimators=[  
        ('lr', LogisticRegression(random_state=42)),  
        ('rf', RandomForestClassifier(random_state=42)),  
        ('svc', SVC(random_state=42))  
    ])
```

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

Out[5]:



```
In [6]: for name, clf in voting_clf.named_estimators_.items():  
    print(name, " = ", clf.score(X_test, y_test))
```

lr = 0.864
rf = 0.896
svc = 0.896

```
In [7]: voting_clf.predict(X_test[:1])
```

Out[7]: array([1])

```
In [8]: [clf.predict(X_test[:1]) for clf in voting_clf.estimators_]
```

Out[8]: [array([1]), array([1]), array([0])]

```
In [9]: voting_clf.score(X_test, y_test)
```

Out[9]: 0.912

Now let's use soft voting:

```
In [10]: voting_clf.voting = "soft"  
voting_clf.named_estimators["svc"].probability = True  
voting_clf.fit(X_train, y_train)  
voting_clf.score(X_test, y_test)
```

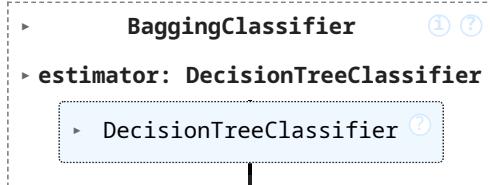
Out[10]: 0.92

Bagging and Pasting

Bagging and Pasting in Scikit-Learn

```
In [11]: from sklearn.ensemble import BaggingClassifier  
from sklearn.tree import DecisionTreeClassifier  
  
bag_clf = BaggingClassifier(DecisionTreeClassifier(), n_estimators=500,  
                           max_samples=100, n_jobs=-1, random_state=42)  
bag_clf.fit(X_train, y_train)
```

Out[11]:



```
In [12]: import numpy as np
```

```
In [13]: # extra code – this cell generates and saves Figure 7–5  
  
def plot_decision_boundary(clf, X, y, alpha=1.0):  
    axes=[-1.5, 2.4, -1, 1.5]  
    x1, x2 = np.meshgrid(np.linspace(axes[0], axes[1], 100),  
                         np.linspace(axes[2], axes[3], 100))  
    X_new = np.c_[x1.ravel(), x2.ravel()]  
    y_pred = clf.predict(X_new).reshape(x1.shape)  
  
    plt.contourf(x1, x2, y_pred, alpha=0.3 * alpha, cmap='Wistia')  
    plt.contour(x1, x2, y_pred, cmap="Greys", alpha=0.8 * alpha)
```

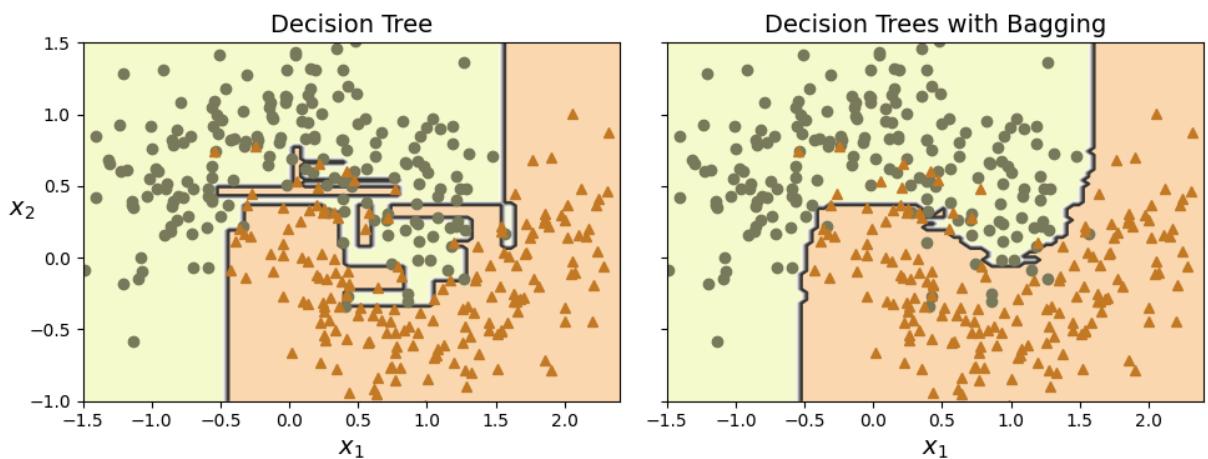
```

colors = ["#78785c", "#c47b27"]
markers = ("o", "^")
for idx in (0, 1):
    plt.plot(X[:, 0][y == idx], X[:, 1][y == idx],
              color=colors[idx], marker=markers[idx], linestyle="none")
plt.axis(axes)
plt.xlabel(r"$x_1$")
plt.ylabel(r"$x_2$", rotation=0)

tree_clf = DecisionTreeClassifier(random_state=42)
tree_clf.fit(X_train, y_train)

fig, axes = plt.subplots(ncols=2, figsize=(10, 4), sharey=True)
plt.sca(axes[0])
plot_decision_boundary(tree_clf, X_train, y_train)
plt.title("Decision Tree")
plt.sca(axes[1])
plot_decision_boundary(bag_clf, X_train, y_train)
plt.title("Decision Trees with Bagging")
plt.ylabel("")
save_fig("decision_tree_without_and_with_bagging_plot")
plt.show()

```



Out-of-Bag evaluation

```
In [14]: bag_clf = BaggingClassifier(DecisionTreeClassifier(), n_estimators=500,
                                  oob_score=True, n_jobs=-1, random_state=42)
bag_clf.fit(X_train, y_train)
bag_clf.oob_score_
```

```
Out[14]: 0.896
```

```
In [15]: bag_clf.oob_decision_function_[:3] # probas for the first 3 instances
```

```
Out[15]: array([[0.32352941, 0.67647059],
                 [0.3375      , 0.6625      ],
                 [1.         , 0.         ]])
```

```
In [16]: from sklearn.metrics import accuracy_score
y_pred = bag_clf.predict(X_test)
accuracy_score(y_test, y_pred)
```

```
Out[16]: 0.92
```

If you randomly draw one instance from a dataset of size m , each instance in the dataset obviously has probability $1/m$ of getting picked, and therefore it has a probability $1 - 1/m$ of *not* getting picked. If you draw m instances with replacement, all draws are independent and therefore each instance has a probability $(1 - 1/m)^m$ of *not* getting picked. Now let's use the fact that $\exp(x)$ is equal to the limit of $(1 + x/m)^m$ as m approaches infinity. So if m is large, the ratio of out-of-bag instances will be about $\exp(-1) \approx 0.37$. So roughly 63% ($1 - 0.37$) will be sampled.

```
In [17]: # extra code - shows how to compute the 63% proba
print(1 - (1 - 1 / 1000) ** 1000)
print(1 - np.exp(-1))

0.6323045752290363
0.6321205588285577
```

Random Forests

```
In [18]: from sklearn.ensemble import RandomForestClassifier

rnd_clf = RandomForestClassifier(n_estimators=500, max_leaf_nodes=16,
                                 n_jobs=-1, random_state=42)
rnd_clf.fit(X_train, y_train)
y_pred_rf = rnd_clf.predict(X_test)
```

A Random Forest is equivalent to a bag of decision trees:

```
In [19]: bag_clf = BaggingClassifier(
    DecisionTreeClassifier(max_features="sqrt", max_leaf_nodes=16),
    n_estimators=500, n_jobs=-1, random_state=42)
```

```
In [20]: # extra code - verifies that the predictions are identical
bag_clf.fit(X_train, y_train)
y_pred_bag = bag_clf.predict(X_test)
np.all(y_pred_bag == y_pred_rf) # same predictions
```

```
Out[20]: np.True_
```

Feature Importance

```
In [21]: from sklearn.datasets import load_iris

iris = load_iris(as_frame=True)
rnd_clf = RandomForestClassifier(n_estimators=500, random_state=42)
rnd_clf.fit(iris.data, iris.target)
for score, name in zip(rnd_clf.feature_importances_, iris.data.columns):
    print(round(score, 2), name)

0.11 sepal length (cm)
0.02 sepal width (cm)
0.44 petal length (cm)
0.42 petal width (cm)
```

```
In [22]: # extra code - this cell generates and saves Figure 7-6

from sklearn.datasets import fetch_openml

X_mnist, y_mnist = fetch_openml('mnist_784', return_X_y=True, as_frame=False)

rnd_clf = RandomForestClassifier(n_estimators=100, random_state=42)
rnd_clf.fit(X_mnist, y_mnist)

heatmap_image = rnd_clf.feature_importances_.reshape(28, 28)
plt.imshow(heatmap_image, cmap="hot")
cbar = plt.colorbar(ticks=[rnd_clf.feature_importances_.min(),
                           rnd_clf.feature_importances_.max()])
cbar.ax.set_yticklabels(['Not important', 'Very important'], fontsize=14)
plt.axis("off")
save_fig("mnist_feature_importance_plot")
plt.show()
```

