

## Chapter 7 – Ensemble Learning and Random Forests

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



[Run in Google Colab \(https://colab.research.google.com/github/ageron/handson-ml2/blob/master/07\\_ensemble\\_learning\\_and\\_random\\_forests.ipynb\)](https://colab.research.google.com/github/ageron/handson-ml2/blob/master/07_ensemble_learning_and_random_forests.ipynb)

## Setup

First, let's import a few common modules, ensure Matplotlib plots figures inline and prepare a function to save the figures. We also check that Python 3.5 or later is installed (although Python 2.x may work, it is deprecated so we strongly recommend you use Python 3 instead), as well as Scikit-Learn  $\geq 0.20$ .

In [1]:

```
# Python ≥3.5 is required
import sys
assert sys.version_info >= (3, 5)

# Scikit-Learn ≥0.20 is required
import sklearn
assert sklearn.__version__ >= "0.20"

# Common imports
import numpy as np
import os

# to make this notebook's output stable across runs
np.random.seed(42)

# To plot pretty figures
%matplotlib inline
import matplotlib as mpl
import matplotlib.pyplot as plt
mpl.rc('axes', labelsz=14)
mpl.rc('xtick', labelsz=12)
mpl.rc('ytick', labelsz=12)

# Where to save the figures
PROJECT_ROOT_DIR = "."
CHAPTER_ID = "ensembles"
IMAGES_PATH = os.path.join(PROJECT_ROOT_DIR, "images", CHAPTER_ID)
os.makedirs(IMAGES_PATH, exist_ok=True)

def save_fig(fig_id, tight_layout=True, fig_extension="png", resolution=
    path = os.path.join(IMAGES_PATH, fig_id + "." + fig_extension)
    print("Saving figure", fig_id)
    if tight_layout:
        plt.tight_layout()
    plt.savefig(path, format=fig_extension, dpi=resolution)
```

## Voting classifiers



In [5]:

```
from sklearn.metrics import accuracy_score

for clf in (log_clf, rnd_clf, svm_clf, voting_clf):
    clf.fit(X_train, y_train)
    y_pred = clf.predict(X_test)
    print(clf.__class__.__name__, accuracy_score(y_test, y_pred))
```

```
LogisticRegression 0.864
RandomForestClassifier 0.896
SVC 0.896
VotingClassifier 0.912
```

Soft voting:

In [6]:

```
log_clf = LogisticRegression(solver="lbfgs", random_state=42)
rnd_clf = RandomForestClassifier(n_estimators=100, random_state=42)
svm_clf = SVC(gamma="scale", probability=True, random_state=42)

voting_clf = VotingClassifier(
    estimators=[('lr', log_clf), ('rf', rnd_clf), ('svc', svm_clf)],
    voting='soft')
voting_clf.fit(X_train, y_train)
```

Out[6]:

```
VotingClassifier(estimators=[('lr', LogisticRegression(ran
dom_state=42)),
                           ('rf', RandomForestClassifier
(random_state=42)),
                           ('svc', SVC(probability=True,
random_state=42))],
                 voting='soft')
```

In [7]:

```
from sklearn.metrics import accuracy_score

for clf in (log_clf, rnd_clf, svm_clf, voting_clf):
    clf.fit(X_train, y_train)
    y_pred = clf.predict(X_test)
    print(clf.__class__.__name__, accuracy_score(y_test, y_pred))
```

LogisticRegression 0.864  
RandomForestClassifier 0.896  
SVC 0.896  
VotingClassifier 0.92

## Bagging ensembles

In [8]:

```
from sklearn.ensemble import BaggingClassifier
from sklearn.tree import DecisionTreeClassifier

bag_clf = BaggingClassifier(
    DecisionTreeClassifier(random_state=42), n_estimators=500,
    max_samples=100, bootstrap=True, random_state=42)
bag_clf.fit(X_train, y_train)
y_pred = bag_clf.predict(X_test)
```

In [9]:

```
from sklearn.metrics import accuracy_score
print(accuracy_score(y_test, y_pred))
```

0.904

In [10]:

```
tree_clf = DecisionTreeClassifier(random_state=42)
tree_clf.fit(X_train, y_train)
y_pred_tree = tree_clf.predict(X_test)
print(accuracy_score(y_test, y_pred_tree))
```

0.856

In [11]:

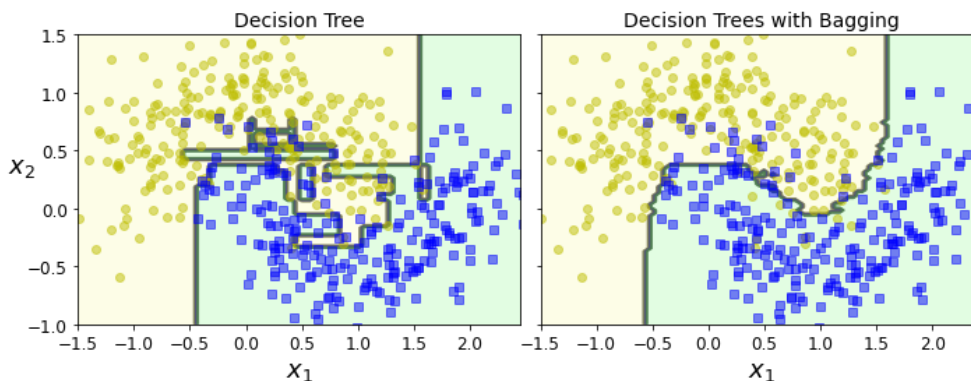
```
from matplotlib.colors import ListedColormap

def plot_decision_boundary(clf, X, y, axes=[-1.5, 2.45, -1, 1.5], alpha=
    x1s = np.linspace(axes[0], axes[1], 100)
    x2s = np.linspace(axes[2], axes[3], 100)
    x1, x2 = np.meshgrid(x1s, x2s)
    X_new = np.c_[x1.ravel(), x2.ravel()]
    y_pred = clf.predict(X_new).reshape(x1.shape)
    custom_cmap = ListedColormap(['#fafab0', '#9898ff', '#a0faa0'])
    plt.contourf(x1, x2, y_pred, alpha=0.3, cmap=custom_cmap)
    if contour:
        custom_cmap2 = ListedColormap(['#7d7d58', '#4c4c7f', '#507d50'])
        plt.contour(x1, x2, y_pred, cmap=custom_cmap2, alpha=0.8)
    plt.plot(X[:, 0][y==0], X[:, 1][y==0], "yo", alpha=alpha)
    plt.plot(X[:, 0][y==1], X[:, 1][y==1], "bs", alpha=alpha)
    plt.axis(axes)
    plt.xlabel(r"$x_1$", fontsize=18)
    plt.ylabel(r"$x_2$", fontsize=18, rotation=0)
```

In [12]:

```
fig, axes = plt.subplots(ncols=2, figsize=(10,4), sharey=True)
plt.sca(axes[0])
plot_decision_boundary(tree_clf, X, y)
plt.title("Decision Tree", fontsize=14)
plt.sca(axes[1])
plot_decision_boundary(bag_clf, X, y)
plt.title("Decision Trees with Bagging", fontsize=14)
plt.ylabel("")
save_fig("decision_tree_without_and_with_bagging_plot")
plt.show()
```

Saving figure decision\_tree\_without\_and\_with\_bagging\_plot



# Random Forests

In [13]:

```
bag_clf = BaggingClassifier(  
    DecisionTreeClassifier(splitter="random", max_leaf_nodes=16, random_  
        n_estimators=500, max_samples=100, bootstrap=True, random_state=42)
```

In [14]:

```
bag_clf.fit(X_train, y_train)  
y_pred = bag_clf.predict(X_test)
```

In [15]:

```
from sklearn.ensemble import RandomForestClassifier  
  
rnd_clf = RandomForestClassifier(n_estimators=500, max_leaf_nodes=16, ra  
rnd_clf.fit(X_train, y_train)  
  
y_pred_rf = rnd_clf.predict(X_test)
```

In [16]:

```
np.sum(y_pred == y_pred_rf) / len(y_pred)  # almost identical prediction
```

Out[16]:

0.976

In [17]:

```
from sklearn.datasets import load_iris  
iris = load_iris()  
rnd_clf = RandomForestClassifier(n_estimators=500, random_state=42)  
rnd_clf.fit(iris["data"], iris["target"])  
for name, score in zip(iris["feature_names"], rnd_clf.feature_importance  
    print(name, score)
```

```
sepal length (cm) 0.11249225099876375  
sepal width (cm) 0.02311928828251033  
petal length (cm) 0.4410304643639577  
petal width (cm) 0.4233579963547682
```

In [18]:

```
rnd_clf.feature_importances_
```

Out[18]:

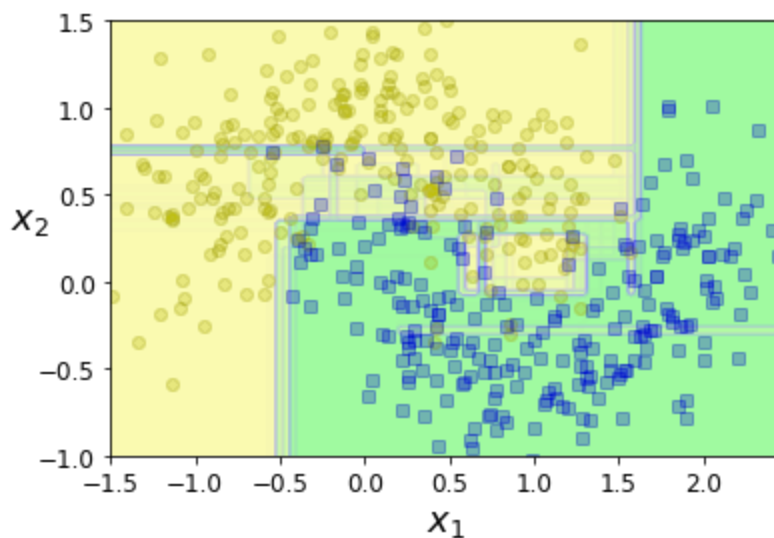
```
array([0.11249225, 0.02311929, 0.44103046, 0.423358  ])
```

In [19]:

```
plt.figure(figsize=(6, 4))

for i in range(15):
    tree_clf = DecisionTreeClassifier(max_leaf_nodes=16, random_state=42)
    indices_with_replacement = np.random.randint(0, len(X_train), len(X_train))
    tree_clf.fit(X_train[indices_with_replacement], y_train[indices_with_replacement])
    plot_decision_boundary(tree_clf, X_train, y_train, axes=[-1.5, 2.45, -1, 1.5], a

plt.show()
```



## Out-of-Bag evaluation



```
bag_clf = BaggingClassifier(
    DecisionTreeClassifier(random_state=42), n_estimators=500,
    bootstrap=True, oob_score=True, random_state=40)
bag_clf.fit(X_train, y_train)
bag_clf.oob_score_
```

0.8986666666666666

```
bag_clf.oob_decision_function_
```

```
[0.          , 1.          ],
[0.62569832, 0.37430168],
[0.          , 1.          ],
[1.          , 0.          ],
[0.          , 1.          ],
[0.          , 1.          ],
[0.13402062, 0.86597938],
[1.          , 0.          ],
[0.          , 1.          ],
[0.38251366, 0.61748634],
[0.          , 1.          ],
[1.          , 0.          ],
[0.27093596, 0.72906404],
[0.34146341, 0.65853659],
[1.          , 0.          ],
[1.          , 0.          ],
[0.          , 1.          ],
[1.          , 0.          ],
[1.          , 0.          ],
[0.          , 1.          ],
[0.          , 1.          ]
```

```
from sklearn.metrics import accuracy_score
y_pred = bag_clf.predict(X_test)
accuracy_score(y_test, y_pred)
```

0.912

## Feature importance

In [23]:

```
from sklearn.datasets import fetch_openml  
  
mnist = fetch_openml('mnist_784', version=1)  
mnist.target = mnist.target.astype(np.uint8)
```

In [24]:

```
rnd_clf = RandomForestClassifier(n_estimators=100, random_state=42)  
rnd_clf.fit(mnist["data"], mnist["target"])
```

Out[24]:

```
RandomForestClassifier(random_state=42)
```

In [25]:

```
def plot_digit(data):  
    image = data.reshape(28, 28)  
    plt.imshow(image, cmap = mpl.cm.hot,  
                interpolation="nearest")  
    plt.axis("off")
```

In [26]:

```
plot_digit(rnd_clf.feature_importances_)

cbar = plt.colorbar(ticks=[rnd_clf.feature_importances_.min(), rnd_clf.f
cbar.ax.set_yticklabels(['Not important', 'Very important'])

save_fig("mnist_feature_importance_plot")
plt.show()
```

Saving figure mnist\_feature\_importance\_plot

