

Lab Module 2 – Classification**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 ≥ 0.20 .

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

# Scikit-Learn  $\geq 0.20$  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', labelsiz=14)
mpl.rc('xtick', labelsiz=12)
mpl.rc('ytick', labelsiz=12)

# Where to save the figures
PROJECT_ROOT_DIR = "."
CHAPTER_ID = "classification"
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=300):
    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)
```

MNIST

```
from sklearn.datasets import fetch_openml
mnist = fetch_openml('mnist_784', version=1)
mnist.keys()

dict_keys(['data', 'target', 'frame', 'categories', 'feature_names', 'target_names', 'DESCR', 'details', 'url'])

X, y = mnist["data"], mnist["target"]
X.shape

(70000, 784)

y.shape

(70000,)

28 * 28

784

%matplotlib inline
import matplotlib as mpl
import matplotlib.pyplot as plt

some_digit = X[0]
some_digit_image = some_digit.reshape(28, 28)
```

```
plt.imshow(some_digit_image, cmap=mpl.cm.binary)
plt.axis("off")
```

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

↗ Saving figure some_digit_plot



```
y[0]
```

↗ '5'

```
y = y.astype(np.uint8)
```

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

```
# EXTRA
```

```
def plot_digits(instances, images_per_row=10, **options):
    size = 28
    images_per_row = min(len(instances), images_per_row)
    images = [instance.reshape(size,size) for instance in instances]
    n_rows = (len(instances) - 1) // images_per_row + 1
    row_images = []
    n_empty = n_rows * images_per_row - len(instances)
    images.append(np.zeros((size, size * n_empty)))
    for row in range(n_rows):
        rimages = images[row * images_per_row : (row + 1) * images_per_row]
        row_images.append(np.concatenate(rimages, axis=1))
    image = np.concatenate(row_images, axis=0)
    plt.imshow(image, cmap = mpl.cm.binary, **options)
    plt.axis("off")
```

```
plt.figure(figsize=(9,9))
example_images = X[:100]
plot_digits(example_images, images_per_row=10)
save_fig("more_digits_plot")
plt.show()
```

↗ Saving figure more_digits_plot



y[0]

↗ 5

```
X_train, X_test, y_train, y_test = X[:60000], X[60000:], y[:60000], y[60000:]
```

▼ Binary classifier

```
y_train_5 = (y_train == 5) # True for all 5s, False for all other digits
y_test_5 = (y_test == 5)
```

Note: some hyperparameters will have a different default value in future versions of Scikit-Learn, such as `max_iter` and `tol`. To be future-proof, we explicitly set these hyperparameters to their future default values. For simplicity, this is not shown in the book.

```
from sklearn.linear_model import SGDClassifier
```

```
sgd_clf = SGDClassifier(max_iter=1000, tol=1e-3, random_state=42)
sgd_clf.fit(X_train, y_train_5)
```

↗ SGDClassifier(random_state=42)

```
sgd_clf.predict([some_digit])
```

↗ array([True])

```
from sklearn.model_selection import cross_val_score
cross_val_score(sgd_clf, X_train, y_train_5, cv=3, scoring="accuracy")
```

↗ array([0.95035, 0.96035, 0.9604])

```
from sklearn.model_selection import StratifiedKFold
from sklearn.base import clone
```

```
skfolds = StratifiedKFold(n_splits=3, shuffle=True, random_state=42)
```

```
for train_index, test_index in skfolds.split(X_train, y_train_5):
    clone_clf = clone(sgd_clf)
    X_train_folds = X_train[train_index]
    y_train_folds = y_train_5[train_index]
```

```

X_test_fold = X_train[test_index]
y_test_fold = y_train_5[test_index]

clone_clf.fit(X_train_folds, y_train_folds)
y_pred = clone_clf.predict(X_test_fold)
n_correct = sum(y_pred == y_test_fold)
print(n_correct / len(y_pred))

```

```

0.9669
0.91625
0.96785

```

Note: shuffle=True was omitted by mistake in previous releases of the book.

```

from sklearn.base import BaseEstimator
class Never5Classifier(BaseEstimator):
    def fit(self, X, y=None):
        pass
    def predict(self, X):
        return np.zeros((len(X), 1), dtype=bool)

```

```

never_5_clf = Never5Classifier()
cross_val_score(never_5_clf, X_train, y_train_5, cv=3, scoring="accuracy")

```

```

array([0.91125, 0.90855, 0.90915])

```

Warning: this output (and many others in this notebook and other notebooks) may differ slightly from those in the book. Don't worry, that's okay! There are several reasons for this:

- first, Scikit-Learn and other libraries evolve, and algorithms get tweaked a bit, which may change the exact result you get. If you use the latest Scikit-Learn version (and in general, you really should), you probably won't be using the exact same version I used when I wrote the book or this notebook, hence the difference. I try to keep this notebook reasonably up to date, but I can't change the numbers on the pages in your copy of the book.
- second, many training algorithms are stochastic, meaning they rely on randomness. In principle, it's possible to get consistent outputs from a random number generator by setting the seed from which it generates the pseudo-random numbers (which is why you will see `random_state=42` or `np.random.seed(42)` pretty often). However, sometimes this does not suffice due to the other factors listed here.
- third, if the training algorithm runs across multiple threads (as do some algorithms implemented in C) or across multiple processes (e.g., when using the `n_jobs` argument), then the precise order in which operations will run is not always guaranteed, and thus the exact result may vary slightly.
- lastly, other things may prevent perfect reproducibility, such as Python maps and sets whose order is not guaranteed to be stable across sessions, or the order of files in a directory which is also not guaranteed.

```

from sklearn.model_selection import cross_val_predict

y_train_pred = cross_val_predict(sgd_clf, X_train, y_train_5, cv=3)

```

```

from sklearn.metrics import confusion_matrix

confusion_matrix(y_train_5, y_train_pred)

```

```

array([[53892,  687],
       [1891, 3530]])

```

```

y_train_perfect_predictions = y_train_5 # pretend we reached perfection
confusion_matrix(y_train_5, y_train_perfect_predictions)

```

```

array([[54579,    0],
       [    0, 5421]])

```

```

from sklearn.metrics import precision_score, recall_score

precision_score(y_train_5, y_train_pred)

```

```

0.8370879772350012

```

```

cm = confusion_matrix(y_train_5, y_train_pred)
cm[1, 1] / (cm[0, 1] + cm[1, 1])

```

```

0.8370879772350012

```

```
recall_score(y_train_5, y_train_pred)
```

```
↪ 0.6511713705958311
```

```
cm[1, 1] / (cm[1, 0] + cm[1, 1])
```

```
↪ 0.6511713705958311
```

```
from sklearn.metrics import f1_score
```

```
f1_score(y_train_5, y_train_pred)
```

```
↪ 0.7325171197343846
```

```
cm[1, 1] / (cm[1, 1] + (cm[1, 0] + cm[0, 1]) / 2)
```

```
↪ 0.7325171197343847
```

```
y_scores = sgd_clf.decision_function([some_digit])
y_scores
```

```
↪ array([2164.22030239])
```

```
threshold = 0
y_some_digit_pred = (y_scores > threshold)
```

```
y_some_digit_pred
```

```
↪ array([ True])
```

```
threshold = 8000
y_some_digit_pred = (y_scores > threshold)
y_some_digit_pred
```

```
↪ array([False])
```

```
y_scores = cross_val_predict(sgd_clf, X_train, y_train_5, cv=3,
                             method="decision_function")
```

```
from sklearn.metrics import precision_recall_curve
```

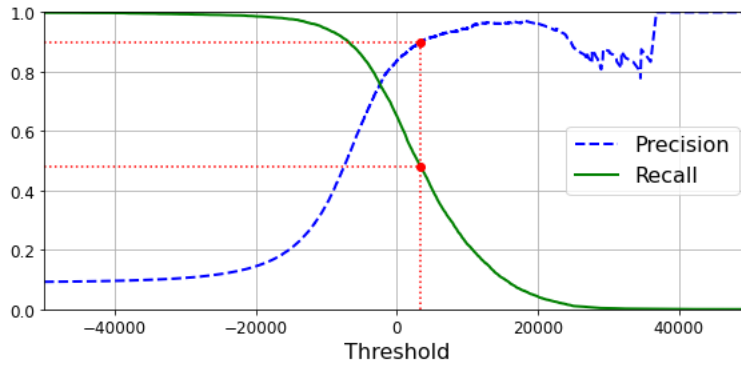
```
precisions, recalls, thresholds = precision_recall_curve(y_train_5, y_scores)
```

```
def plot_precision_recall_vs_threshold(precisions, recalls, thresholds):
    plt.plot(thresholds, precisions[:-1], "b--", label="Precision", linewidth=2)
    plt.plot(thresholds, recalls[:-1], "g-", label="Recall", linewidth=2)
    plt.legend(loc="center right", fontsize=16) # Not shown in the book
    plt.xlabel("Threshold", fontsize=16)       # Not shown
    plt.grid(True)                             # Not shown
    plt.axis([-50000, 50000, 0, 1])           # Not shown
```

```
recall_90_precision = recalls[np.argmax(precisions >= 0.90)]
threshold_90_precision = thresholds[np.argmax(precisions >= 0.90)]
```

```
plt.figure(figsize=(8, 4)) # Not shown
plot_precision_recall_vs_threshold(precisions, recalls, thresholds)
plt.plot([threshold_90_precision, threshold_90_precision], [0., 0.9], "r:") # Not shown
plt.plot([-50000, threshold_90_precision], [0.9, 0.9], "r:") # Not shown
plt.plot([-50000, threshold_90_precision], [recall_90_precision, recall_90_precision], "r:") # Not shown
plt.plot([threshold_90_precision], [0.9], "ro") # Not shown
plt.plot([threshold_90_precision], [recall_90_precision], "ro") # Not shown
save_fig("precision_recall_vs_threshold_plot") # Not shown
plt.show()
```

↗ Saving figure precision_recall_vs_threshold_plot



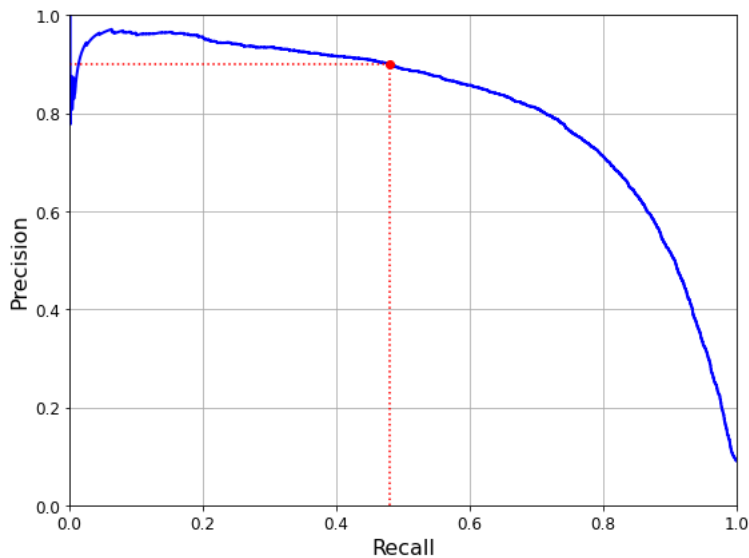
```
(y_train_pred == (y_scores > 0)).all()
```

↗ True

```
def plot_precision_vs_recall(precisions, recalls):
    plt.plot(recalls, precisions, "b-", linewidth=2)
    plt.xlabel("Recall", fontsize=16)
    plt.ylabel("Precision", fontsize=16)
    plt.axis([0, 1, 0, 1])
    plt.grid(True)

plt.figure(figsize=(8, 6))
plot_precision_vs_recall(precisions, recalls)
plt.plot([recall_90_precision, recall_90_precision], [0., 0.9], "r:")
plt.plot([0.0, recall_90_precision], [0.9, 0.9], "r:")
plt.plot([recall_90_precision], [0.9], "ro")
save_fig("precision_vs_recall_plot")
plt.show()
```

↗ Saving figure precision_vs_recall_plot



```
threshold_90_precision = thresholds[np.argmax(precisions >= 0.90)]
```

```
threshold_90_precision
```

↗ 3370.0194991439557

```
y_train_pred_90 = (y_scores >= threshold_90_precision)
```

```
precision_score(y_train_5, y_train_pred_90)
```

↗ 0.9000345901072293

```
recall_score(y_train_5, y_train_pred_90)
```

↗ 0.4799852425751706

✓ ROC curves

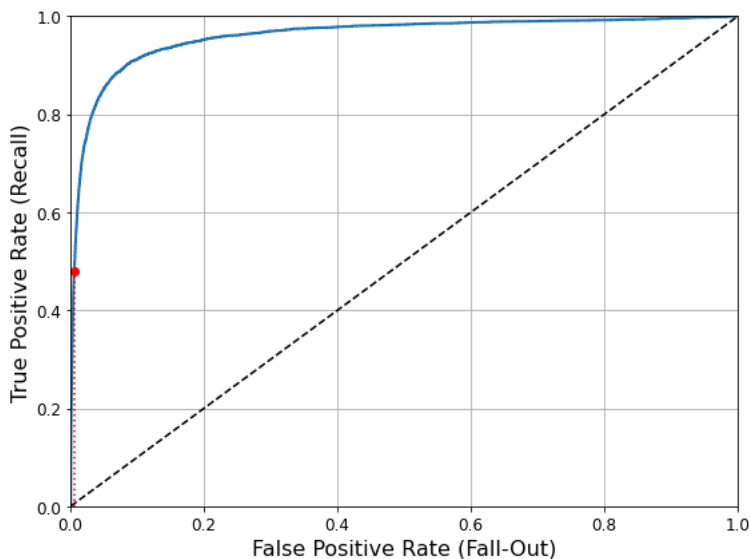
```
from sklearn.metrics import roc_curve
```

```
fpr, tpr, thresholds = roc_curve(y_train_5, y_scores)
```

```
def plot_roc_curve(fpr, tpr, label=None):
    plt.plot(fpr, tpr, linewidth=2, label=label)
    plt.plot([0, 1], [0, 1], 'k--') # dashed diagonal
    plt.axis([0, 1, 0, 1]) # Not shown in the book
    plt.xlabel('False Positive Rate (Fall-Out)', fontsize=16) # Not shown
    plt.ylabel('True Positive Rate (Recall)', fontsize=16) # Not shown
    plt.grid(True) # Not shown

plt.figure(figsize=(8, 6)) # Not shown
plot_roc_curve(fpr, tpr)
fpr_90 = fpr[np.argmax(tpr >= recall_90_precision)] # Not shown
plt.plot([fpr_90, fpr_90], [0., recall_90_precision], "r:") # Not shown
plt.plot([0.0, fpr_90], [recall_90_precision, recall_90_precision], "r:") # Not shown
plt.plot([fpr_90], [recall_90_precision], "ro") # Not shown
save_fig("roc_curve_plot") # Not shown
plt.show()
```

→ Saving figure roc_curve_plot



```
from sklearn.metrics import roc_auc_score
```

```
roc_auc_score(y_train_5, y_scores)
```

→ 0.9604938554008616

Note: we set `n_estimators=100` to be future-proof since this will be the default value in Scikit-Learn 0.22.

```
from sklearn.ensemble import RandomForestClassifier
forest_clf = RandomForestClassifier(n_estimators=100, random_state=42)
y_proba_forest = cross_val_predict(forest_clf, X_train, y_train_5, cv=3,
                                   method="predict_proba")
```

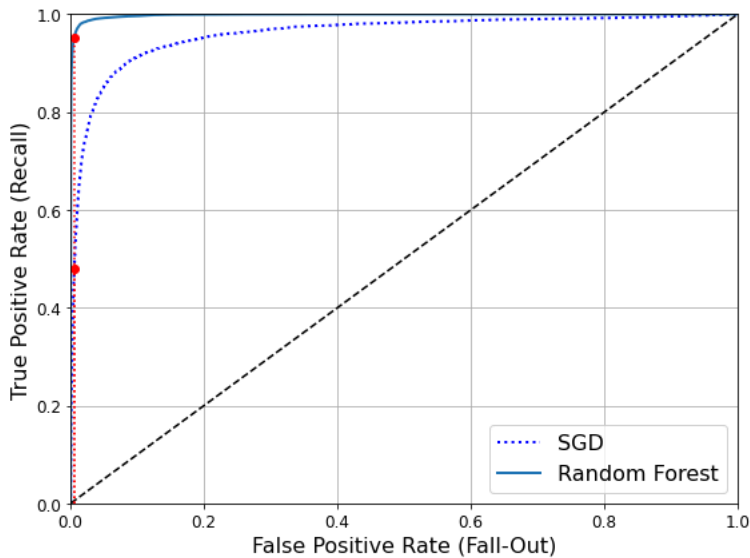
```
y_scores_forest = y_proba_forest[:, 1] # score = proba of positive class
fpr_forest, tpr_forest, thresholds_forest = roc_curve(y_train_5, y_scores_forest)
```

```
recall_for_forest = tpr_forest[np.argmax(fpr_forest >= fpr_90)]
```

```
plt.figure(figsize=(8, 6))
plt.plot(fpr, tpr, "b:", linewidth=2, label="SGD")
plot_roc_curve(fpr_forest, tpr_forest, "Random Forest")
plt.plot([fpr_90, fpr_90], [0., recall_90_precision], "r:")
plt.plot([0.0, fpr_90], [recall_90_precision, recall_90_precision], "r:")
plt.plot([fpr_90], [recall_90_precision], "ro")
plt.plot([fpr_90, fpr_90], [0., recall_for_forest], "r:")
plt.plot([fpr_90], [recall_for_forest], "ro")
plt.grid(True)
plt.legend(loc="lower right", fontsize=16)
```

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

↗ Saving figure roc_curve_comparison_plot



```
roc_auc_score(y_train_5, y_scores_forest)
```

↗ 0.9983436731328145

```
y_train_pred_forest = cross_val_predict(forest_clf, X_train, y_train_5, cv=3)
precision_score(y_train_5, y_train_pred_forest)
```

↗ 0.9905083315756169

```
recall_score(y_train_5, y_train_pred_forest)
```

↗ 0.8662608374838591

✓ Multiclass classification

```
from sklearn.svm import SVC
```

```
svm_clf = SVC(gamma="auto", random_state=42)
svm_clf.fit(X_train[:1000], y_train[:1000]) # y_train, not y_train_5
svm_clf.predict([some_digit])
```

↗ array([5], dtype=uint8)

```
some_digit_scores = svm_clf.decision_function([some_digit])
some_digit_scores
```

↗ array([[2.81585438, 7.09167958, 3.82972099, 0.79365551, 5.8885703 ,
 9.29718395, 1.79862509, 8.10392157, -0.228207 , 4.83753243]])

```
np.argmax(some_digit_scores)
```

↗ 5

```
svm_clf.classes_
```

↗ array([0, 1, 2, 3, 4, 5, 6, 7, 8, 9], dtype=uint8)

```
svm_clf.classes_[5]
```

↗ 5

```
from sklearn.multiclass import OneVsRestClassifier
ovr_clf = OneVsRestClassifier(SVC(gamma="auto", random_state=42))
ovr_clf.fit(X_train[:1000], y_train[:1000])
ovr_clf.predict([some_digit])
```

↗ array([5], dtype=uint8)


```
len(ovr_clf.estimators_)
```

```
10
```

```
sgd_clf.fit(X_train, y_train)
sgd_clf.predict([some_digit])
```

```
array([3], dtype=uint8)
```

```
sgd_clf.decision_function([some_digit])
```

```
array([[ -31893.03095419, -34419.69069632, -9530.63950739,
        1823.73154031, -22320.14822878, -1385.80478895,
        -26188.91070951, -16147.51323997, -4604.35491274,
        -12050.767298   ]])
```

```
cross_val_score(sgd_clf, X_train, y_train, cv=3, scoring="accuracy")
```

```
array([0.87365, 0.85835, 0.8689 ])
```

```
from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
```

```
X_train_scaled = scaler.fit_transform(X_train.astype(np.float64))
```

```
cross_val_score(sgd_clf, X_train_scaled, y_train, cv=3, scoring="accuracy")
```

```
array([0.8983, 0.891 , 0.9018])
```

```
y_train_pred = cross_val_predict(sgd_clf, X_train_scaled, y_train, cv=3)
```

```
conf_mx = confusion_matrix(y_train, y_train_pred)
```

```
conf_mx
```

```
array([[5577,    0,   22,    5,    8,   43,   36,    6,  225,    1],
       [    0, 6400,   37,   24,    4,   44,    4,    7,  212,   10],
       [   27,   27, 5220,   92,   73,   27,   67,   36,  378,   11],
       [   22,   17, 1117, 5227,    2,  203,   27,   40,  403,   73],
       [   12,   14,   41,    9, 5182,   12,   34,   27,  347,  164],
       [   27,   15,   30,  168,   53, 4444,   75,   14,  535,   60],
       [   30,   15,   42,    3,   44,   97, 5552,    3,  131,    1],
       [   21,   10,   51,   30,   49,   12,    3, 5684,  195,  210],
       [   17,   63,   48,   86,    3,  126,   25,   10, 5429,   44],
       [   25,   18,   30,   64,  118,   36,    1,  179,  371, 5107]])
```

```
# since sklearn 0.22, you can use sklearn.metrics.plot_confusion_matrix()
```

```
def plot_confusion_matrix(matrix):
```

```
    """If you prefer color and a colorbar"""
```

```
    fig = plt.figure(figsize=(8,8))
```

```
    ax = fig.add_subplot(111)
```

```
    cax = ax.matshow(matrix)
```

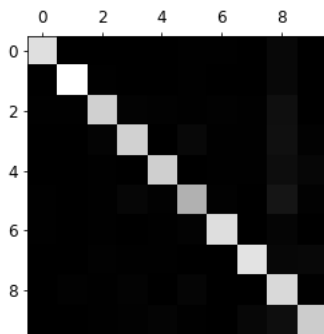
```
    fig.colorbar(cax)
```

```
plt.matshow(conf_mx, cmap=plt.cm.gray)
```

```
save_fig("confusion_matrix_plot", tight_layout=False)
```

```
plt.show()
```

```
Saving figure confusion_matrix_plot
```



```
row_sums = conf_mx.sum(axis=1, keepdims=True)
```

```
norm_conf_mx = conf_mx / row_sums
```

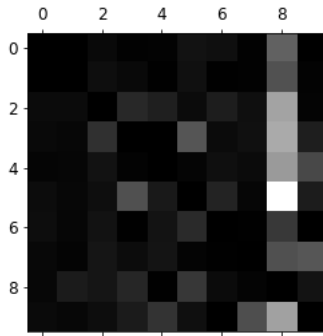
```
np.fill_diagonal(norm_conf_mx, 0)
```

```
plt.matshow(norm_conf_mx, cmap=plt.cm.gray)
```

```
save_fig("confusion_matrix_errors_plot", tight_layout=False)
```

```
plt.show()
```

↗ Saving figure confusion_matrix_errors_plot



```
cl_a, cl_b = 3, 5
X_aa = X_train[(y_train == cl_a) & (y_train_pred == cl_a)]
X_ab = X_train[(y_train == cl_a) & (y_train_pred == cl_b)]
X_ba = X_train[(y_train == cl_b) & (y_train_pred == cl_a)]
X_bb = X_train[(y_train == cl_b) & (y_train_pred == cl_b)]
```

```
plt.figure(figsize=(8,8))
plt.subplot(221); plot_digits(X_aa[:25], images_per_row=5)
plt.subplot(222); plot_digits(X_ab[:25], images_per_row=5)
plt.subplot(223); plot_digits(X_ba[:25], images_per_row=5)
plt.subplot(224); plot_digits(X_bb[:25], images_per_row=5)
save_fig("error_analysis_digits_plot")
plt.show()
```

↗ Saving figure error_analysis_digits_plot



✓ Multilabel classification

```
from sklearn.neighbors import KNeighborsClassifier
```

```
y_train_large = (y_train >= 7)
y_train_odd = (y_train % 2 == 1)
y_multilabel = np.c_[y_train_large, y_train_odd]
```

```
knn_clf = KNeighborsClassifier()
knn_clf.fit(X_train, y_multilabel)
```

↗ KNeighborsClassifier()

```
knn_clf.predict([some_digit])
```

↗ array([[False, True]])

Warning: the following cell may take a very long time (possibly hours depending on your hardware).

```
y_train_knn_pred = cross_val_predict(knn_clf, X_train, y_multilabel, cv=3)
f1_score(y_multilabel, y_train_knn_pred, average="macro")
```

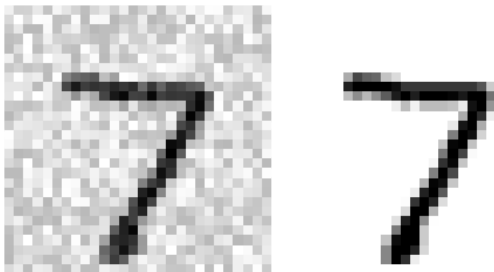
↗ 0.976410265560605

✓ Multioutput classification

```
noise = np.random.randint(0, 100, (len(X_train), 784))
X_train_mod = X_train + noise
noise = np.random.randint(0, 100, (len(X_test), 784))
X_test_mod = X_test + noise
y_train_mod = X_train
y_test_mod = X_test
```

```
some_index = 0
plt.subplot(121); plot_digit(X_test_mod[some_index])
plt.subplot(122); plot_digit(y_test_mod[some_index])
save_fig("noisy_digit_example_plot")
plt.show()
```

↗ Saving figure noisy_digit_example_plot



```
knn_clf.fit(X_train_mod, y_train_mod)
clean_digit = knn_clf.predict([X_test_mod[some_index]])
plot_digit(clean_digit)
save_fig("cleaned_digit_example_plot")
```