

## Chapter 3 – Classification

This notebook contains all the sample code and solutions to the exercises in chapter 3.



# Setup - Importing necessary Modules

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 [ ]:

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

# Is this notebook running on Colab or Kaggle?
IS_COLAB = "google.colab" in sys.modules
IS_KAGGLE = "kaggle_secrets" in sys.modules

# 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', labelsize=14)
mpl.rc('xtick', labelsize=12)
mpl.rc('ytick', labelsize=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)
```

## Importing MNIST dataset

Scikit-Learn provides many helper functions to download popular datasets. MNIST is one of them.

MNIST dataset has a set of 70,000 small images of digits handwritten by high school students and employees of the US Census Bureau.

**Warning:** since Scikit-Learn 0.24, `fetch_openml()` returns a Pandas DataFrame by default. To avoid this and keep the same code as in the book, we use `as_frame=False`.

```
In [ ]: from sklearn.datasets import fetch_openml
mnist = fetch_openml('mnist_784', version=1, as_frame=False)
mnist.keys()
```

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

```
In [ ]: # Looking at the array shapes
X, y = mnist["data"], mnist["target"]
X.shape
```

```
Out[ ]: (70000, 784)
```

```
In [ ]: y.shape
```

```
Out[ ]: (70000,)
```

There are 70000 images and each image has 784 features. This is because each image is  $28 \times 28$  pixels, and each feature simply represents one pixel's intensity, from 0 (white) to 255 (black).

```
In [ ]: 28 * 28
```

```
Out[ ]: 784
```

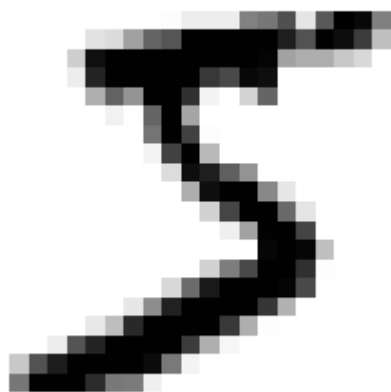
Let's take a peek at one digit from the dataset. All you need to do is grab an instance's feature vector, reshape it to a  $28 \times 28$  array, and display it using Matplotlib's `imshow()` function:

```
In [ ]: %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



This looks like a 5, and indeed that's what the label tells us:

```
In [ ]: y[0]
```

```
Out[ ]: '5'
```

Note that the label is a string. Most ML algorithms expect numbers, so let's cast y to integer:

```
In [ ]: y = y.astype(np.uint8)
```

Plotting the images

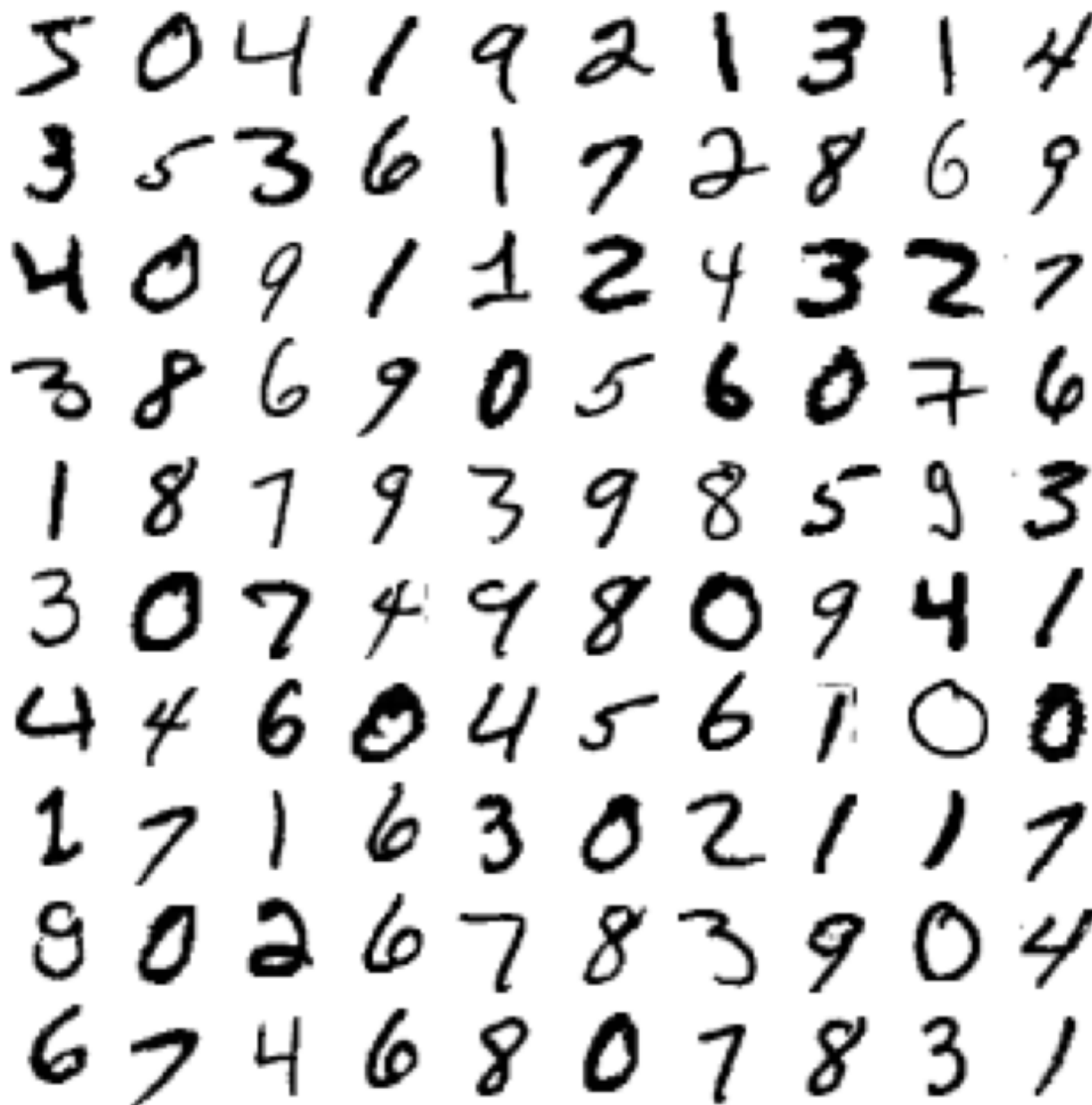
```
In [ ]: def plot_digit(data):
         image = data.reshape(28, 28)
         plt.imshow(image, cmap = mpl.cm.binary,
                    interpolation="nearest")
         plt.axis("off")
```

```
In [ ]: # 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")
```

```
In [ ]: 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



```
In [ ]: # To check if string is converting to integer
        y[0]
```

Out[ ]: 5

You should always create a test set and set it aside before inspecting the data closely. The MNIST dataset is actually already split into a training set (the first 60,000 images) and a test set (the last 10,000 images):

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

The training set is already shuffled for us, which is good because this guarantees that all cross-

validation folds will be similar (you don't want one fold to be missing some digits). Moreover, some learning algorithms are sensitive to the order of the training instances, and they perform poorly if they get many similar instances in a row. Shuffling the dataset ensures that this won't happen.

## Binary classifier

Let's simplify the problem for now and only try to identify one digit—for example, the number 5. This "5-detector" will be an example of a binary classifier, capable of distinguishing between just two classes, 5 and not-5. Let's create the target vectors for this classification task:

```
In [ ]: 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.

Now let's pick a classifier and train it. A good place to start is with a Stochastic Gradient Descent (SGD) classifier, using Scikit-Learn's `SGDClassifier` class. This classifier has the advantage of being capable of handling very large datasets efficiently. This is in part because SGD deals with training instances independently, one at a time (which also makes SGD well suited for online learning), as we will see later. Let's create an `SGDClassifier` and train it on the whole training set:

```
In [ ]: 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)
```

```
Out[ ]: SGDClassifier(alpha=0.0001, average=False, class_weight=None,
                    early_stopping=False, epsilon=0.1, eta0=0.0, fit_intercept=True,
                    l1_ratio=0.15, learning_rate='optimal', loss='hinge',
                    max_iter=1000, n_iter_no_change=5, n_jobs=None, penalty='l2',
                    power_t=0.5, random_state=42, shuffle=True, tol=0.001,
                    validation_fraction=0.1, verbose=0, warm_start=False)
```

Now we can use it to detect images of the number 5:

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

```
Out[ ]: array([ True])
```

The classifier guesses that this image represents a 5 (True). Looks like it guessed right in this particular case! Now, let's evaluate this model's performance.

## Model Performance

Occasionally you will need more control over the cross-validation process than what Scikit-Learn provides off the shelf. In these cases, you can implement cross-validation yourself. The following code does roughly the same thing as Scikit-Learn's `cross_val_score()` function, and it prints the same result:

```
In [ ]: 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
```

```
In [ ]: # Using Cross-Validation
        from sklearn.model_selection import cross_val_score
        cross_val_score(sgd_clf, X_train, y_train_5, cv=3, scoring="accuracy")
```

```
Out[ ]: array([0.95035, 0.96035, 0.9604 ])
```

Above 93% accuracy (ratio of correct predictions) on all cross-validation folds. Now, let's look at a very dumb classifier that just classifies every single image in the "not-5" class:

```
In [ ]: 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)
```

```
In [ ]: never_5_clf = Never5Classifier()
        cross_val_score(never_5_clf, X_train, y_train_5, cv=3, scoring="accuracy")
```

```
Out[ ]: array([0.91125, 0.90855, 0.90915])
```

it has over 90% accuracy! This is simply because only about 10% of the images are 5s, so if you always guess that an image is not a 5, you will be right about 90% of the time.

This demonstrates why accuracy is generally not the preferred performance measure for classifiers, especially when you are dealing with skewed datasets (i.e., when some classes are much more frequent than others)

## Confusion Matrix

A much better way to evaluate the performance of a classifier is to look at the confusion matrix.

To compute the confusion matrix, you first need to have a set of predictions so that they can be compared to the actual targets. We are using the `cross_val_predict()` function

```
In [ ]: from sklearn.model_selection import cross_val_predict

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

We can get the confusion matrix using the `confusion_matrix()` function. Just pass it the target classes (`y_train_5`) and the predicted classes (`y_train_pred`):

```
In [ ]: from sklearn.metrics import confusion_matrix

confusion_matrix(y_train_5, y_train_pred)
```

```
Out[ ]: array([[53892,   687],
               [ 1891,  3530]])
```

The first row of this matrix considers non-5 images (the negative class): 53,892 of them were correctly classified as non-5s (they are called true negatives), while the remaining 687 were wrongly classified as 5s (false positives).

The second row considers the images of 5s (the positive class): 1,891 were wrongly classified as non-5s (false negatives), while the remaining 3,530 were correctly classified as 5s (true positives). A perfect classifier would have only true positives and true negatives, so its confusion matrix would have nonzero values only on its main diagonal (top left to bottom right):

```
In [ ]: y_train_perfect_predictions = y_train_5 # pretend we reached perfection
confusion_matrix(y_train_5, y_train_perfect_predictions)
```

```
Out[ ]: array([[54579,    0],
               [    0,  5421]])
```

## Calculating Precision and Recall

```
In [ ]: from sklearn.metrics import precision_score, recall_score

precision_score(y_train_5, y_train_pred)
```

```
Out[ ]: 0.8370879772350012
```

$\text{precision} = \text{TP} / (\text{TP} + \text{FP})$

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

```
Out[ ]: 0.8370879772350012
```

```
In [ ]: recall_score(y_train_5, y_train_pred)
```

```
Out[ ]: 0.6511713705958311
```

$\text{recall} = \text{TP} / (\text{TP} + \text{FN})$

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

```
Out[ ]: 0.6511713705958311
```

Calculating F1 score - harmonic mean of precision and recall

$\text{F1} = \text{TP} / (\text{TP} + (\text{FN} + \text{FP}) / 2)$

```
In [ ]: from sklearn.metrics import f1_score
        f1_score(y_train_5, y_train_pred)
```

```
Out[ ]: 0.7325171197343846
```

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

```
Out[ ]: 0.7325171197343847
```

Precision/Recall Trade-off

Scikit-Learn does not let you set the threshold directly, but it does give you access to the decision scores that it uses to make predictions. Instead of calling the classifier's `predict()` method, you can call its `decision_function()` method, which returns a score for each instance, and then use any threshold you want to make predictions based on those scores:

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

```
Out[ ]: array([2164.22030239])
```

The `SGDClassifier` uses a threshold equal to 0, so the previous code returns the same result as the `predict()` method (i.e., `True`).

```
In [ ]: threshold = 0
        y_some_digit_pred = (y_scores > threshold)
```

```
In [ ]: y_some_digit_pred
```

```
Out[ ]: array([ True])
```

Let's raise the threshold:

```
In [ ]: threshold = 8000
        y_some_digit_pred = (y_scores > threshold)
        y_some_digit_pred
```

```
array([False])
```



Out[ ]:

This confirms that raising the threshold decreases recall. The image actually represents a 5, and the classifier detects it when the threshold is 0, but it misses it when the threshold is increased to 8,000.

How do you decide which threshold to use? First, use the `cross_val_predict()` function to get the scores of all instances in the training set, but this time specify that you want to return decision scores instead of predictions:

```
In [ ]: y_scores = cross_val_predict(sgd_clf, X_train, y_train_5, cv=3,
                                   method="decision_function")
```

With these scores, we can use the `precision_recall_curve()` function to compute precision and recall for all possible thresholds:

```
In [ ]: from sklearn.metrics import precision_recall_curve

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

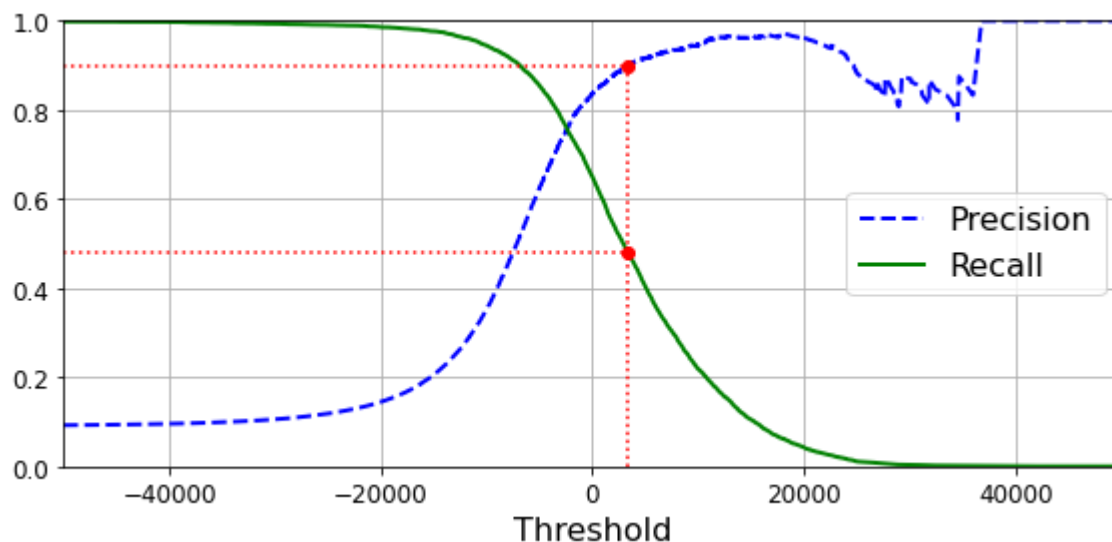
Using Matplotlib to plot precision and recall as functions of the threshold value

```
In [ ]: 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))
plot_precision_recall_vs_threshold(precisions, recalls, thresholds)
plt.plot([threshold_90_precision, threshold_90_precision], [0., 0.9], "r:")
plt.plot([-50000, threshold_90_precision], [0.9, 0.9], "r:")
plt.plot([-50000, threshold_90_precision], [recall_90_precision, recall_90_precision],
plt.plot([threshold_90_precision], [0.9], "ro")
plt.plot([threshold_90_precision], [recall_90_precision], "ro")
save_fig("precision_recall_vs_threshold_plot")
plt.show()
```

Saving figure `precision_recall_vs_threshold_plot`



```
In [ ]: (y_train_pred == (y_scores > 0)).all()
```

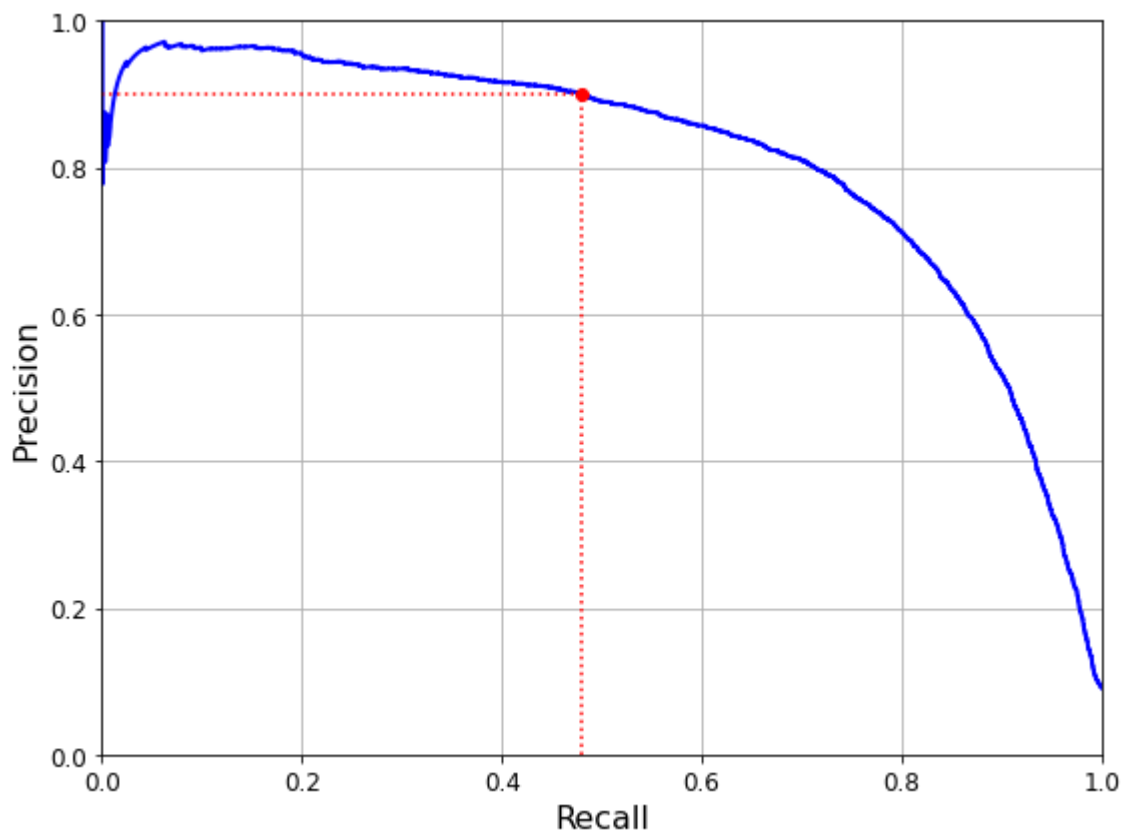
```
Out[ ]: True
```

We can select a good precision/recall trade-off is to plot precision directly against recall.

```
In [ ]: 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



Suppose you decide to aim for 90% precision. You look up the first plot and find that you need to use a threshold of about 8,000. To be more precise you can search for the lowest threshold that gives you at least 90% precision (`np.argmax()` will give you the first index of the maximum value, which in this case means the first True value):

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

```
In [ ]: threshold_90_precision
```

```
Out[ ]: 3370.0194991439557
```

To make predictions (on the training set for now), instead of calling the classifier's `predict()` method, you can run this code:

```
In [ ]: y_train_pred_90 = (y_scores >= threshold_90_precision)
```

Let's check these predictions' precision and recall:

```
In [ ]: precision_score(y_train_5, y_train_pred_90)
```

```
Out[ ]: 0.9000345901072293
```

```
In [ ]: recall_score(y_train_5, y_train_pred_90)
```

```
Out[ ]: 0.4799852425751706
```

# ROC curves

The receiver operating characteristic (ROC) curve is another common tool used with binary classifiers. It is very similar to the precision/recall curve, but instead of plotting precision versus recall, the ROC curve plots the true positive rate (another name for recall) against the false positive rate (FPR).

The ROC curve plots sensitivity (recall) versus  $1 - \text{specificity}$ .

To plot the ROC curve, you first use the `roc_curve()` function to compute the TPR and FPR for various threshold values:

```
In [ ]: from sklearn.metrics import roc_curve

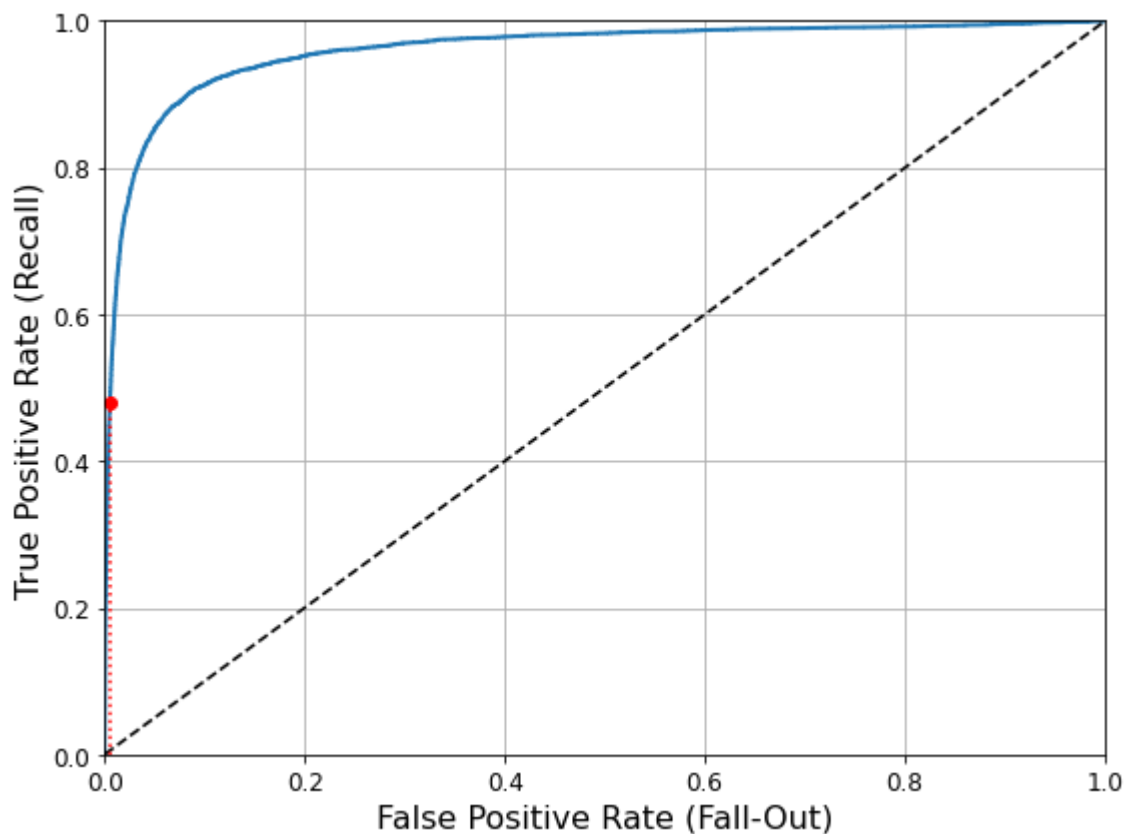
        fpr, tpr, thresholds = roc_curve(y_train_5, y_scores)
```

Then you can plot the FPR against the TPR using Matplotlib

```
In [ ]: 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



One way to compare classifiers is to measure the area under the curve (AUC). A perfect classifier will have a ROC AUC equal to 1, whereas a purely random classifier will have a ROC AUC equal to 0.5. Scikit-Learn provides a function to compute the ROC AUC:

```
In [ ]: from sklearn.metrics import roc_auc_score

        roc_auc_score(y_train_5, y_scores)
```

```
Out[ ]: 0.9604938554008616
```

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

Let's now train a `RandomForestClassifier` and compare its ROC curve and ROC AUC score to those of the `SGDClassifier`

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

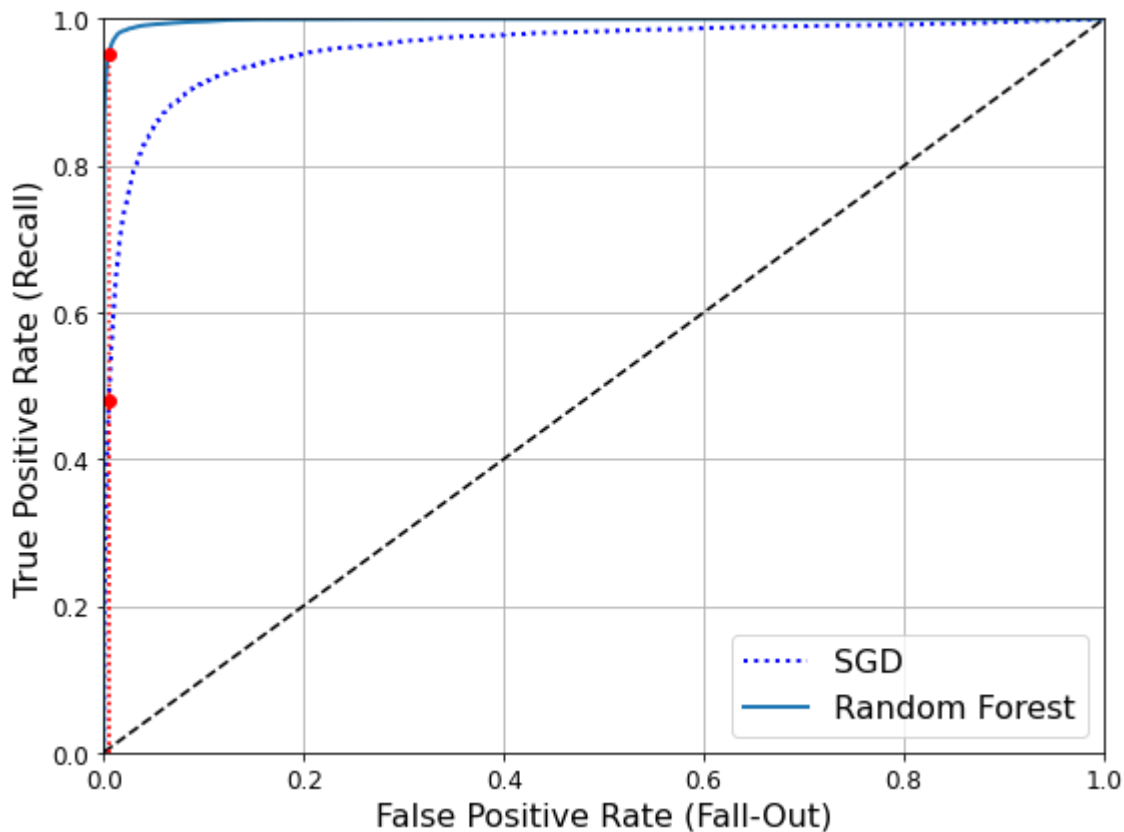
```
In [ ]: y_scores_forest = y_probas_forest[:, 1] # score = proba of positive class
        fpr_forest, tpr_forest, thresholds_forest = roc_curve(y_train_5, y_scores_forest)
```

```
In [ ]: 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



As you can see, the RandomForestClassifier's ROC curve looks much better than the SGDClassifier's: it comes much closer to the top-left corner. As a result, its ROC AUC score is also significantly better

```
In [ ]: roc_auc_score(y_train_5, y_scores_forest)
```

```
Out[ ]: 0.9983436731328145
```

Measuring Precision and Recall

```
In [ ]: 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)
```

```
Out[ ]: 0.9905083315756169
```

```
In [ ]: recall_score(y_train_5, y_train_pred_forest)
```

```
Out[ ]: 0.8662608374838591
```

# Multiclass classification

Binary classifiers distinguish between two classes, multiclass classifiers (also called multinomial classifiers) can distinguish between more than two classes.

Scikit-Learn detects when you try to use a binary classification algorithm for a multiclass classification task, and it automatically runs OvR (One-versus-the-rest) or OvO (One-versus-One), depending on the algorithm.

```
In [ ]: 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])
```

```
Out[ ]: array([5], dtype=uint8)
```

When you call the `decision_function()` method, you will see that it returns 10 scores per instance (instead of just 1). That's one score per class.

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

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

The highest score is indeed the one corresponding to class 5:

```
In [ ]: np.argmax(some_digit_scores)
```

```
Out[ ]: 5
```

```
In [ ]: svm_clf.classes_
```

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

```
In [ ]: svm_clf.classes_[5]
```

```
Out[ ]: 5
```

If you want to force Scikit-Learn to use one-versus-one or one-versus-the-rest, you can use the `OneVsOneClassifier` or `OneVsRestClassifier` classes.

```
In [ ]: 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])
```

```
Out[ ]: array([5], dtype=uint8)
```

```
In [ ]: len(ovr_clf.estimators_)
```

```
Out[ ]: 10
```

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

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

This time Scikit-Learn used the OvR strategy under the hood: since there are 10 classes, it trained 10 binary classifiers. The `decision_function()` method now returns one value per class.

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

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

**Warning:** the following two cells may take close to 30 minutes to run, or more depending on your hardware.

Using the `cross_val_score()` function to evaluate the `SGDClassifier`'s accuracy:

```
In [ ]: cross_val_score(sgd_clf, X_train, y_train, cv=3, scoring="accuracy")
```

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

It gets over 84% on all test folds. If you used a random classifier, you would get 10% accuracy, so this is not such a bad score, but you can still do much better. Simply scaling the inputs (as discussed in Chapter 2) increases accuracy above 89%:

```
In [ ]: 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")
```

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

## Error analysis

First, look at the confusion matrix. You need to make predictions using the `cross_val_predict()` function, then call the `confusion_matrix()` function

```
In [ ]: 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
```



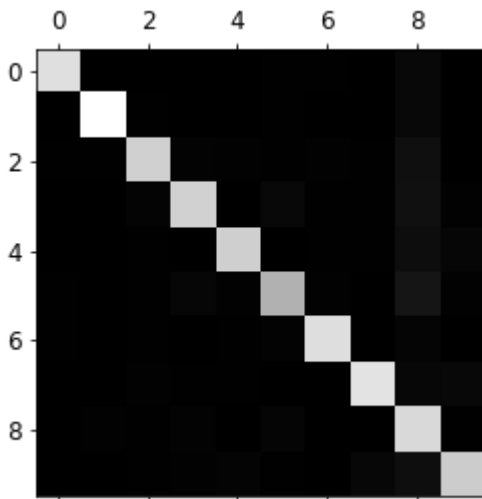
```
Out[ ]: array([[5577,  0, 22,  5,  8, 43, 36,  6, 225,  1],
               [  0, 6400, 37, 24,  4, 44,  4,  7, 212, 10],
               [ 27, 27, 5220, 92, 73, 27, 67, 36, 378, 11],
               [ 22, 17, 117, 5227,  2, 203, 27, 40, 403, 73],
               [ 12, 14, 41,  9, 5182, 12, 34, 27, 347, 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]])
```

Image representation of the confusion matrix, using Matplotlib's `matshow()` function:

```
In [ ]: # 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)
```

```
In [ ]: plt.matshow(conf_mx, cmap=plt.cm.gray)
save_fig("confusion_matrix_plot", tight_layout=False)
plt.show()
```

Saving figure confusion\_matrix\_plot



This confusion matrix looks pretty good, since most images are on the main diagonal, which means that they were classified correctly. The 5s look slightly darker than the other digits, which could mean that there are fewer images of 5s in the dataset or that the classifier does not perform as well on 5s as on other digits. In fact, you can verify that both are the case.

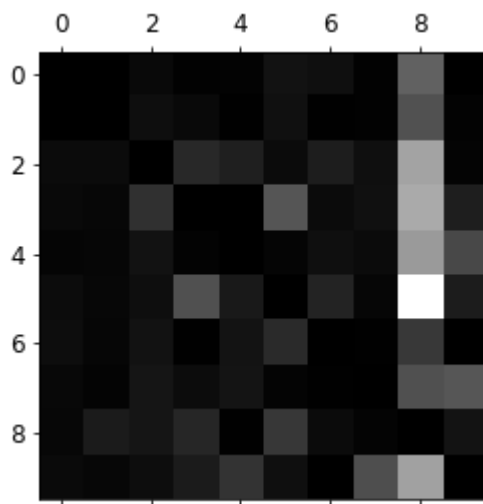
Let's focus the plot on the errors. First, you need to divide each value in the confusion matrix by the number of images in the corresponding class so that you can compare error rates instead of absolute numbers of errors.

```
In [ ]: row_sums = conf_mx.sum(axis=1, keepdims=True)
norm_conf_mx = conf_mx / row_sums
```

Fill the diagonal with zeros to keep only the errors, and plot the result:

```
In [ ]: 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



Analyzing individual errors can also be a good way to gain insights on what your classifier is doing and why it is failing, but it is more difficult and time-consuming. For example, let's plot examples of 3s and 5s - the plot\_digits() function just uses Matplotlib's imshow() function.

```
In [ ]: 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

```
In [ ]: 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)
```

```
Out[ ]: KNeighborsClassifier(algorithm='auto', leaf_size=30, metric='minkowski',
                             metric_params=None, n_jobs=None, n_neighbors=5, p=2,
                             weights='uniform')
```

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

```
Out[ ]: array([[False,  True]])
```

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

hardware).

There are many ways to evaluate a multilabel classifier, and selecting the right metric really depends on your project. One approach is to measure the F1 score for each individual label, then simply compute the average score. This code computes the average F1 score across all labels:

```
In [ ]: 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")
```

```
Out[ ]: 0.976410265560605
```

## Multioutput classification

Let's start by creating the training and test sets by taking the MNIST images and adding noise to their pixel intensities with NumPy's `randint()` function. The target images will be the original images:

```
In [ ]: 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
```

Plotting the images

```
In [ ]: 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()
```

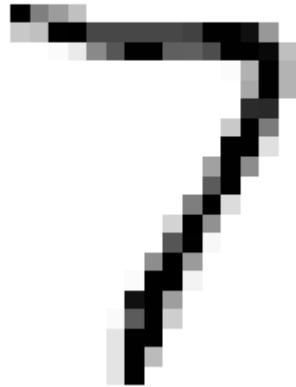
Saving figure noisy\_digit\_example\_plot



On the left is the noisy input image, and on the right is the clean target image. Now let's train the classifier and make it clean this image:

```
In [ ]: 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")
```

Saving figure cleaned\_digit\_example\_plot



## Extra material

### Dummy (ie. random) classifier

```
In [ ]: from sklearn.dummy import DummyClassifier
dmy_clf = DummyClassifier(strategy="prior")
y_probas_dmy = cross_val_predict(dmy_clf, X_train, y_train_5, cv=3, method="predict_proba")
y_scores_dmy = y_probas_dmy[:, 1]
```

```
In [ ]: fpr, tpr, threshold = roc_curve(y_train_5, y_scores_dmy)
plot_roc_curve(fpr, tpr)
```



### KNN classifier

```
In [ ]: from sklearn.neighbors import KNeighborsClassifier
knn_clf = KNeighborsClassifier(weights='distance', n_neighbors=4)
```

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

```
Out[ ]: KNeighborsClassifier(algorithm='auto', leaf_size=30, metric='minkowski',
                             metric_params=None, n_jobs=None, n_neighbors=4, p=2,
                             weights='distance')
```

```
In [ ]: y_knn_pred = knn_clf.predict(X_test)
```

```
In [ ]: from sklearn.metrics import accuracy_score
accuracy_score(y_test, y_knn_pred)
```

```
Out[ ]: 0.9714
```

```
In [ ]: from scipy.ndimage.interpolation import shift
def shift_digit(digit_array, dx, dy, new=0):
    return shift(digit_array.reshape(28, 28), [dy, dx], cval=new).reshape(784)

plot_digit(shift_digit(some_digit, 5, 1, new=100))
```



```
In [ ]: X_train_expanded = [X_train]
y_train_expanded = [y_train]
for dx, dy in ((1, 0), (-1, 0), (0, 1), (0, -1)):
    shifted_images = np.apply_along_axis(shift_digit, axis=1, arr=X_train, dx=dx, dy=dy)
    X_train_expanded.append(shifted_images)
    y_train_expanded.append(y_train)

X_train_expanded = np.concatenate(X_train_expanded)
y_train_expanded = np.concatenate(y_train_expanded)
X_train_expanded.shape, y_train_expanded.shape
```

```
Out[ ]: ((300000, 784), (300000,))
```

```
In [ ]: knn_clf.fit(X_train_expanded, y_train_expanded)
```

```
Out[ ]: KNeighborsClassifier(algorithm='auto', leaf_size=30, metric='minkowski',
                             metric_params=None, n_jobs=None, n_neighbors=4, p=2,
                             weights='distance')
```

```
In [ ]:
```

```
y_knn_expanded_pred = knn_clf.predict(X_test)
```

```
In [ ]: accuracy_score(y_test, y_knn_expanded_pred)
```

```
Out[ ]: 0.9763
```

```
In [ ]: ambiguous_digit = X_test[2589]
knn_clf.predict_proba([ambiguous_digit])
```

```
Out[ ]: array([[0.24579675, 0.          , 0.          , 0.          , 0.          ,
                0.          , 0.          , 0.          , 0.75420325]])
```

```
In [ ]: plot_digit(ambiguous_digit)
```



## Exercise solutions

### 1. An MNIST Classifier With Over 97% Accuracy

**Warning:** the next cell may take close to 16 hours to run, or more depending on your hardware.

```
In [ ]: from sklearn.model_selection import GridSearchCV

param_grid = [{'weights': ["uniform", "distance"], 'n_neighbors': [3, 4, 5]}]

knn_clf = KNeighborsClassifier()
grid_search = GridSearchCV(knn_clf, param_grid, cv=5, verbose=3)
grid_search.fit(X_train, y_train)
```

```
Fitting 5 folds for each of 6 candidates, totalling 30 fits
[CV] n_neighbors=3, weights=uniform .....
[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.
```

```
In [ ]: grid_search.best_params_
```

```
In [ ]: grid_search.best_score_
```

```
In [ ]: from sklearn.metrics import accuracy_score

y_pred = grid_search.predict(X_test)
accuracy_score(y_test, y_pred)
```

## 2. Data Augmentation

```
In [ ]: from scipy.ndimage.interpolation import shift
```

```
In [ ]: def shift_image(image, dx, dy):
        image = image.reshape((28, 28))
        shifted_image = shift(image, [dy, dx], cval=0, mode="constant")
        return shifted_image.reshape([-1])
```

```
In [ ]: image = X_train[1000]
        shifted_image_down = shift_image(image, 0, 5)
        shifted_image_left = shift_image(image, -5, 0)

        plt.figure(figsize=(12,3))
        plt.subplot(131)
        plt.title("Original", fontsize=14)
        plt.imshow(image.reshape(28, 28), interpolation="nearest", cmap="Greys")
        plt.subplot(132)
        plt.title("Shifted down", fontsize=14)
        plt.imshow(shifted_image_down.reshape(28, 28), interpolation="nearest", cmap="Greys")
        plt.subplot(133)
        plt.title("Shifted left", fontsize=14)
        plt.imshow(shifted_image_left.reshape(28, 28), interpolation="nearest", cmap="Greys")
        plt.show()
```

```
In [ ]: X_train_augmented = [image for image in X_train]
        y_train_augmented = [label for label in y_train]

        for dx, dy in ((1, 0), (-1, 0), (0, 1), (0, -1)):
            for image, label in zip(X_train, y_train):
                X_train_augmented.append(shift_image(image, dx, dy))
                y_train_augmented.append(label)

        X_train_augmented = np.array(X_train_augmented)
        y_train_augmented = np.array(y_train_augmented)
```

```
In [ ]: shuffle_idx = np.random.permutation(len(X_train_augmented))
        X_train_augmented = X_train_augmented[shuffle_idx]
        y_train_augmented = y_train_augmented[shuffle_idx]
```

```
In [ ]: knn_clf = KNeighborsClassifier(**grid_search.best_params_)
```

```
In [ ]: knn_clf.fit(X_train_augmented, y_train_augmented)
```



**Warning:** the following cell may take close to an hour to run, depending on your hardware.

```
In [ ]: y_pred = knn_clf.predict(X_test)
        accuracy_score(y_test, y_pred)
```

By simply augmenting the data, we got a 0.5% accuracy boost. :)

### 3. Tackle the Titanic dataset

The goal is to predict whether or not a passenger survived based on attributes such as their age, sex, passenger class, where they embarked and so on.

First, login to [Kaggle](#) and go to the [Titanic challenge](#) to download `train.csv` and `test.csv`. Save them to the `datasets/titanic` directory.

Next, let's load the data:

```
In [ ]: import os

        TITANIC_PATH = os.path.join("datasets", "titanic")
```

```
In [ ]: import pandas as pd

        def load_titanic_data(filename, titanic_path=TITANIC_PATH):
            csv_path = os.path.join(titanic_path, filename)
            return pd.read_csv(csv_path)
```

```
In [ ]: train_data = load_titanic_data("train.csv")
        test_data = load_titanic_data("test.csv")
```

The data is already split into a training set and a test set. However, the test data does *not* contain the labels: your goal is to train the best model you can using the training data, then make your predictions on the test data and upload them to Kaggle to see your final score.

Let's take a peek at the top few rows of the training set:

```
In [ ]: train_data.head()
```

The attributes have the following meaning:

- **Survived:** that's the target, 0 means the passenger did not survive, while 1 means he/she survived.
- **Pclass:** passenger class.
- **Name, Sex, Age:** self-explanatory
- **SibSp:** how many siblings & spouses of the passenger aboard the Titanic.
- **Parch:** how many children & parents of the passenger aboard the Titanic.
- **Ticket:** ticket id
- **Fare:** price paid (in pounds)
- **Cabin:** passenger's cabin number

- **Embarked**: where the passenger embarked the Titanic

Let's get more info to see how much data is missing:

```
In [ ]: train_data.info()
```

Okay, the **Age**, **Cabin** and **Embarked** attributes are sometimes null (less than 891 non-null), especially the **Cabin** (77% are null). We will ignore the **Cabin** for now and focus on the rest. The **Age** attribute has about 19% null values, so we will need to decide what to do with them. Replacing null values with the median age seems reasonable.

The **Name** and **Ticket** attributes may have some value, but they will be a bit tricky to convert into useful numbers that a model can consume. So for now, we will ignore them.

Let's take a look at the numerical attributes:

```
In [ ]: train_data.describe()
```

- Yikes, only 38% **Survived**. :( That's close enough to 40%, so accuracy will be a reasonable metric to evaluate our model.
- The mean **Fare** was £32.20, which does not seem so expensive (but it was probably a lot of money back then).
- The mean **Age** was less than 30 years old.

Let's check that the target is indeed 0 or 1:

```
In [ ]: train_data["Survived"].value_counts()
```

Now let's take a quick look at all the categorical attributes:

```
In [ ]: train_data["Pclass"].value_counts()
```

```
In [ ]: train_data["Sex"].value_counts()
```

```
In [ ]: train_data["Embarked"].value_counts()
```

The Embarked attribute tells us where the passenger embarked: C=Cherbourg, Q=Queenstown, S=Southampton.

**Note:** the code below uses a mix of `Pipeline`, `FeatureUnion` and a custom `DataFrameSelector` to preprocess some columns differently. Since Scikit-Learn 0.20, it is preferable to use a `ColumnTransformer`, like in the previous chapter.

Now let's build our preprocessing pipelines. We will reuse the `DataframeSelector` we built in the previous chapter to select specific attributes from the `DataFrame` :

```
In [ ]: from sklearn.base import BaseEstimator, TransformerMixin
```

```
class DataFrameSelector(BaseEstimator, TransformerMixin):
    def __init__(self, attribute_names):
        self.attribute_names = attribute_names
    def fit(self, X, y=None):
        return self
    def transform(self, X):
        return X[self.attribute_names]
```

Let's build the pipeline for the numerical attributes:

```
In [ ]: from sklearn.pipeline import Pipeline
        from sklearn.impute import SimpleImputer

        num_pipeline = Pipeline([
            ("select_numeric", DataFrameSelector(["Age", "SibSp", "Parch", "Fare"])),
            ("imputer", SimpleImputer(strategy="median")),
        ])
```

```
In [ ]: num_pipeline.fit_transform(train_data)
```

We will also need an imputer for the string categorical columns (the regular `SimpleImputer` does not work on those):

```
In [ ]: # Inspired from stackoverflow.com/questions/25239958
        class MostFrequentImputer(BaseEstimator, TransformerMixin):
            def fit(self, X, y=None):
                self.most_frequent_ = pd.Series([X[c].value_counts().index[0] for c in X],
                                                index=X.columns)
                return self
            def transform(self, X, y=None):
                return X.fillna(self.most_frequent_)
```

```
In [ ]: from sklearn.preprocessing import OneHotEncoder
```

Now we can build the pipeline for the categorical attributes:

```
In [ ]: cat_pipeline = Pipeline([
        ("select_cat", DataFrameSelector(["Pclass", "Sex", "Embarked"])),
        ("imputer", MostFrequentImputer()),
        ("cat_encoder", OneHotEncoder(sparse=False)),
    ])
```

```
In [ ]: cat_pipeline.fit_transform(train_data)
```

Finally, let's join the numerical and categorical pipelines:

```
In [ ]: from sklearn.pipeline import FeatureUnion
        preprocess_pipeline = FeatureUnion(transformer_list=[
            ("num_pipeline", num_pipeline),
            ("cat_pipeline", cat_pipeline),
        ])
```

Cool! Now we have a nice preprocessing pipeline that takes the raw data and outputs numerical input features that we can feed to any Machine Learning model we want.

```
In [ ]: X_train = preprocess_pipeline.fit_transform(train_data)
X_train
```

Let's not forget to get the labels:

```
In [ ]: y_train = train_data["Survived"]
```

We are now ready to train a classifier. Let's start with an SVC :

```
In [ ]: from sklearn.svm import SVC

svm_clf = SVC(gamma="auto")
svm_clf.fit(X_train, y_train)
```

Great, our model is trained, let's use it to make predictions on the test set:

```
In [ ]: X_test = preprocess_pipeline.transform(test_data)
y_pred = svm_clf.predict(X_test)
```

And now we could just build a CSV file with these predictions (respecting the format expected by Kaggle), then upload it and hope for the best. But wait! We can do better than hope. Why don't we use cross-validation to have an idea of how good our model is?

```
In [ ]: from sklearn.model_selection import cross_val_score

svm_scores = cross_val_score(svm_clf, X_train, y_train, cv=10)
svm_scores.mean()
```

Okay, over 73% accuracy, clearly better than random chance, but it's not a great score. Looking at the [leaderboard](#) for the Titanic competition on Kaggle, you can see that you need to reach above 80% accuracy to be within the top 10% Kagglers. Some reached 100%, but since you can easily find the [list of victims](#) of the Titanic, it seems likely that there was little Machine Learning involved in their performance! ;-) So let's try to build a model that reaches 80% accuracy.

Let's try a RandomForestClassifier :

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

forest_clf = RandomForestClassifier(n_estimators=100, random_state=42)
forest_scores = cross_val_score(forest_clf, X_train, y_train, cv=10)
forest_scores.mean()
```

That's much better!

Instead of just looking at the mean accuracy across the 10 cross-validation folds, let's plot all 10 scores for each model, along with a box plot highlighting the lower and upper quartiles, and "whiskers" showing the extent of the scores (thanks to Nevin Yilmaz for suggesting this

visualization). Note that the `boxplot()` function detects outliers (called "fliers") and does not include them within the whiskers. Specifically, if the lower quartile is  $Q_1$  and the upper quartile is  $Q_3$ , then the interquartile range  $IQR = Q_3 - Q_1$  (this is the box's height), and any score lower than  $Q_1 - 1.5 \times IQR$  is a flier, and so is any score greater than  $Q_3 + 1.5 \times IQR$ .

```
In [ ]: plt.figure(figsize=(8, 4))
plt.plot([1]*10, svm_scores, ".")
plt.plot([2]*10, forest_scores, ".")
plt.boxplot([svm_scores, forest_scores], labels=("SVM", "Random Forest"))
plt.ylabel("Accuracy", fontsize=14)
plt.show()
```

To improve this result further, you could:

- Compare many more models and tune hyperparameters using cross validation and grid search,
- Do more feature engineering, for example:
  - replace **SibSp** and **Parch** with their sum,
  - try to identify parts of names that correlate well with the **Survived** attribute (e.g. if the name contains "Countess", then survival seems more likely),
- try to convert numerical attributes to categorical attributes: for example, different age groups had very different survival rates (see below), so it may help to create an age bucket category and use it instead of the age. Similarly, it may be useful to have a special category for people traveling alone since only 30% of them survived (see below).

```
In [ ]: train_data["AgeBucket"] = train_data["Age"] // 15 * 15
train_data[["AgeBucket", "Survived"]].groupby(['AgeBucket']).mean()
```

```
In [ ]: train_data["RelativesOnboard"] = train_data["SibSp"] + train_data["Parch"]
train_data[["RelativesOnboard", "Survived"]].groupby(['RelativesOnboard']).mean()
```

## 4. Spam classifier

First, let's fetch the data:

```
In [ ]: import os
import tarfile
import urllib.request

DOWNLOAD_ROOT = "http://spamassassin.apache.org/old/publiccorpus/"
HAM_URL = DOWNLOAD_ROOT + "20030228_easy_ham.tar.bz2"
SPAM_URL = DOWNLOAD_ROOT + "20030228_spam.tar.bz2"
SPAM_PATH = os.path.join("datasets", "spam")

def fetch_spam_data(ham_url=HAM_URL, spam_url=SPAM_URL, spam_path=SPAM_PATH):
    if not os.path.isdir(spam_path):
        os.makedirs(spam_path)
    for filename, url in (("ham.tar.bz2", ham_url), ("spam.tar.bz2", spam_url)):
        path = os.path.join(spam_path, filename)
        if not os.path.isfile(path):
            urllib.request.urlretrieve(url, path)
        tar_bz2_file = tarfile.open(path)
```

```
tar_bz2_file.extractall(path=spam_path)
tar_bz2_file.close()
```

```
In [ ]: fetch_spam_data()
```

Next, let's load all the emails:

```
In [ ]: HAM_DIR = os.path.join(SPAM_PATH, "easy_ham")
        SPAM_DIR = os.path.join(SPAM_PATH, "spam")
        ham_filenames = [name for name in sorted(os.listdir(HAM_DIR)) if len(name) > 20]
        spam_filenames = [name for name in sorted(os.listdir(SPAM_DIR)) if len(name) > 20]
```

```
In [ ]: len(ham_filenames)
```

```
In [ ]: len(spam_filenames)
```

We can use Python's `email` module to parse these emails (this handles headers, encoding, and so on):

```
In [ ]: import email
        import email.policy

        def load_email(is_spam, filename, spam_path=SPAM_PATH):
            directory = "spam" if is_spam else "easy_ham"
            with open(os.path.join(spam_path, directory, filename), "rb") as f:
                return email.parser.BytesParser(policy=email.policy.default).parse(f)
```

```
In [ ]: ham_emails = [load_email(is_spam=False, filename=name) for name in ham_filenames]
        spam_emails = [load_email(is_spam=True, filename=name) for name in spam_filenames]
```

Let's look at one example of ham and one example of spam, to get a feel of what the data looks like:

```
In [ ]: print(ham_emails[1].get_content().strip())
```

```
In [ ]: print(spam_emails[6].get_content().strip())
```

Some emails are actually multipart, with images and attachments (which can have their own attachments). Let's look at the various types of structures we have:

```
In [ ]: def get_email_structure(email):
        if isinstance(email, str):
            return email
        payload = email.get_payload()
        if isinstance(payload, list):
            return "multipart({})".format(", ".join([
                get_email_structure(sub_email)
                for sub_email in payload
            ]))
```

```

else:
    return email.get_content_type()

```

```

In [ ]: from collections import Counter

def structures_counter(emails):
    structures = Counter()
    for email in emails:
        structure = get_email_structure(email)
        structures[structure] += 1
    return structures

```

```

In [ ]: structures_counter(ham_emails).most_common()

```

```

In [ ]: structures_counter(spam_emails).most_common()

```

It seems that the ham emails are more often plain text, while spam has quite a lot of HTML.

Moreover, quite a few ham emails are signed using PGP, while no spam is. In short, it seems that the email structure is useful information to have.

Now let's take a look at the email headers:

```

In [ ]: for header, value in spam_emails[0].items():
        print(header, ":", value)

```

There's probably a lot of useful information in there, such as the sender's email address (12a1mailbot1@web.de looks fishy), but we will just focus on the `Subject` header:

```

In [ ]: spam_emails[0]["Subject"]

```

Okay, before we learn too much about the data, let's not forget to split it into a training set and a test set:

```

In [ ]: import numpy as np
        from sklearn.model_selection import train_test_split

        X = np.array(ham_emails + spam_emails, dtype=object)
        y = np.array([0] * len(ham_emails) + [1] * len(spam_emails))

        X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=4

```

Okay, let's start writing the preprocessing functions. First, we will need a function to convert HTML to plain text. Arguably the best way to do this would be to use the great [BeautifulSoup](#) library, but I would like to avoid adding another dependency to this project, so let's hack a quick & dirty solution using regular expressions (at the risk of [unholy radiance destroying all enlightenment](#)). The following function first drops the `<head>` section, then converts all `<a>` tags to the word `HYPERLINK`, then it gets rid of all HTML tags, leaving only the plain text. For readability, it also replaces multiple newlines with single newlines, and finally it unescapes html entities (such as `&gt;` or `&nbsp;`):

```
In [ ]: import re
        from html import unescape

        def html_to_plain_text(html):
            text = re.sub('<head.*?>.*?</head>', '', html, flags=re.M | re.S | re.I)
            text = re.sub('<a\s.*?>', ' HYPERLINK ', text, flags=re.M | re.S | re.I)
            text = re.sub('<.*?>', '', text, flags=re.M | re.S)
            text = re.sub(r'(\s*\n)+', '\n', text, flags=re.M | re.S)
            return unescape(text)
```

Let's see if it works. This is HTML spam:

```
In [ ]: html_spam_emails = [email for email in X_train[y_train==1]
                        if get_email_structure(email) == "text/html"]
        sample_html_spam = html_spam_emails[7]
        print(sample_html_spam.get_content().strip()[:1000], "...")
```

And this is the resulting plain text:

```
In [ ]: print(html_to_plain_text(sample_html_spam.get_content())[:1000], "...")
```

Great! Now let's write a function that takes an email as input and returns its content as plain text, whatever its format is:

```
In [ ]: def email_to_text(email):
        html = None
        for part in email.walk():
            ctype = part.get_content_type()
            if not ctype in ("text/plain", "text/html"):
                continue
            try:
                content = part.get_content()
            except: # in case of encoding issues
                content = str(part.get_payload())
            if ctype == "text/plain":
                return content
            else:
                html = content
        if html:
            return html_to_plain_text(html)
```

```
In [ ]: print(email_to_text(sample_html_spam)[:100], "...")
```

Let's throw in some stemming! For this to work, you need to install the Natural Language Toolkit (NLTK). It's as simple as running the following command (don't forget to activate your virtualenv first; if you don't have one, you will likely need administrator rights, or use the `--user` option):

```
$ pip3 install nltk
```

```
In [ ]: try:
        import nltk
```



```

stemmer = nltk.PorterStemmer()
for word in ("Computations", "Computation", "Computing", "Computed", "Compute", "Co
    print(word, "=>", stemmer.stem(word))
except ImportError:
    print("Error: stemming requires the NLTK module.")
stemmer = None

```

We will also need a way to replace URLs with the word "URL". For this, we could use hard core [regular expressions](#) but we will just use the [urlextract](#) library. You can install it with the following command (don't forget to activate your virtualenv first; if you don't have one, you will likely need administrator rights, or use the `--user` option):

```
$ pip3 install urlextract
```

```

In [ ]: # if running this notebook on Colab or Kaggle, we just pip install urlextract
if IS_COLAB or IS_KAGGLE:
    !pip install -q -U urlextract

```

```

In [ ]: try:
import urlextract # may require an Internet connection to download root domain name

url_extractor = urlextract.URLExtract()
print(url_extractor.find_urls("Will it detect github.com and https://youtu.be/7Pq-S
except ImportError:
    print("Error: replacing URLs requires the urlextract module.")
    url_extractor = None

```

We are ready to put all this together into a transformer that we will use to convert emails to word counters. Note that we split sentences into words using Python's `split()` method, which uses whitespaces for word boundaries. This works for many written languages, but not all. For example, Chinese and Japanese scripts generally don't use spaces between words, and Vietnamese often uses spaces even between syllables. It's okay in this exercise, because the dataset is (mostly) in English.

```

In [ ]: from sklearn.base import BaseEstimator, TransformerMixin

class EmailToWordCounterTransformer(BaseEstimator, TransformerMixin):
    def __init__(self, strip_headers=True, lower_case=True, remove_punctuation=True,
        replace_urls=True, replace_numbers=True, stemming=True):
        self.strip_headers = strip_headers
        self.lower_case = lower_case
        self.remove_punctuation = remove_punctuation
        self.replace_urls = replace_urls
        self.replace_numbers = replace_numbers
        self.stemming = stemming
    def fit(self, X, y=None):
        return self
    def transform(self, X, y=None):
        X_transformed = []
        for email in X:
            text = email_to_text(email) or ""
            if self.lower_case:
                text = text.lower()
            if self.replace_urls and url_extractor is not None:
                urls = list(set(url_extractor.find_urls(text)))

```

```

        urls.sort(key=lambda url: len(url), reverse=True)
        for url in urls:
            text = text.replace(url, " URL ")
    if self.replace_numbers:
        text = re.sub(r'\d+(?:\.\d*)?(?:[eE][+-]?\d+)?', 'NUMBER', text)
    if self.remove_punctuation:
        text = re.sub(r'\W+', ' ', text, flags=re.M)
    word_counts = Counter(text.split())
    if self.stemming and stemmer is not None:
        stemmed_word_counts = Counter()
        for word, count in word_counts.items():
            stemmed_word = stemmer.stem(word)
            stemmed_word_counts[stemmed_word] += count
        word_counts = stemmed_word_counts
    X_transformed.append(word_counts)
    return np.array(X_transformed)

```

Let's try this transformer on a few emails:

```

In [ ]: X_few = X_train[:3]
        X_few_wordcounts = EmailToWordCounterTransformer().fit_transform(X_few)
        X_few_wordcounts

```

This looks about right!

Now we have the word counts, and we need to convert them to vectors. For this, we will build another transformer whose `fit()` method will build the vocabulary (an ordered list of the most common words) and whose `transform()` method will use the vocabulary to convert word counts to vectors. The output is a sparse matrix.

```

In [ ]: from scipy.sparse import csr_matrix

class WordCounterToVectorTransformer(BaseEstimator, TransformerMixin):
    def __init__(self, vocabulary_size=1000):
        self.vocabulary_size = vocabulary_size
    def fit(self, X, y=None):
        total_count = Counter()
        for word_count in X:
            for word, count in word_count.items():
                total_count[word] += min(count, 10)
        most_common = total_count.most_common()[0:self.vocabulary_size]
        self.vocabulary_ = {word: index + 1 for index, (word, count) in enumerate(most_common)}
        return self
    def transform(self, X, y=None):
        rows = []
        cols = []
        data = []
        for row, word_count in enumerate(X):
            for word, count in word_count.items():
                rows.append(row)
                cols.append(self.vocabulary_.get(word, 0))
                data.append(count)
        return csr_matrix((data, (rows, cols)), shape=(len(X), self.vocabulary_size + 1))

```

```

In [ ]: vocab_transformer = WordCounterToVectorTransformer(vocabulary_size=10)
        X_few_vectors = vocab_transformer.fit_transform(X_few_wordcounts)

```

```
X_few_vectors
```

```
In [ ]: X_few_vectors.toarray()
```

What does this matrix mean? Well, the 99 in the second row, first column, means that the second email contains 99 words that are not part of the vocabulary. The 11 next to it means that the first word in the vocabulary is present 11 times in this email. The 9 next to it means that the second word is present 9 times, and so on. You can look at the vocabulary to know which words we are talking about. The first word is "the", the second word is "of", etc.

```
In [ ]: vocab_transformer.vocabulary_
```

We are now ready to train our first spam classifier! Let's transform the whole dataset:

```
In [ ]: from sklearn.pipeline import Pipeline

preprocess_pipeline = Pipeline([
    ("email_to_wordcount", EmailToWordCounterTransformer()),
    ("wordcount_to_vector", WordCounterToVectorTransformer()),
])

X_train_transformed = preprocess_pipeline.fit_transform(X_train)
```

**Note:** to be future-proof, we set `solver="lbfgs"` since this will be the default value in Scikit-Learn 0.22.

```
In [ ]: from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import cross_val_score

log_clf = LogisticRegression(solver="lbfgs", max_iter=1000, random_state=42)
score = cross_val_score(log_clf, X_train_transformed, y_train, cv=3, verbose=3)
score.mean()
```

Over 98.5%, not bad for a first try! :) However, remember that we are using the "easy" dataset. You can try with the harder datasets, the results won't be so amazing. You would have to try multiple models, select the best ones and fine-tune them using cross-validation, and so on.

But you get the picture, so let's stop now, and just print out the precision/recall we get on the test set:

```
In [ ]: from sklearn.metrics import precision_score, recall_score

X_test_transformed = preprocess_pipeline.transform(X_test)

log_clf = LogisticRegression(solver="lbfgs", max_iter=1000, random_state=42)
log_clf.fit(X_train_transformed, y_train)

y_pred = log_clf.predict(X_test_transformed)

print("Precision: {:.2f}%".format(100 * precision_score(y_test, y_pred)))
print("Recall: {:.2f}%".format(100 * recall_score(y_test, y_pred)))
```

In [ ]: