

# 03\_classification

January 22, 2021

## Chapter 3 – Classification

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

### 1 Setup

First, let's make sure this notebook works well in both python 2 and 3, import a few common modules, ensure Matplotlib plots figures inline and prepare a function to save the figures:

```
[1]: # To support both python 2 and python 3
from __future__ import division, print_function, unicode_literals

# 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"

def save_fig(fig_id, tight_layout=True):
    path = os.path.join(PROJECT_ROOT_DIR, "images", CHAPTER_ID, fig_id + ".png")
    print("Saving figure", fig_id)
    if tight_layout:
        plt.tight_layout()
    plt.savefig(path, format='png', dpi=300)
```

## 2 MNIST

**Warning:** `fetch_mldata()` is deprecated since Scikit-Learn 0.20. You should use `fetch_openml()` instead. However, it returns the unsorted MNIST dataset, whereas `fetch_mldata()` returned the dataset sorted by target (the training set and the test set were sorted separately). In general, this is fine, but if you want to get the exact same results as before, you need to sort the dataset using the following function:

```
[2]: def sort_by_target(mnist):
    reorder_train = np.array(sorted([(target, i) for i, target in
    → enumerate(mnist.target[:60000])]))[:, 1]
    reorder_test = np.array(sorted([(target, i) for i, target in
    → enumerate(mnist.target[60000:])]))[:, 1]
    mnist.data[:60000] = mnist.data[reorder_train]
    mnist.target[:60000] = mnist.target[reorder_train]
    mnist.data[60000:] = mnist.data[reorder_test + 60000]
    mnist.target[60000:] = mnist.target[reorder_test + 60000]

[3]: try:
    from sklearn.datasets import fetch_openml
    mnist = fetch_openml('mnist_784', version=1, cache=True)
    mnist.target = mnist.target.astype(np.int8) # fetch_openml() returns
    → targets as strings
    sort_by_target(mnist) # fetch_openml() returns an unsorted dataset
except ImportError:
    from sklearn.datasets import fetch_mldata
    mnist = fetch_mldata('MNIST original')
mnist["data"], mnist["target"]

[3]: (array([[0., 0., 0., ..., 0., 0., 0.],
    [0., 0., 0., ..., 0., 0., 0.],
    [0., 0., 0., ..., 0., 0., 0.],
    ...,
    [0., 0., 0., ..., 0., 0., 0.],
    [0., 0., 0., ..., 0., 0., 0.],
    [0., 0., 0., ..., 0., 0., 0.]]),
    array([0, 0, 0, ..., 9, 9, 9], dtype=int8))

[4]: mnist.data.shape

[4]: (70000, 784)

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

[5]: (70000, 784)

[6]: y.shape

[6]: (70000, )

[7]: 28*28
```

[7]: 784

```
[8]: some_digit = X[36000]
some_digit_image = some_digit.reshape(28, 28)
plt.imshow(some_digit_image, cmap = mpl.cm.binary,
           interpolation="nearest")
plt.axis("off")

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

Saving figure some\_digit\_plot



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

```
[10]: # 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")

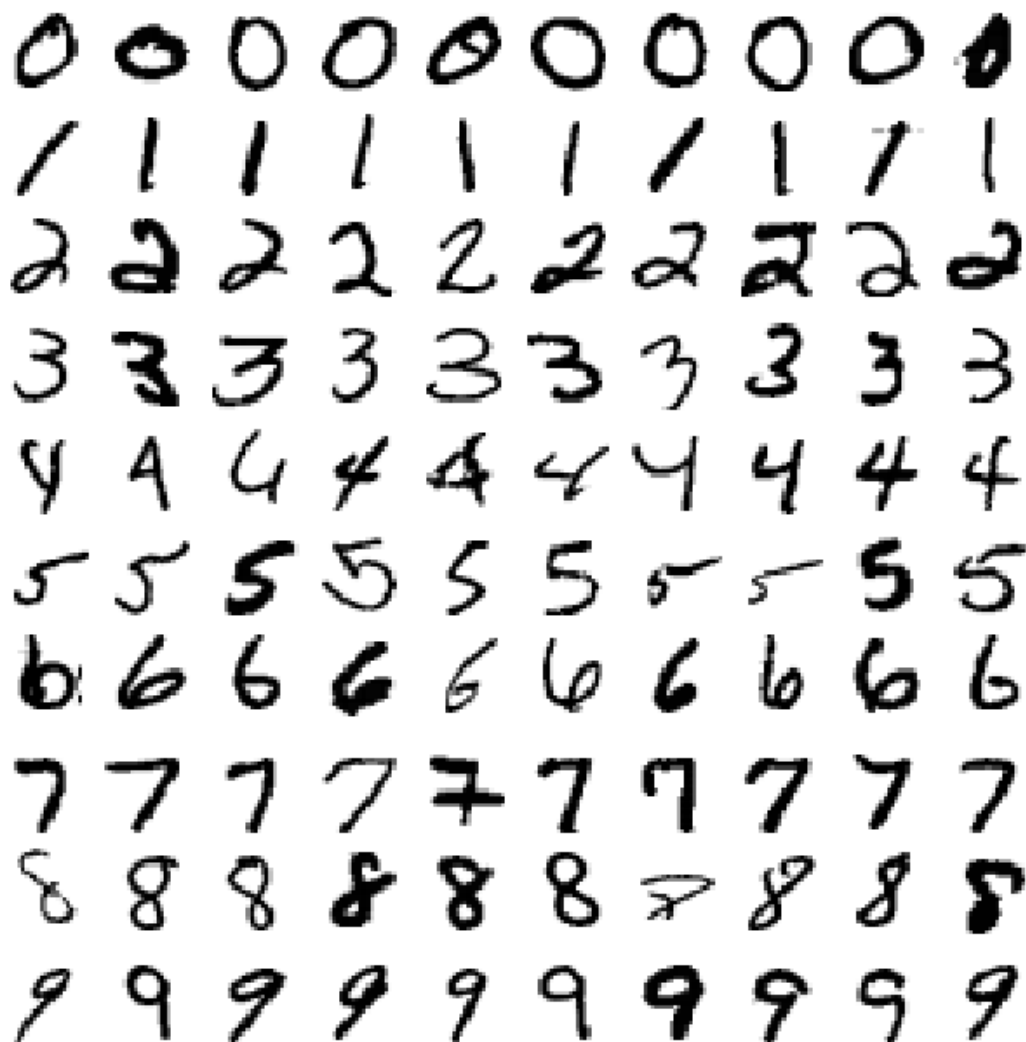
```

```

[11]: plt.figure(figsize=(9,9))
example_images = np.r_[X[:12000:600], X[13000:30600:600], X[30600:60000:590]]
plot_digits(example_images, images_per_row=10)
save_fig("more_digits_plot")
plt.show()

```

Saving figure more\_digits\_plot



```
[12]: y[36000]
```

```
[12]: 5
```

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

```
[14]: import numpy as np
```

```
shuffle_index = np.random.permutation(60000)  
X_train, y_train = X_train[shuffle_index], y_train[shuffle_index]
```

### 3 Binary classifier

```
[15]: y_train_5 = (y_train == 5)
      y_test_5 = (y_test == 5)
```

**Note:** a few hyperparameters will have a different default value in future versions of Scikit-Learn, so a warning is issued if you do not set them explicitly. This is why we set `max_iter=5` and `tol=-np.infty`, to get the same results as in the book, while avoiding the warnings.

```
[16]: from sklearn.linear_model import SGDClassifier

sgd_clf = SGDClassifier(max_iter=5, tol=-np.infty, random_state=42)
sgd_clf.fit(X_train, y_train_5)
```

```
[16]: 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=5,
                  n_iter=None, n_iter_no_change=5, n_jobs=None, penalty='l2',
                  power_t=0.5, random_state=42, shuffle=True, tol=-inf,
                  validation_fraction=0.1, verbose=0, warm_start=False)
```

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

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

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

```
[18]: array([0.9502 , 0.96565, 0.96495])
```

```
[19]: from sklearn.model_selection import StratifiedKFold
      from sklearn.base import clone

skfolds = StratifiedKFold(n_splits=3, 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))
```

```
0.9502
0.96565
0.96495
```

```
[20]: 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)
```

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

```
[21]: array([0.909 , 0.90715, 0.9128 ])
```

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

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

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

confusion_matrix(y_train_5, y_train_pred)
```

```
[23]: array([[53272, 1307],
        [ 1077, 4344]])
```

```
[24]: y_train_perfect_predictions = y_train_5
```

```
[25]: confusion_matrix(y_train_5, y_train_perfect_predictions)
```

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

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

precision_score(y_train_5, y_train_pred)
```

```
[26]: 0.7687135020350381
```

```
[27]: 4344 / (4344 + 1307)
```

```
[27]: 0.7687135020350381
```

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

```
[28]: 0.801328168234643
```

```
[29]: 4344 / (4344 + 1077)
```

```
[29]: 0.801328168234643
```

```
[30]: from sklearn.metrics import f1_score

f1_score(y_train_5, y_train_pred)
```

```
[30]: 0.7846820809248555
```

```
[31]: 4344 / (4344 + (1077 + 1307)/2)
```

```
[31]: 0.7846820809248555
```

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

```
[32]: array([161855.74572176])
```

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

```
[34]: y_some_digit_pred
```

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

```
[35]: threshold = 200000
y_some_digit_pred = (y_scores > threshold)
y_some_digit_pred
```

```
[35]: array([False])
```

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

Note: there was an [issue](#) in Scikit-Learn 0.19.0 (fixed in 0.19.1) where the result of `cross_val_predict()` was incorrect in the binary classification case when using `method="decision_function"`, as in the code above. The resulting array had an extra first dimension full of 0s. Just in case you are using 0.19.0, we need to add this small hack to work around this issue:

```
[37]: y_scores.shape
```

```
[37]: (60000,)
```

```
[38]: # hack to work around issue #9589 in Scikit-Learn 0.19.0
if y_scores.ndim == 2:
    y_scores = y_scores[:, 1]
```

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

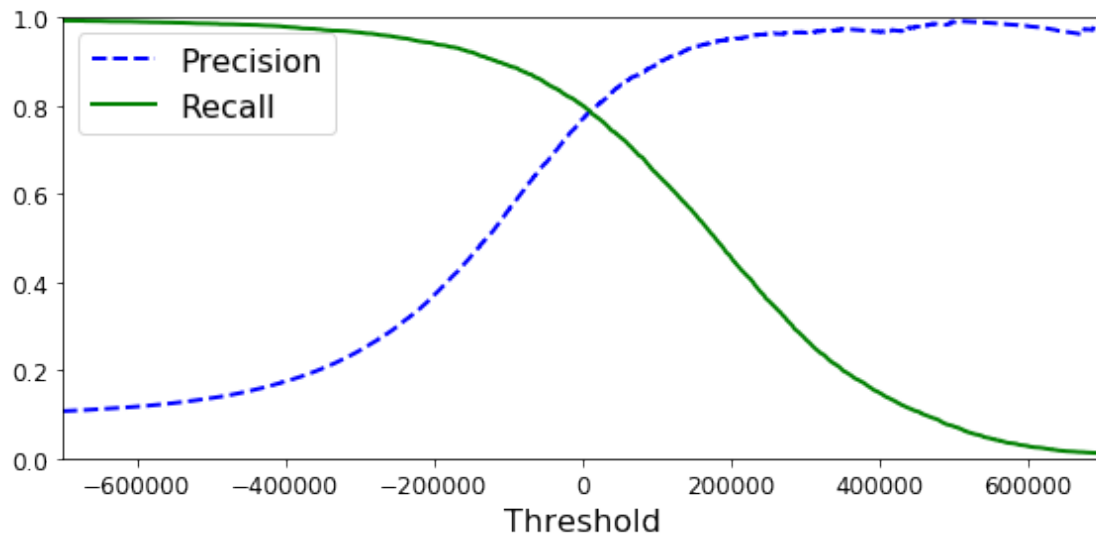
precisions, recalls, thresholds = precision_recall_curve(y_train_5, y_scores)
```

```
[40]: 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.xlabel("Threshold", fontsize=16)
    plt.legend(loc="upper left", fontsize=16)
    plt.ylim([0, 1])

plt.figure(figsize=(8, 4))
plot_precision_recall_vs_threshold(precisions, recalls, thresholds)
plt.xlim([-700000, 700000])
save_fig("precision_recall_vs_threshold_plot")
plt.show()
```

Saving figure `precision_recall_vs_threshold_plot`





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

```
[41]: True
```

```
[42]: y_train_pred_90 = (y_scores > 70000)
```

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

```
[43]: 0.8659205116491548
```

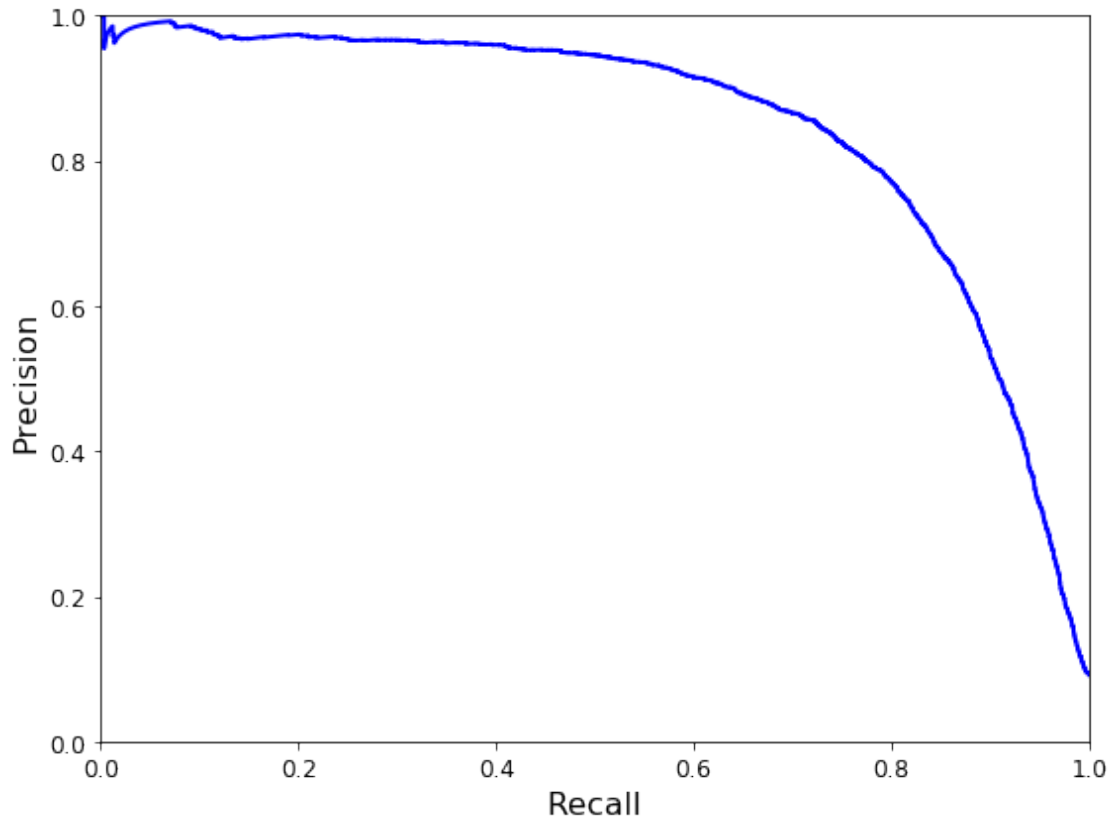
```
[44]: recall_score(y_train_5, y_train_pred_90)
```

```
[44]: 0.6993174691016417
```

```
[45]: 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.figure(figsize=(8, 6))
    plot_precision_vs_recall(precisions, recalls)
    save_fig("precision_vs_recall_plot")
    plt.show()
```

Saving figure precision\_vs\_recall\_plot



## 4 ROC curves

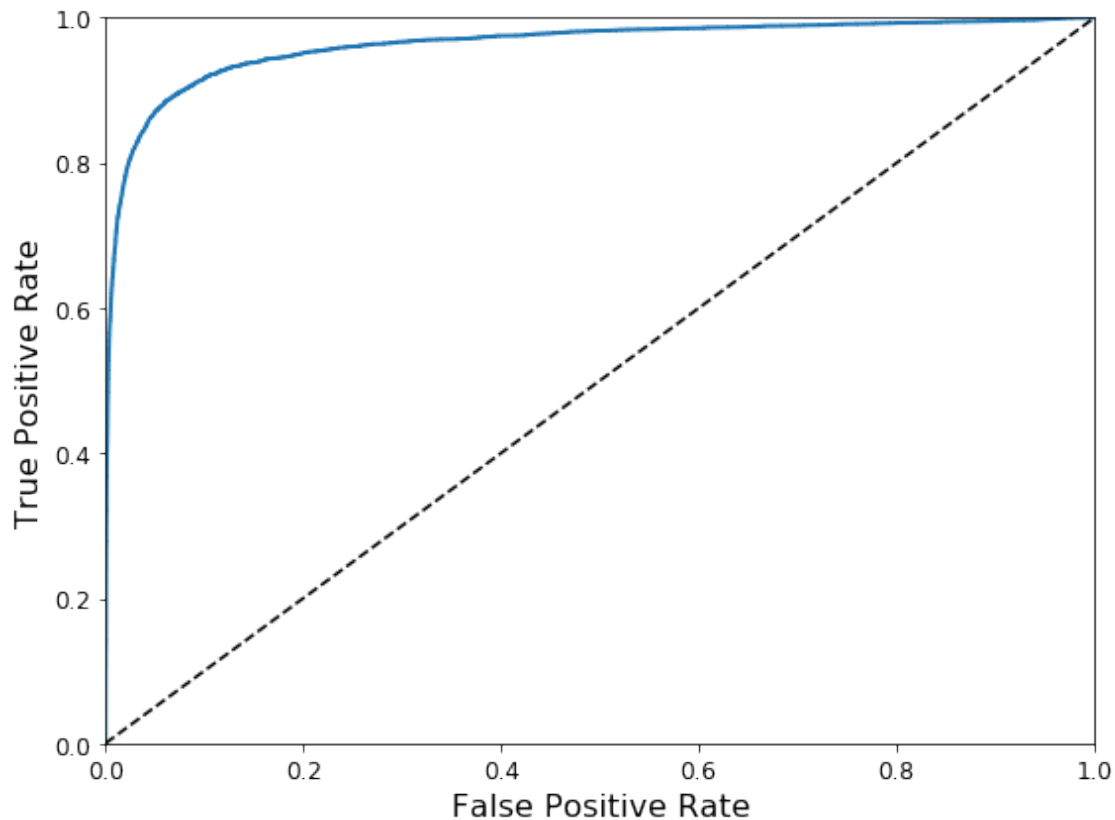
```
[46]: from sklearn.metrics import roc_curve

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

```
[47]: def plot_roc_curve(fpr, tpr, label=None):
    plt.plot(fpr, tpr, linewidth=2, label=label)
    plt.plot([0, 1], [0, 1], 'k--')
    plt.axis([0, 1, 0, 1])
    plt.xlabel('False Positive Rate', fontsize=16)
    plt.ylabel('True Positive Rate', fontsize=16)

plt.figure(figsize=(8, 6))
plot_roc_curve(fpr, tpr)
save_fig("roc_curve_plot")
plt.show()
```

Saving figure roc\_curve\_plot



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

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

```
[48]: 0.9624496555967156
```

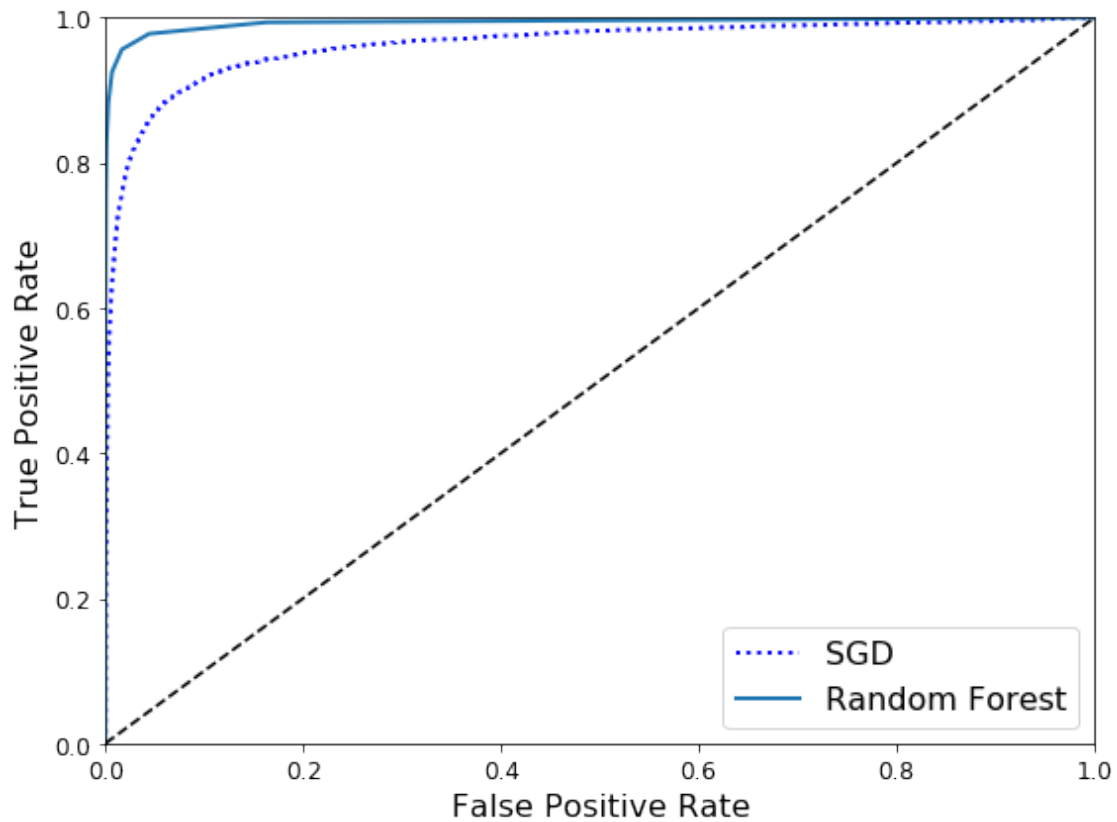
**Note:** we set `n_estimators=10` to avoid a warning about the fact that its default value will be set to 100 in Scikit-Learn 0.22.

```
[49]: from sklearn.ensemble import RandomForestClassifier
forest_clf = RandomForestClassifier(n_estimators=10, random_state=42)
y_proba_forest = cross_val_predict(forest_clf, X_train, y_train_5, cv=3,
                                   method="predict_proba")
```

```
[50]: 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)
```

```
[51]: plt.figure(figsize=(8, 6))
plt.plot(fpr, tpr, "b:", linewidth=2, label="SGD")
plot_roc_curve(fpr_forest, tpr_forest, "Random Forest")
plt.legend(loc="lower right", fontsize=16)
save_fig("roc_curve_comparison_plot")
plt.show()
```

Saving figure `roc_curve_comparison_plot`



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

```
[52]: 0.9931243366003829
```

```
[53]: 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)
```

```
[53]: 0.9852973447443494
```

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

```
[54]: 0.8282604685482383
```

## 5 Multiclass classification

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

```
[55]: array([5], dtype=int8)
```

```
[56]: some_digit_scores = sgd_clf.decision_function([some_digit])
some_digit_scores
```

```
[56]: array([[ -311402.62954431, -363517.28355739, -446449.5306454 ,  
          -183226.61023518, -414337.15339485,  161855.74572176,  
          -452576.39616343, -471957.14962573, -518542.33997148,  
          -536774.63961222]])
```

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

```
[57]: 5
```

```
[58]: sgd_clf.classes_
```

```
[58]: array([0, 1, 2, 3, 4, 5, 6, 7, 8, 9], dtype=int8)
```

```
[59]: sgd_clf.classes_[5]
```

```
[59]: 5
```

```
[60]: from sklearn.multiclass import OneVsOneClassifier  
ovo_clf = OneVsOneClassifier(SGDClassifier(max_iter=5, tol=-np.infty,  
    ↪random_state=42))  
ovo_clf.fit(X_train, y_train)  
ovo_clf.predict([some_digit])
```

```
[60]: array([5], dtype=int8)
```

```
[61]: len(ovo_clf.estimators_)
```

```
[61]: 45
```

```
[62]: forest_clf.fit(X_train, y_train)  
forest_clf.predict([some_digit])
```

```
[62]: array([5], dtype=int8)
```

```
[63]: forest_clf.predict_proba([some_digit])
```

```
[63]: array([[0.1, 0. , 0. , 0.1, 0. , 0.8, 0. , 0. , 0. , 0. ]])
```

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

```
[64]: array([0.84063187, 0.84899245, 0.86652998])
```

```
[65]: 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")
```

```
[65]: array([0.91011798, 0.90874544, 0.906636  ])
```

```
[66]: 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
```

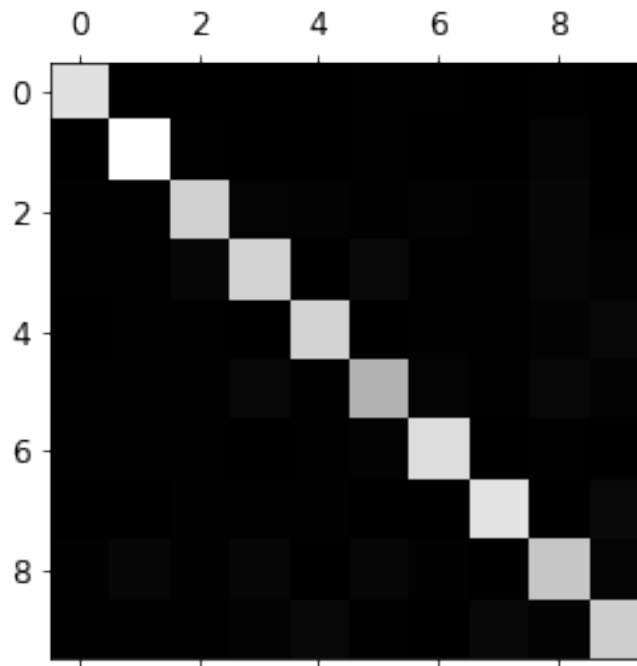
```
[66]: array([[5725,    3,   24,    9,   10,   49,   50,   10,   39,    4],  
        [    2, 6493,   43,   25,    7,   40,    5,   10,  109,    8],  
        [   51,   41, 5321,  104,   89,   26,   87,   60,  166,   13],  
        [   47,   46,  141, 5342,    1,  231,   40,   50,  141,   92],
```

```
[ 19,  29,  41,  10, 5366,   9,  56,  37,  86, 189],
[ 73,  45,  36, 193,  64, 4582, 111,  30, 193,  94],
[ 29,  34,  44,   2,  42,  85, 5627,  10,  45,   0],
[ 25,  24,  74,  32,  54,  12,   6, 5787,  15, 236],
[ 52, 161,  73, 156,  10, 163,  61,  25, 5027, 123],
[ 43,  35,  26,  92, 178,  28,   2, 223,  82, 5240]])
```

```
[67]: 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)
```

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

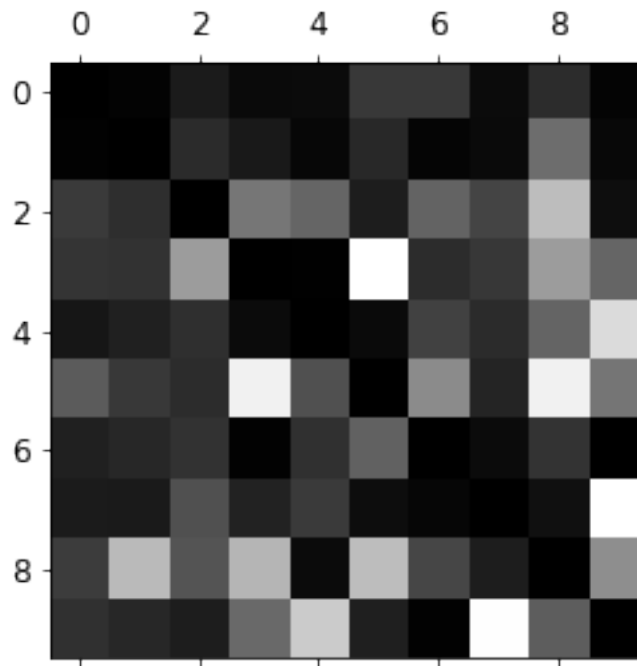
Saving figure confusion\_matrix\_plot



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

```
[70]: 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()
```

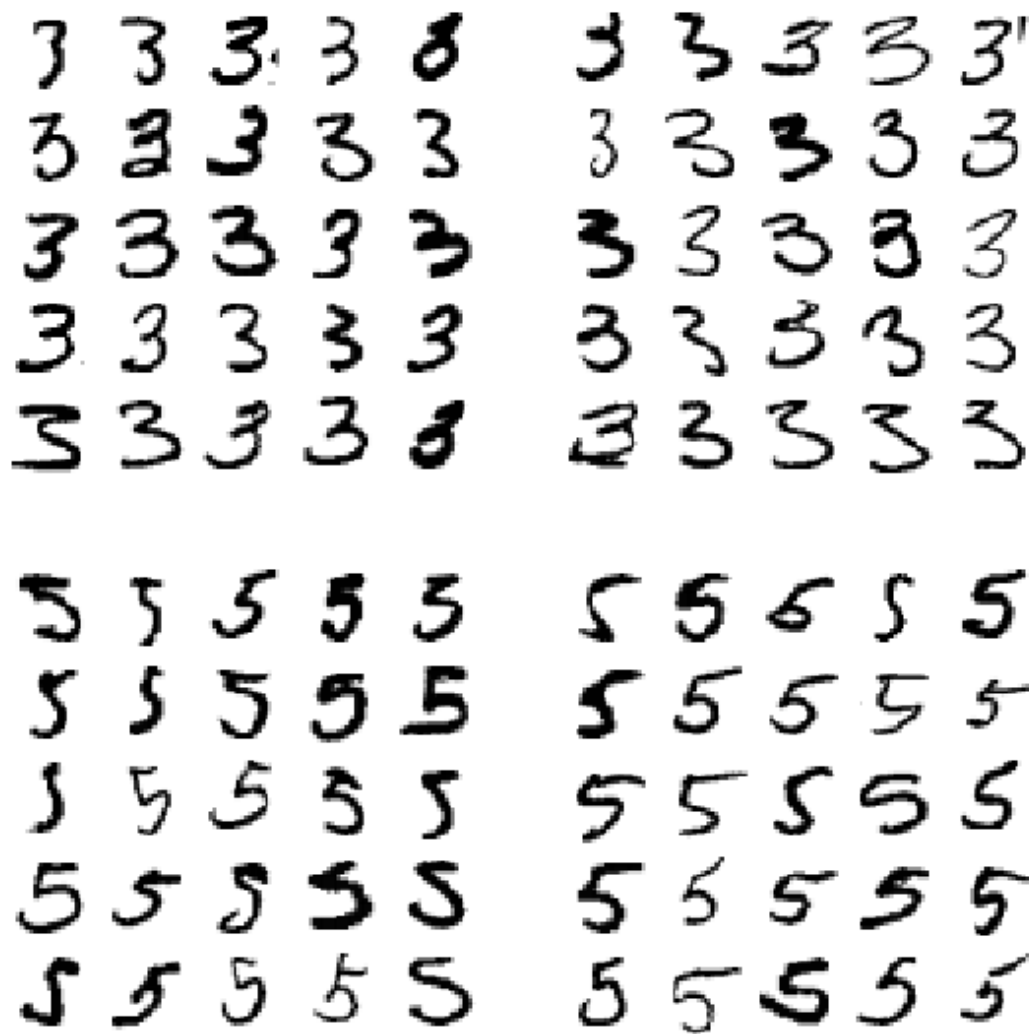
Saving figure confusion\_matrix\_errors\_plot



```
[71]: 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



## 6 Multilabel classification

```
[72]: 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)
```



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

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

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

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

```
[74]: y_train_knn_pred = cross_val_predict(knn_clf, X_train, y_multilabel, cv=3,  
    ↪n_jobs=-1)  
f1_score(y_multilabel, y_train_knn_pred, average="macro")
```

```
[74]: 0.97709078477525
```

## 7 Multioutput classification

```
[75]: 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
```

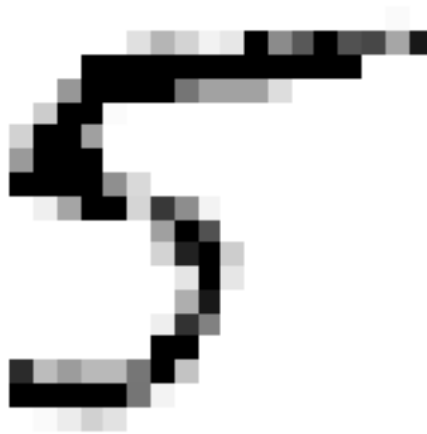
```
[76]: some_index = 5500  
plt.subplot(121); plot_digit(X_test_mod[some_index])  
plt.subplot(122); plot_digit(y_test_mod[some_index])  
save_fig("noisy_digit_example_plot")  
plt.show()
```

Saving figure noisy\_digit\_example\_plot



```
[77]: 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

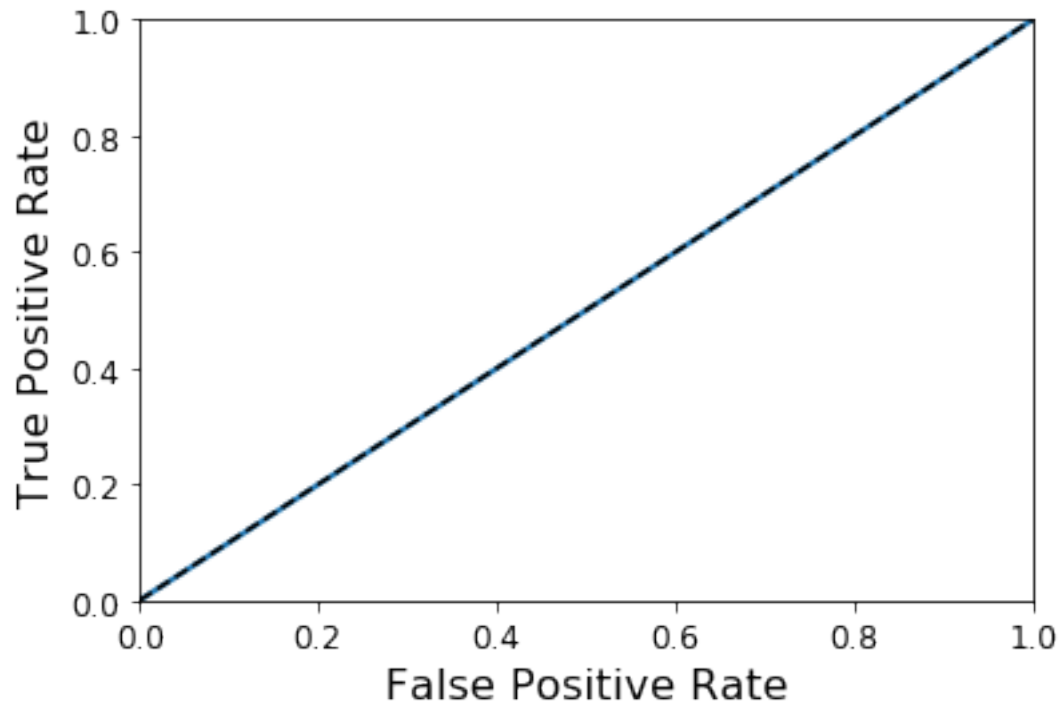


## 8 Extra material

### 8.1 Dummy (ie. random) classifier

```
[78]: from sklearn.dummy import DummyClassifier
      dmy_clf = DummyClassifier()
      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]
```

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



## 8.2 KNN classifier

```
[80]: from sklearn.neighbors import KNeighborsClassifier
knn_clf = KNeighborsClassifier(n_jobs=-1, weights='distance', n_neighbors=4)
knn_clf.fit(X_train, y_train)
```

```
[80]: KNeighborsClassifier(algorithm='auto', leaf_size=30, metric='minkowski',
metric_params=None, n_jobs=-1, n_neighbors=4, p=2,
weights='distance')
```

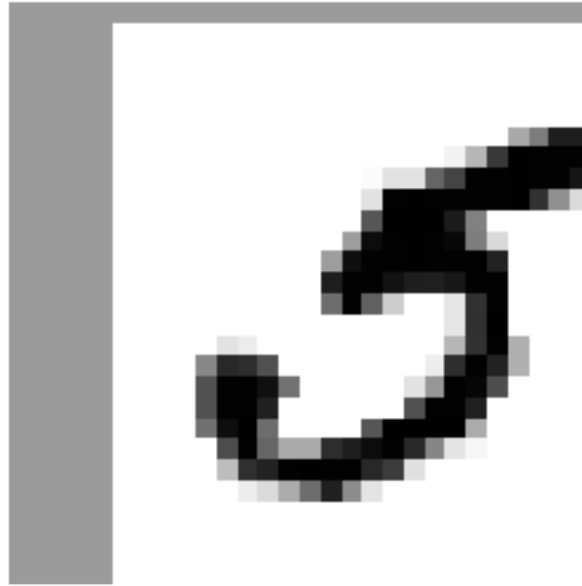
```
[81]: y_knn_pred = knn_clf.predict(X_test)
```

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

```
[82]: 0.9714
```

```
[83]: 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))
```



```
[84]: 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
```

```
[84]: ((300000, 784), (300000,))
```

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

```
[85]: KNeighborsClassifier(algorithm='auto', leaf_size=30, metric='minkowski',
    metric_params=None, n_jobs=-1, n_neighbors=4, p=2,
    weights='distance')
```

```
[86]: y_knn_expanded_pred = knn_clf.predict(X_test)
```

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

```
[87]: 0.9763
```

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

```
[88]: array([[0.          , 0.          , 0.5053645, 0.          , 0.          , 0.          ,
            0.          , 0.4946355, 0.          , 0.          ]])
```

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



## 9 Exercise solutions

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

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

```
[90]: 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, n_jobs=-1)
grid_search.fit(X_train, y_train)
```

Fitting 5 folds for each of 6 candidates, totalling 30 fits

```
[Parallel(n_jobs=-1)]: Using backend LokyBackend with 8 concurrent workers.
[Parallel(n_jobs=-1)]: Done 26 out of 30 | elapsed: 640.0min remaining:
98.5min
[Parallel(n_jobs=-1)]: Done 30 out of 30 | elapsed: 640.1min finished
```

```
[90]: GridSearchCV(cv=5, error_score='raise-deprecating',
                estimator=KNeighborsClassifier(algorithm='auto', leaf_size=30,
                metric='minkowski',
                metric_params=None, n_jobs=None, n_neighbors=5, p=2,
                weights='uniform'),
                fit_params=None, iid='warn', n_jobs=-1,
                param_grid=[{'weights': ['uniform', 'distance'], 'n_neighbors': [3, 4,
                5]}],
                pre_dispatch='2*n_jobs', refit=True, return_train_score='warn',
                scoring=None, verbose=3)
```

```
[91]: grid_search.best_params_
```

```
[91]: {'n_neighbors': 4, 'weights': 'distance'}
```

```
[92]: grid_search.best_score_
```

```
[92]: 0.97325
```

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

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

```
[93]: 0.9714
```

## 9.2 2. Data Augmentation

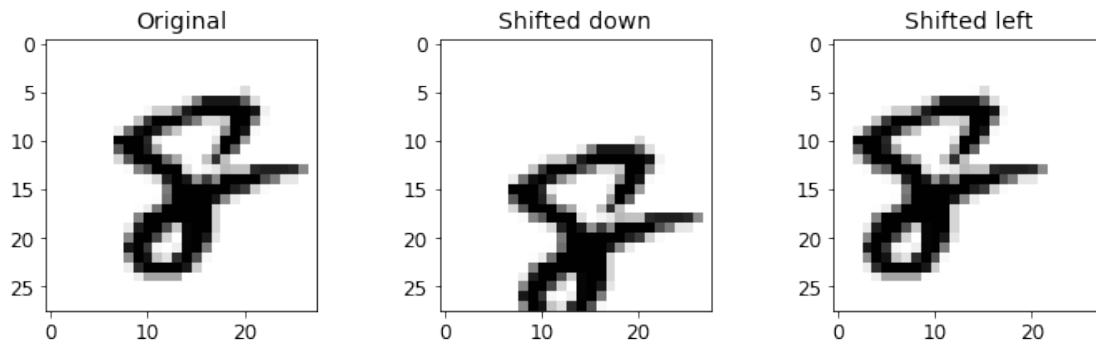
```
[94]: from scipy.ndimage.interpolation import shift
```

```
[95]: 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])
```

```
[96]: 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()
```



```
[97]: 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)
```

```
[98]: 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]
```

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

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

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

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

```
[101]: 0.9763
```

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

### 9.3 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:

```
[102]: import os

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

[103]: 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)

[104]: 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:

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

```
[105]: PassengerId  Survived  Pclass  \
0             1         0         3
1             2         1         1
2             3         1         3
3             4         1         1
4             5         0         3
```

```

                                Name      Sex  Age  SibSp  \
0                Braund, Mr. Owen Harris   male  22.0     1
1  Cumings, Mrs. John Bradley (Florence Briggs Th... female  38.0     1
2                Heikkinen, Miss. Laina   female  26.0     0
3  Futrelle, Mrs. Jacques Heath (Lily May Peel)   female  35.0     1
4                Allen, Mr. William Henry   male  35.0     0
```

```

   Parch      Ticket    Fare Cabin Embarked
0      0   A/5 21171    7.2500   NaN        S
1      0   PC 17599   71.2833   C85        C
2      0  STON/O2. 3101282   7.9250   NaN        S
3      0    113803   53.1000  C123        S
4      0    373450    8.0500   NaN        S
```

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:



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

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 891 entries, 0 to 890
Data columns (total 12 columns):
PassengerId      891 non-null int64
Survived         891 non-null int64
Pclass           891 non-null int64
Name             891 non-null object
Sex              891 non-null object
Age              714 non-null float64
SibSp            891 non-null int64
Parch           891 non-null int64
Ticket          891 non-null object
Fare             891 non-null float64
Cabin           204 non-null object
Embarked        889 non-null object
dtypes: float64(2), int64(5), object(5)
memory usage: 83.6+ KB
```

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:

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

```
[107]:
```

	PassengerId	Survived	Pclass	Age	SibSp \
count	891.000000	891.000000	891.000000	714.000000	891.000000
mean	446.000000	0.383838	2.308642	29.699118	0.523008
std	257.353842	0.486592	0.836071	14.526497	1.102743
min	1.000000	0.000000	1.000000	0.420000	0.000000
25%	223.500000	0.000000	2.000000	20.125000	0.000000
50%	446.000000	0.000000	3.000000	28.000000	0.000000
75%	668.500000	1.000000	3.000000	38.000000	1.000000
max	891.000000	1.000000	3.000000	80.000000	8.000000

	Parch	Fare
count	891.000000	891.000000
mean	0.381594	32.204208
std	0.806057	49.693429
min	0.000000	0.000000
25%	0.000000	7.910400
50%	0.000000	14.454200
75%	0.000000	31.000000
max	6.000000	512.329200

- 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:

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

```
[108]: 0    549
      1    342
      Name: Survived, dtype: int64
```

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

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

```
[109]: 3    491
      1    216
      2    184
      Name: Pclass, dtype: int64
```

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

```
[110]: male    577
      female  314
      Name: Sex, dtype: int64
```

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

```
[111]: S    644
      C    168
      Q     77
      Name: Embarked, dtype: int64
```

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

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:

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

      # A class to select numerical or categorical columns
      # since Scikit-Learn doesn't handle DataFrames yet
      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:

**Warning:** Since Scikit-Learn 0.20, the `sklearn.preprocessing.Imputer` class was replaced by the `sklearn.impute.SimpleImputer` class.

```
[113]: from sklearn.pipeline import Pipeline
try:
    from sklearn.impute import SimpleImputer # Scikit-Learn 0.20+
except ImportError:
    from sklearn.preprocessing import Imputer as SimpleImputer

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

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

```
[114]: array([[22.    ,  1.    ,  0.    ,  7.25  ],
       [38.    ,  1.    ,  0.    , 71.2833],
       [26.    ,  0.    ,  0.    ,  7.925  ],
       ...,
       [28.    ,  1.    ,  2.    , 23.45  ],
       [26.    ,  0.    ,  0.    , 30.    ],
       [32.    ,  0.    ,  0.    ,  7.75  ]])
```

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

```
[115]: # 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
    ↪X],
                                         index=X.columns)

        return self
    def transform(self, X, y=None):
        return X.fillna(self.most_frequent_)
```

**Warning:** earlier versions of the book used the `LabelBinarizer` or `CategoricalEncoder` classes to convert each categorical value to a one-hot vector. It is now preferable to use the `OneHotEncoder` class. Since Scikit-Learn 0.20 it can handle string categorical inputs (see [PR #10521](#)), not just integer categorical inputs. If you are using an older version of Scikit-Learn, you can import the new version from `future_encoders.py`:

```
[116]: try:
    from sklearn.preprocessing import OrdinalEncoder # just to raise an
    ↪ImportError if Scikit-Learn < 0.20
    from sklearn.preprocessing import OneHotEncoder
except ImportError:
    from future_encoders import OneHotEncoder # Scikit-Learn < 0.20
```

Now we can build the pipeline for the categorical attributes:

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

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

```
[118]: array([[0., 0., 1., ..., 0., 0., 1.],
 [1., 0., 0., ..., 1., 0., 0.],
 [0., 0., 1., ..., 0., 0., 1.],
 ...,
 [0., 0., 1., ..., 0., 0., 1.],
 [1., 0., 0., ..., 1., 0., 0.],
 [0., 0., 1., ..., 0., 1., 0.]])
```

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

```
[119]: 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.

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

```
[120]: array([[22., 1., 0., ..., 0., 0., 1.],
 [38., 1., 0., ..., 1., 0., 0.],
 [26., 0., 0., ..., 0., 0., 1.],
 ...,
 [28., 1., 2., ..., 0., 0., 1.],
 [26., 0., 0., ..., 1., 0., 0.],
 [32., 0., 0., ..., 0., 1., 0.]])
```

Let's not forget to get the labels:

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

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

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

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

```
[122]: SVC(C=1.0, cache_size=200, class_weight=None, coef0=0.0,
    decision_function_shape='ovr', degree=3, gamma='auto', kernel='rbf',
    max_iter=-1, probability=False, random_state=None, shrinking=True,
    tol=0.001, verbose=False)
```

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

```
[123]: 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?

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

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

```
[124]: 0.7365250822835092
```

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:

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

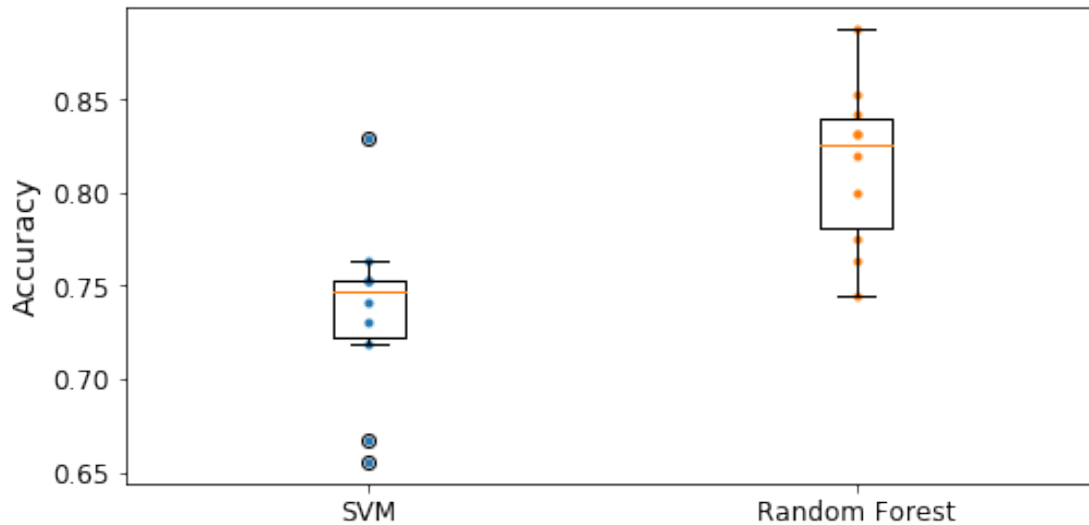
forest_clf = RandomForestClassifier(n_estimators=100, random_state=42)
forest_scores = cross_val_score(forest_clf, X_train, y_train, cv=10)
forest_scores.mean()
```

```
[125]: 0.8149526160481217
```

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  $Q_3 + 1.5 \times IQR$ .

```
[126]: 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).

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

```
[127]:      Survived
AgeBucket
0.0      0.576923
15.0     0.362745
30.0     0.423256
45.0     0.404494
60.0     0.240000
75.0     1.000000
```

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

```
[128]:      Survived
RelativesOnboard
0      0.303538
1      0.552795
2      0.578431
3      0.724138
```

4	0.200000
5	0.136364
6	0.333333
7	0.000000
10	0.000000

## 9.4 4. Spam classifier

First, let's fetch the data:

```
[129]: import os
import tarfile
from six.moves import urllib

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(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()

[130]: fetch_spam_data()
```

Next, let's load all the emails:

```
[131]: 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]

[132]: len(ham_filenames)

[132]: 2500

[133]: len(spam_filenames)

[133]: 500
```

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

```
[134]: 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)

[135]: 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:

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

Martin A posted:

Tassos Papadopoulos, the Greek sculptor behind the plan, judged that the limestone of Mount Kerdylio, 70 miles east of Salonika and not far from the Mount Athos monastic community, was ideal for the patriotic sculpture.

As well as Alexander's granite features, 240 ft high and 170 ft wide, a museum, a restored amphitheatre and car park for admiring crowds are planned

-----

So is this mountain limestone or granite?

If it's limestone, it'll weather pretty fast.

----- Yahoo! Groups Sponsor ----->

4 DVDs Free +s&p Join Now

<http://us.click.yahoo.com/pt6YBB/NXiEAA/mG3HAA/7gSolB/TM>

----->

To unsubscribe from this group, send an email to:

forteana-unsubscribe@egroups.com

Your use of Yahoo! Groups is subject to <http://docs.yahoo.com/info/terms/>

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

Help wanted. We are a 14 year old fortune 500 company, that is growing at a tremendous rate. We are looking for individuals who want to work from home.

This is an opportunity to make an excellent income. No experience



is required. We will train you.

So if you are looking to be employed from home with a career that has vast opportunities, then go:

<http://www.basetel.com/wealthnow>

We are looking for energetic and self motivated people. If that is you than click on the link and fill out the form, and one of our employment specialist will contact you.

To be removed from our link simple go to:

<http://www.basetel.com/remove.html>

4139vOLW7-758DoDY1425FRhM1-764SMFc8513fCsL140

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:

```
[138]: 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()
```

```
[139]: from collections import Counter

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

```
[140]: structures_counter(ham_emails).most_common()
```

```
[140]: [('text/plain', 2408),
        ('multipart(text/plain, application/pgp-signature)', 66),
        ('multipart(text/plain, text/html)', 8),
        ('multipart(text/plain, text/plain)', 4),
        ('multipart(text/plain)', 3),
        ('multipart(text/plain, application/octet-stream)', 2),
```

```

('multipart(text/plain, text/enriched)', 1),
('multipart(text/plain, application/ms-tnef, text/plain)', 1),
('multipart(multipart(text/plain, text/plain, text/plain), application/pgp-
signature)',
1),
('multipart(text/plain, video/mng)', 1),
('multipart(text/plain, multipart(text/plain))', 1),
('multipart(text/plain, application/x-pkcs7-signature)', 1),
('multipart(text/plain, multipart(text/plain, text/plain),
text/rfc822-headers)',
1),
('multipart(text/plain, multipart(text/plain, text/plain),
multipart(multipart(text/plain, application/x-pkcs7-signature)))',
1),
('multipart(text/plain, application/x-java-applet)', 1)]

```

```
[141]: structures_counter(spam_emails).most_common()
```

```

[141]: [('text/plain', 218),
('text/html', 183),
('multipart(text/plain, text/html)', 45),
('multipart(text/html)', 20),
('multipart(text/plain)', 19),
('multipart(multipart(text/html))', 5),
('multipart(text/plain, image/jpeg)', 3),
('multipart(text/html, application/octet-stream)', 2),
('multipart(text/plain, application/octet-stream)', 1),
('multipart(text/html, text/plain)', 1),
('multipart(multipart(text/html), application/octet-stream, image/jpeg)', 1),
('multipart(multipart(text/plain, text/html), image/gif)', 1),
('multipart/alternative', 1)]

```

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:

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

```

Return-Path : <12a1mailbot1@web.de>
Delivered-To : zzzz@localhost.spamassassin.taint.org
Received : from localhost (localhost [127.0.0.1])      by
phobos.labs.spamassassin.taint.org (Postfix) with ESMTP id 136B943C32      for
<zzzz@localhost>; Thu, 22 Aug 2002 08:17:21 -0400 (EDT)
Received : from mail.webnote.net [193.120.211.219]      by localhost with POP3
(fetchmail-5.9.0)      for zzzz@localhost (single-drop); Thu, 22 Aug 2002
13:17:21 +0100 (IST)
Received : from dd_it7 ([210.97.77.167])      by webnote.net (8.9.3/8.9.3)
with ESMTP id NAA04623      for <zzzz@spamassassin.taint.org>; Thu, 22 Aug 2002

```

```
13:09:41 +0100
From : 12a1mailbot1@web.de
Received : from r-smtp.korea.com - 203.122.2.197 by dd_it7 with Microsoft
SMTPSVC(5.5.1775.675.6); Sat, 24 Aug 2002 09:42:10 +0900
To : dcek1a1@netsgo.com
Subject : Life Insurance - Why Pay More?
Date : Wed, 21 Aug 2002 20:31:57 -1600
MIME-Version : 1.0
Message-ID : <0103c1042001882DD_IT7@dd_it7>
Content-Type : text/html; charset="iso-8859-1"
Content-Transfer-Encoding : quoted-printable
```

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:

```
[143]: spam_emails[0]["Subject"]
```

```
[143]: 'Life Insurance - Why Pay More?'
```

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

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

X = np.array(ham_emails + spam_emails)
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=42)
```

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;):

```
[145]: 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:

```
[146]: 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], "...")
```

```
<HTML><HEAD><TITLE></TITLE><META http-equiv="Content-Type" content="text/html;
charset=windows-1252"><STYLE>A:link {TEX-DECORATION: none}A:active {TEXT-
DECORATION: none}A:visited {TEXT-DECORATION: none}A:hover {COLOR: #0033ff; TEXT-
DECORATION: underline}</STYLE><META content="MSHTML 6.00.2713.1100"
name="GENERATOR"></HEAD>
<BODY text="#000000" vLink="#0033ff" link="#0033ff" bgColor="#CCCC99"><TABLE
borderColor="#660000" cellSpacing="0" cellPadding="0" border="0"
width="100%"><TR><TD bgColor="#CCCC99" valign="top" colspan="2" height="27">
<font size="6" face="Arial, Helvetica, sans-serif" color="#660000">
<b>OTC</b></font></TD></TR><TR><TD height="2" bgcolor="#6a694f">
<font size="5" face="Times New Roman, Times, serif" color="#FFFFFF">
<b>&nbsp;   Newsletter</b></font></TD><TD height="2" bgcolor="#6a694f"><div
align="right"><font color="#FFFFFF">
<b>Discover Tomorrow's Winners&nbsp;   </b></font></div></TD></TR><TR><TD
height="25" colspan="2" bgcolor="#CCCC99"><table width="100%" border="0" ...
```

And this is the resulting plain text:

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

```
OTC
ãNewsletter
Discover Tomorrow's Winners
For Immediate Release
Cal-Bay (Stock Symbol: CBYI)
Watch for analyst "Strong Buy Recommendations" and several advisory newsletters
picking CBYI. CBYI has filed to be traded on the OTCBB, share prices
historically INCREASE when companies get listed on this larger trading exchange.
CBYI is trading around 25 cents and should skyrocket to $2.66 - $3.25 a share in
the near future.
Put CBYI on your watch list, acquire a position TODAY.
REASONS TO INVEST IN CBYI
A profitable company and is on track to beat ALL earnings estimates!
One of the FASTEST growing distributors in environmental & safety equipment
instruments.
Excellent management team, several EXCLUSIVE contracts. IMPRESSIVE client list
including the U.S. Air Force, Anheuser-Busch, Chevron Refining and Mitsubishi
Heavy Industries, GE-Energy & Environmental Research.
RAPIDLY GROWING INDUSTRY
Industry revenues exceed $900 million, estimates indicate that there could be as
much as $25 billi ...
```

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

```
[148]: 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)

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

```
OTC
ãNewsletter
Discover Tomorrow's Winners
For Immediate Release
Cal-Bay (Stock Symbol: CBYI)
Wat ...
```

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
```

```
[150]: try:
        import nltk

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

```
Computations => comput
Computation => comput
```

```
Computing => comput
Computed => comput
Compute => comput
Compulsive => compuls
```

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
```

```
[151]: try:
        import urlextract # may require an Internet connection to download root
        ↪domain names

        url_extractor = urlextract.URLExtract()
        print(url_extractor.find_urls("Will it detect github.com and https://youtu.
        ↪be/7Pq-S557XQU?t=3m32s"))
    except ImportError:
        print("Error: replacing URLs requires the urlextract module.")
        url_extractor = None
```

```
['github.com', 'https://youtu.be/7Pq-S557XQU?t=3m32s']
```

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.

```
[152]: 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:

```

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

```

```

[153]: array([Counter({'chuck': 1, 'murcko': 1, 'wrote': 1, 'stuff': 1, 'yawn': 1, 'r':
1}),
Counter({'the': 11, 'of': 9, 'and': 8, 'all': 3, 'christian': 3, 'to': 3,
'by': 3, 'jefferson': 2, 'i': 2, 'have': 2, 'superstit': 2, 'one': 2, 'on': 2,
'been': 2, 'ha': 2, 'half': 2, 'rogueri': 2, 'teach': 2, 'jesu': 2, 'some': 1,
'interest': 1, 'quot': 1, 'url': 1, 'thoma': 1, 'examin': 1, 'known': 1, 'word':
1, 'do': 1, 'not': 1, 'find': 1, 'in': 1, 'our': 1, 'particular': 1, 'redeem':
1, 'featur': 1, 'they': 1, 'are': 1, 'alike': 1, 'found': 1, 'fabl': 1,
'mytholog': 1, 'million': 1, 'innoc': 1, 'men': 1, 'women': 1, 'children': 1,
'sinc': 1, 'introduc': 1, 'burnt': 1, 'tortur': 1, 'fine': 1, 'imprison': 1,
'what': 1, 'effect': 1, 'thi': 1, 'coercion': 1, 'make': 1, 'world': 1, 'fool':
1, 'other': 1, 'hypocrit': 1, 'support': 1, 'error': 1, 'over': 1, 'earth': 1,
'six': 1, 'histor': 1, 'american': 1, 'john': 1, 'e': 1, 'remsburg': 1,
'letter': 1, 'william': 1, 'short': 1, 'again': 1, 'becom': 1, 'most': 1,
'pervert': 1, 'system': 1, 'that': 1, 'ever': 1, 'shone': 1, 'man': 1, 'absurd':
1, 'untruth': 1, 'were': 1, 'perpetr': 1, 'upon': 1, 'a': 1, 'larg': 1, 'band':
1, 'dupe': 1, 'import': 1, 'led': 1, 'paul': 1, 'first': 1, 'great': 1,
'corrupt': 1}),
Counter({'url': 5, 's': 3, 'group': 3, 'to': 3, 'in': 2, 'fortean': 2,
'martin': 2, 'an': 2, 'and': 2, 'we': 2, 'is': 2, 'yahoo': 2, 'unsubscribe': 2,
'y': 1, 'adamson': 1, 'wrote': 1, 'for': 1, 'altern': 1, 'rather': 1, 'more': 1,
'factual': 1, 'base': 1, 'rundown': 1, 'on': 1, 'hamza': 1, 'career': 1,

```

```
'includ': 1, 'hi': 1, 'belief': 1, 'that': 1, 'all': 1, 'non': 1, 'muslim': 1,
'yemen': 1, 'should': 1, 'be': 1, 'murder': 1, 'outright': 1, 'know': 1, 'how':
1, 'unbias': 1, 'memri': 1, 'don': 1, 't': 1, 'html': 1, 'rob': 1, 'sponsor': 1,
'number': 1, 'dvd': 1, 'free': 1, 'p': 1, 'join': 1, 'now': 1, 'from': 1, 'thi':
1, 'send': 1, 'email': 1, 'your': 1, 'use': 1, 'of': 1, 'subject': 1}]],
dtype=object)
```

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.

```
[154]: 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()[ :self.vocabulary_size]
        self.most_common_ = most_common
        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))
```

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

```
[155]: <3x11 sparse matrix of type '<class 'numpy.int64'>'
        with 20 stored elements in Compressed Sparse Row format>
```

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



```
[156]: array([[ 6,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0],
              [99, 11,  9,  8,  1,  3,  3,  1,  3,  2,  3],
              [65,  0,  1,  2,  5,  3,  1,  2,  0,  1,  0]], dtype=int64)
```

What does this matrix mean? Well, the 64 in the third row, first column, means that the third email contains 64 words that are not part of the vocabulary. The 1 next to it means that the first word in the vocabulary is present once in this email. The 2 next to it means that the second word is present twice, and so on. You can look at the vocabulary to know which words we are talking about. The first word is "of", the second word is "and", etc.

```
[157]: vocab_transformer.vocabulary_
```

```
[157]: {'the': 1,
        'of': 2,
        'and': 3,
        'url': 4,
        'to': 5,
        'all': 6,
        'in': 7,
        'christian': 8,
        'on': 9,
        'by': 10}
```

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

```
[158]: 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)
```

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

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

```
[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.
[Parallel(n_jobs=1)]: Done 1 out of 1 | elapsed: 0.0s remaining: 0.0s
[Parallel(n_jobs=1)]: Done 2 out of 2 | elapsed: 0.1s remaining: 0.0s
```

```
[CV] ...
[CV] ... , score=0.98375, total= 0.0s
[CV] ...
[CV] ... , score=0.985, total= 0.1s
[CV] ...
[CV] ... , score=0.9925, total= 0.1s
```

```
[Parallel(n_jobs=1)]: Done   3 out of   3 | elapsed:   0.2s finished
```

```
[159]: 0.9870833333333334
```

Over 98.7%, 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:

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

X_test_transformed = preprocess_pipeline.transform(X_test)

log_clf = LogisticRegression(solver="liblinear", 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)))
```

```
Precision: 94.90%
```

```
Recall: 97.89%
```

```
[ ]:
```