

# Classification

This notebook is heavily inspired by Andre Guernon work, that can be found here:

[https://github.com/ageron/handson-ml/blob/master/04\\_training\\_linear\\_models.ipynb](https://github.com/ageron/handson-ml/blob/master/04_training_linear_models.ipynb)  
([https://github.com/ageron/handson-ml/blob/master/04\\_training\\_linear\\_models.ipynb](https://github.com/ageron/handson-ml/blob/master/04_training_linear_models.ipynb))

## Setup

```
In [7]: 1 # Python ≥3.8 is required
2 import sys
3 assert sys.version_info >= (3, 8)
4
5 # Scikit-Learn ≥1.0 is required
6 import sklearn
7 assert sklearn.__version__ >= "1.0"
8
9 # Common imports
10 import numpy as np
11 import pandas as pd
12 import os
13
14 # To plot pretty figures
15 %matplotlib inline
16 import matplotlib as mpl
17 import matplotlib.pyplot as plt
18 mpl.rcParams['axes', labelsizes=14)
19 mpl.rcParams['xtick', labelsizes=12)
20 mpl.rcParams['ytick', labelsizes=12)
21
22 from time import time
23
24 # Ignore useless warnings (see SciPy issue #5998)
25 import warnings
26 warnings.filterwarnings(action="ignore", message="^internal gel
```

## The MNIST Dataset

We will be using the MNIST dataset, which is a set of 70,000 small images of digits handwritten by high school students and employees of the US Census Bureau. Each image is labeled with the digit it represents.

We will use `sklearn.datasets.fetch_openml()` to fetch dataset from openml by name or dataset id.

```
In [8]: 1 from sklearn.datasets import fetch_openml
2 mnist = fetch_openml(
3     'mnist_784',
4     version=1,
5     as_frame=False # we want the dataset as NumPy ndarray not a
6 )
7 mnist.keys()
```

```
Out[8]: dict_keys(['data', 'target', 'frame', 'categories', 'feature_names',
', 'target_names', 'DESCR', 'details', 'url'])
```

```
In [3]: 1 mnist['DESCR']
```

```
Out[3]: """Author*: Yann LeCun, Corinna Cortes, Christopher J.C. Burges
\n**Source*: [MNIST Website](http://yann.lecun.com/exdb/mnist/) -
Date unknown \n**Please cite*: \n\nThe MNIST database of handwr
itten digits with 784 features, raw data available at: http://yann.lecun.com/exdb/mnist/. (http://yann.lecun.com/exdb/mnist/.) It ca
n be split in a training set of the first 60,000 examples, and a t
est set of 10,000 examples \n\nIt is a subset of a larger set ava
ilable from NIST. The digits have been size-normalized and centere
d in a fixed-size image. It is a good database for people who want
to try learning techniques and pattern recognition methods on real
-world data while spending minimal efforts on preprocessing and fo
rmatting. The original black and white (bilevel) images from NIST
were size normalized to fit in a 20x20 pixel box while preserving
their aspect ratio. The resulting images contain grey levels as a
result of the anti-aliasing technique used by the normalization al
gorithm. the images were centered in a 28x28 image by computing th
e center of mass of the pixels, and translating the image so as to
position this point at the center of the 28x28 field. \n\nWith so
me classification methods (particularly template-based methods, su
ch as SVM and K-nearest neighbors), the error rate improves when t
he digits are centered by bounding box rather than center of mass.
If you do this kind of pre-processing, you should report it in you
r publications. The MNIST database was constructed from NIST's NIS
T originally designated SD-3 as their training set and SD-1 as the
ir test set. However, SD-3 is much cleaner and easier to recognize
than SD-1. The reason for this can be found on the fact that SD-3
was collected among Census Bureau employees, while SD-1 was collec
ted among high-school students. Drawing sensible conclusions from
learning experiments requires that the result be independent of th
e choice of training set and test among the complete set of sample
s. Therefore it was necessary to build a new database by mixing NI
ST's datasets. \n\nThe MNIST training set is composed of 30,000 p
atterns from SD-3 and 30,000 patterns from SD-1. Our test set was
composed of 5,000 patterns from SD-3 and 5,000 patterns from SD-1.
The 60,000 pattern training set contained examples from approximat
ely 250 writers. We made sure that the sets of writers of the trai
ning set and test set were disjoint. SD-1 contains 58,527 digit im
ages written by 500 different writers. In contrast to SD-3, where
blocks of data from each writer appeared in sequence, the data in
SD-1 is scrambled. Writer identities for SD-1 is available and we
```

used this information to unscramble the writers. We then split SD-1 in two: characters written by the first 250 writers went into our new training set. The remaining 250 writers were placed in our test set. Thus we had two sets with nearly 30,000 examples each. The new training set was completed with enough examples from SD-3, starting at pattern # 0, to make a full set of 60,000 training patterns. Similarly, the new test set was completed with SD-3 examples starting at pattern # 35,000 to make a full set with 60,000 test patterns. Only a subset of 10,000 test images (5,000 from SD-1 and 5,000 from SD-3) is available on this site. The full 60,000 sample training set is available.\n\nDownloaded from openml.org."

Let's import the dataset, inputs and labels:

```
In [9]: 1 X, y = mnist['data'], mnist['target']
```

```
In [10]: 1 type(X)
```

```
Out[10]: numpy.ndarray
```

```
In [6]: 1 X.shape # 28x28 = 784
```

```
Out[6]: (70000, 784)
```

**If we print out the values for these columns, we see the edges - which are mainly 'white' (0)**

```
In [11]: 1 df = pd.DataFrame(X)
          2 df[:10]
```

Out[11]:

	0	1	2	3	4	5	6	7	8	9	...	774	775	776	777	778	779	780	781
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	0.0
1	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
2	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
3	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
4	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
5	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
6	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
7	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
8	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
9	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

10 rows × 784 columns

**However, choosing a pixel in the middle of the image - shows the range of colour (greyscale)**

```
In [12]: 1 df[:10][456]
```

Out[12]:

0	0.0
1	230.0
2	177.0
3	0.0
4	235.0
5	78.0
6	0.0
7	0.0
8	0.0
9	247.0

Name: 456, dtype: float64

**Y, on the other hand, is a one-dimensional array**

```
In [13]: 1 y.shape
```

```
Out[13]: (70000,)
```

```
In [14]: 1 y
```

```
Out[14]: array(['5', '0', '4', ..., '4', '5', '6'], dtype=object)
```

```
In [15]: 1 y[78]
```

```
Out[15]: '1'
```

```
In [16]: 1 y[26362]
```

```
Out[16]: '8'
```

X contains 70,000 images each of them contains 784 features, because each of them is a 28x28 picture. Each feature is a pixel intensity encoded in an 8-bit scale: from 0 (white) to 255 (black)

Let's display one or more images using matplotlib `imshow()`

```
In [17]: 1 digit = X[9]
2 digit_img = digit.reshape(28, 28)
3
4 plt.imshow(digit_img, cmap='binary')
5 plt.axis('off')
6 plt.show()
```



```
In [18]: 1 y[9]
```

```
Out[18]: '4'
```

The label is a string. We must convert it to a number for it to work on a Machine Learning algorithm.

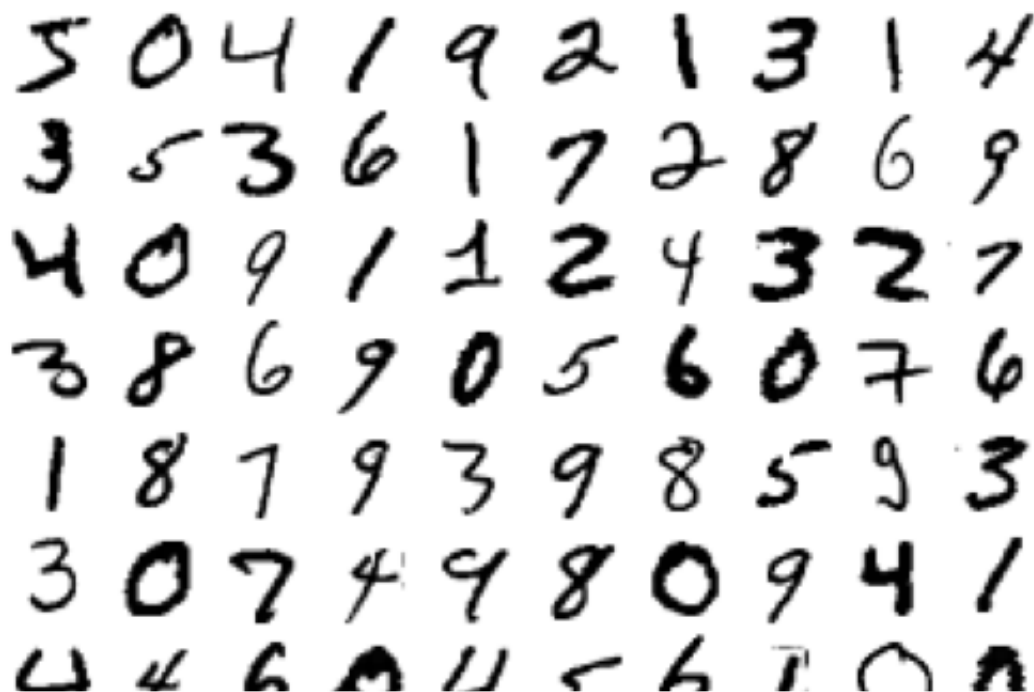
```
In [19]: 1 y = y.astype(np.uint8)
```

```
In [20]: 1 y[9]
```

```
Out[20]: 4
```

```
In [21]: 1 def show_digits(instances, images_per_row=10, **opts):
2         """
3         Utility function to display the MNIST digits on a grid
4         """
5         # the size of our images (28x28)
6         size = 28
7         images_per_row = min(len(instances), images_per_row)
8         # convert images from 1-D to 2-D arrays
9         images = [instance.reshape(size, size) for instance in inst
10        # compute how many rows you need in the grid
11        n_rows = (len(instances) - 1) // images_per_row + 1
12        row_images = []
13        # create empty "dummy" images to fill potential remaining s
14        n_empty = n_rows * images_per_row - len(instances)
15        images.append(np.zeros((size, size * n_empty)))
16        # concatenate all the images in a single grid image
17        for row in range(n_rows):
18            rimages = images[row * images_per_row : (row + 1) * ima
19            row_images.append(np.concatenate(rimages, axis=1))
20        image = np.concatenate(row_images, axis=0)
21        # plot the grid image
22        plt.imshow(image, cmap = mpl.cm.binary, **opts)
23        plt.axis("off")
```

```
In [22]: 1 plt.figure(figsize=(9,9))
          2 example_images = X[:100]
          3 show_digits(example_images, images_per_row=10)
          4 plt.show()
```



## 4.1 Split the dataset in training and test set

We'll set aside 10,000 samples for testing purposes. The data set is already shuffled for us so we can just take the last 10,000 samples for our test set.

```
In [23]: 1 # use slicing to create training and test set
          2 X_train, X_test, y_train, y_test = X[:60000], X[60000:], y[:60000], y[60000:]
```

```
In [24]: 1 X_train.shape, X_test.shape, y_train.shape, y_test.shape
```

```
Out[24]: ((60000, 784), (10000, 784), (60000,), (10000,))
```

## 4.1 Training a binary classifier

As a first goal, we will train a binary classifier, reducing our classes (classifications/categories) from 10 (values: 0-9) to 2 (first example: either 8 or not, and Ex1: even or odd).

Let's define two set of labels for the training and test set, named `y_train_8` and `y_test_8`. These must contain the value `True` whenever the original label is an 8, `False` otherwise

```
In [25]: 1 # Implemented using boolean masking
          2 y_train_8 = y_train == 8
          3 y_test_8 = y_test == 8
```

```
In [26]: 1 y_train_8[:20], y_train[:20]
```

```
Out[26]: (array([False, False, False, False, False, False, False, False, False, False,
                False, False, False, False, False, False, False, False, False,
                False, False]),
          array([5, 0, 4, 1, 9, 2, 1, 3, 1, 4, 3, 5, 3, 6, 1, 7, 2, 8, 6, 9],
                dtype=uint8))
```

**Exercise 1:** Let's suppose we want to implement a binary classifier to classify even vs odd digits. Define two set of labels for the training and test set, named `y_train_even` and `y_test_even`. These must contain the value `True` whenever the original label is a digit representing an even number, `False` if it's an odd number

```
In [27]: 1 # Write your solution here:
          2
          3 y_train_even = y_train % 2 == 0
          4 y_test_even = y_test % 2 == 0
          5
```

```
In [28]: 1 y_train_even = (y_train + 1) % 2
```

```
In [29]: 1 y_train_even[:20]
```

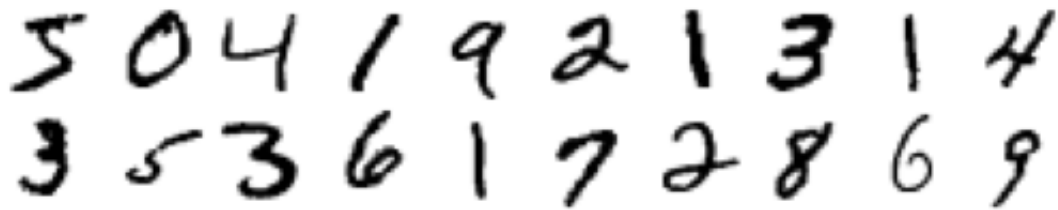
```
Out[29]: array([0, 1, 1, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 1, 1, 1, 0],
                dtype=uint8)
```

```
In [30]: 1 y[:20]
```

```
Out[30]: array([5, 0, 4, 1, 9, 2, 1, 3, 1, 4, 3, 5, 3, 6, 1, 7, 2, 8, 6, 9],
                dtype=uint8)
```



```
In [31]: 1 plt.figure(figsize=(9,9))
          2 example_images = X[:20]
          3 show_digits(example_images, images_per_row=10)
          4 plt.show()
```



### *Logistic regression classifier*

**(NOTE: This is for the Y\_TRAIN\_8 - NOT even and odd)**

We can train a logistic regression classifier by either using `sklearn.linear_model.SGDClassifier` with `loss` argument set as `log`.

```
In [32]: 1 from sklearn.linear_model import LogisticRegression
          2 lin_cl = LogisticRegression()
          3 lin_cl.fit(X_train, y_train_8)
```

/Users/nick/opt/anaconda3/lib/python3.9/site-packages/sklearn/linear\_model/\_logistic.py:814: ConvergenceWarning: lbfgs failed to converge (status=1):  
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.

Increase the number of iterations (`max_iter`) or scale the data as shown in:

<https://scikit-learn.org/stable/modules/preprocessing.html>  
(<https://scikit-learn.org/stable/modules/preprocessing.html>)

Please also refer to the documentation for alternative solver options:

[https://scikit-learn.org/stable/modules/linear\\_model.html#logistic-regression](https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression) ([https://scikit-learn.org/stable/modules/linear\\_model.html#logistic-regression](https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression))

```
n_iter_i = _check_optimize_result(
```

Out[32]: LogisticRegression()

```
In [33]: 1 from sklearn.linear_model import SGDClassifier
          2
          3 import time
          4 start_time = time.time()
          5
          6 sgd_cl = SGDClassifier(random_state=77, loss="log")
          7 sgd_cl.fit(X_train, y_train_8)
          8
          9 print("---- %s seconds ----" % (time.time() - start_time))

---- 21.016982793807983 seconds ----
```

```
In [34]: 1 sgd_cl.predict(X_train[:20])
```

```
Out [34]: array([False, False, False, False, False, False, False, False, False, False,
                False, False, True, False, False, False, False, False, True,
                False, False])
```

```
In [35]: 1 y_train[:20]
```

```
Out [35]: array([5, 0, 4, 1, 9, 2, 1, 3, 1, 4, 3, 5, 3, 6, 1, 7, 2, 8, 6, 9],
                dtype=uint8)
```

```
In [36]: 1 y[17]
```

```
Out [36]: 8
```

It has correctly predicted the "8" at index 17. However this belongs to the data it used during the training phase. We need a validation set to fairly evaluate the performance of our logistic regression classifier.

### 4.1.1 Performance Measures: measuring Accuracy Using Cross-Validation

We'll now use `cross_val_score()` to assess the accuracy of our Classifier `sgd_cl` on `(X_train, y_train)`, using 3-fold cross-validation.

```
In [37]: 1 from sklearn.model_selection import cross_val_score
          2
          3 import time
          4 start_time = time.time()
          5
          6 scores = cross_val_score(
          7     sgd_cl, X_train, y_train_8, cv=3, scoring="accuracy"
          8 )
          9 print("---- %s seconds ----" % (time.time() - start_time))
         10 #%time

---- 29.494776964187622 seconds ----
```

```
In [38]: 1 scores
```

```
Out[38]: array([0.93855, 0.93255, 0.84255])
```

```
In [39]: 1 np.mean(scores)
```

```
Out[39]: 0.90455
```

Our accuracy is 93 % in the first two runs and 84% in the third. The classifier looks very performant, but is it really the case?

**Exercise 2:** Use `cross_val_score()` to assess the accuracy of an SGD classifier implementing an online support vector machine (SVM), on `(X_train, y_train_8)`, using 5-fold cross-validation.

It is more or less accurate than the SGD classifier implementing logistic regression?

## Limited the dataset to a few thousand to reduce time to train

```

In [40]: 1 # Write your solution here
          2 from sklearn.svm import SVC
          3 svm_cl = SVC(gamma='auto')
          4 #start = time()
          5 #svm_cl.fit(X_train[:1000], y_train[:1000])
          6 #print('Duration: {} s'.format(time() - start))
          7 #svm_cl.predict(X_train[:10])
          8
          9
         10 import time
         11 start_time = time.time()
         12
         13 from sklearn.model_selection import cross_val_score
         14 start = time.time()
         15 scores = cross_val_score(
         16     svm_cl, # model we want to train
         17     X_train[:10000], # features
         18     y_train_8[:10000], # labels
         19     scoring='accuracy', # accuracy
         20     cv=5 #cross val checks (should be 5)
         21 )
         22 print("--- %s seconds ---" % (time.time() - start_time))
         23

```

--- 161.52036571502686 seconds ---

## More consistency between validation checks - doesn't drop off

```
In [41]: 1 scores
```

```
Out[41]: array([0.906 , 0.9055, 0.9055, 0.9055, 0.9055])
```

NOTE: Sometimes you need more control over cross-validation than what is offered out of the box with `cross_val_score()`. In the example in the cell below we are going to use the `StratifiedKFold` class to implement cross-validation

**NOTE: do not run this in the class, it takes way too much time.**

```
In [111]:
```

```

1  from sklearn.model_selection import StratifiedKFold
2  from sklearn.base import clone
3
4  skfolds = StratifiedKFold(n_splits=3, shuffle=True)
5
6  import time
7  start_time = time.time()
8
9  for train_index, val_index in skfolds.split(X_train, y_train_8)
10     # make a clone (copy) of our Stochastic Gradient Classifier
11     clone_sgd_cl = clone(sgd_cl)
12
13     # get training and validation set for current CV iteration
14     X_train_f = X_train[train_index]
15     X_val_f = X_train[val_index]
16     y_train_f = y_train[train_index]
17     y_val_f = y_train[val_index]
18
19     # train the SGD classifier
20     clone_sgd_cl.fit(X_train_f, y_train_f)
21     # make predictions on validation set
22     y_pred = clone_sgd_cl.predict(X_val_f)
23     # count number of correct predictions
24     n_correct = sum(y_pred == y_val_f)
25     # print out accuracy score
26     print(n_correct / len(y_val_f))
27
28 print("---- %s seconds ----" % (time.time() - start_time))

```

-----  
 ValueError

Traceback (most recent c

all last)

Input In [111], in <cell line: 4>()

```

1  from sklearn.model_selection import StratifiedKFold
2  from sklearn.base import clone
----> 4  skfolds = StratifiedKFold(n_splits=3, random_state=77)
6  import time
7  start_time = time.time()

```

File ~/opt/anaconda3/lib/python3.9/site-packages/sklearn/model\_selection/\_split.py:644, in StratifiedKFold.\_\_init\_\_(self, n\_splits, shuffle, random\_state)

```

643 def __init__(self, n_splits=5, *, shuffle=False, random_state=None):
--> 644     super().__init__(n_splits=n_splits, shuffle=shuffle, random_state=random_state)

```

File ~/opt/anaconda3/lib/python3.9/site-packages/sklearn/model\_selection/\_split.py:296, in \_BaseKFold.\_\_init\_\_(self, n\_splits, shuffle, random\_state)

```

293     raise TypeError("shuffle must be True or False; got {0}").format(shuffle))
295 if not shuffle and random_state is not None: # None is th

```

```

295 if not shuffle and random_state is not None: # None is th
e default
--> 296     raise ValueError(
297         "Setting a random_state has no effect since shuffl
e is "
298         "False. You should leave "
299         "random_state to its default (None), or set shuffl
e=True.",
300     )
302 self.n_splits = n_splits
303 self.shuffle = shuffle

```

**ValueError:** Setting a random\_state has no effect since shuffle is False. You should leave random\_state to its default (None), or set shuffle=True.

Let's go back to our SGD classifier `sgd_cl` trained as a logistic regressor. Our accuracy was 93 % in the first two runs and 84% in the third. The classifier looked very performant, but is it really the case?

Let's create a dummy classifier that never predicts that a digit is an "8". It will just always return False (i.e. 0) as a predicted label.

```

In [42]: 1 from sklearn.base import BaseEstimator
2 class Never8Classifier(BaseEstimator):
3
4     def fit(self, X, y=None):
5         pass
6
7     def predict(self, X):
8         return np.zeros((len(X), 1), dtype=bool)
9
10 never_8_clf = Never8Classifier()

```

Use `cross_val_score()` to assess the accuracy of our Classifier `never_8_clf` on `(X_train, y_train)`, using 3-fold cross-validation. Which accuracy do you expect?

In [43]:

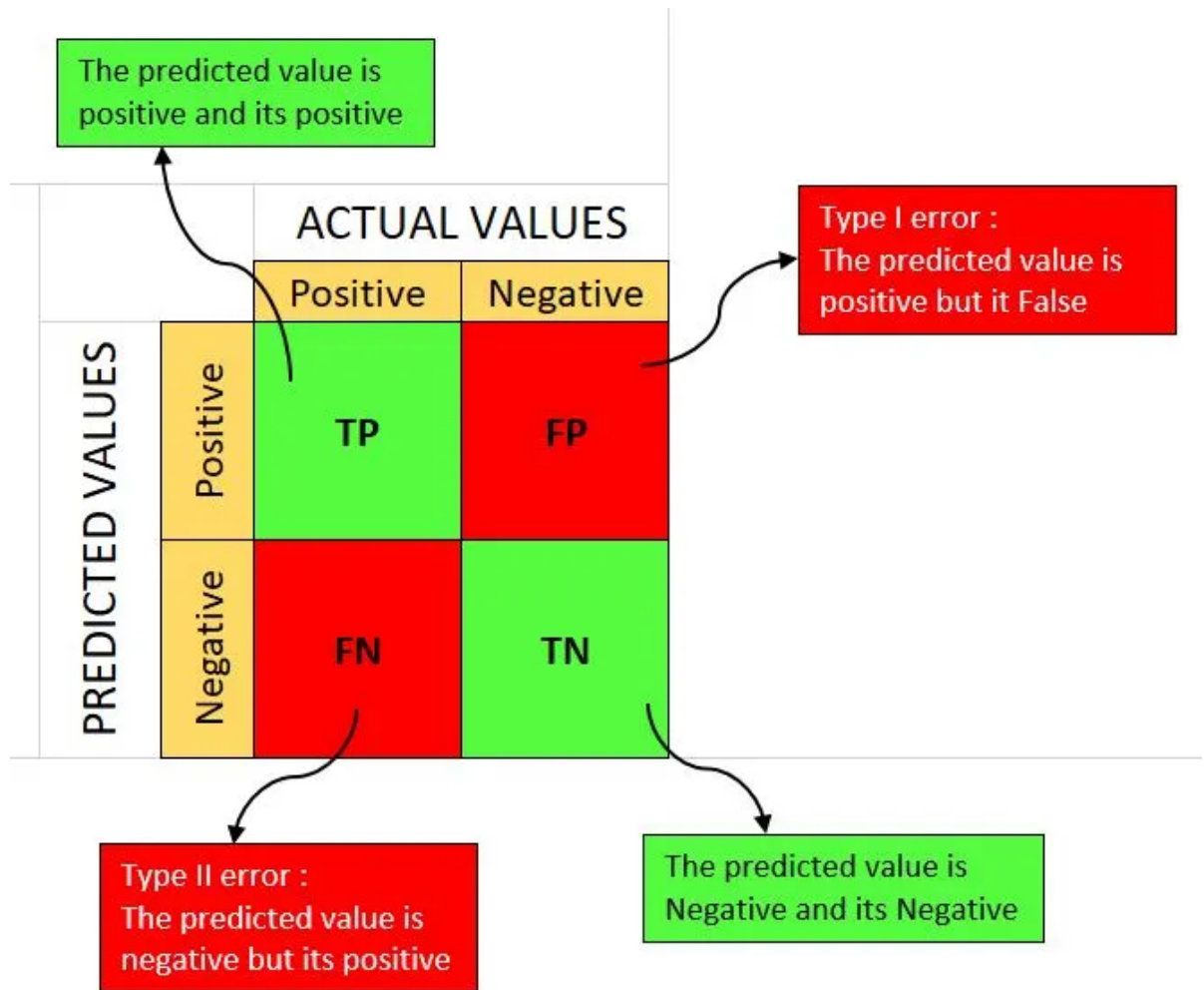
```
1 # Write your solution here
2
3 cross_val_score(
4     never_8_clf,
5     X_train,
6     y_train_8,
7     cv=3,
8     scoring="accuracy"
9 )
10
```

Out[43]: array([0.9039 , 0.9031 , 0.90045])

It has over 90% accuracy! This is simply because only about 10% of the images are 8s, so if you always guess that an image is not a 8, you will be right about 90% of the time.

Accuracy *per se* is not the preferred metrics when dealing with classifiers. This is even more true in this case, as we are dealing with a skewed dataset.

### 4.1.2 Performance Measures: Confusion Matrix



Left to right:

TP: "this is 8" and it is actually an 8

FP: "this is 8" when it is **not actually** an 8 => **TYPE 1 ERROR**

FN: "this is not an 8" when **it is actually** an 8 => **TYPE 2 ERROR**

TN: "this is not an 8" when it is not actually an 8

A more reliable way to measure the performance of a classifier is to look at the so-called *confusion matrix*. The aim is to quantify how many times members of a class C1 are misclassified as members of the class C2. To do that we will use the `cross_value_predict()` (rather the CV scores, it returns the predictions) function together with the `confusion_matrix()` metric.



```
In [44]: 1 from sklearn.model_selection import cross_val_predict
2 from sklearn.metrics import confusion_matrix
3
4 y_train_pred = cross_val_predict(
5     sgd_cl,
6     X_train,
7     y_train_8,
8     cv=3
9 )
10 y_train_pred
```

```
Out[44]: array([False, False, False, ..., False, False,  True])
```

```
In [45]: 1 confusion_matrix(y_train_8, y_train_pred)
```

```
Out[45]: array([[50274,  3875],
               [ 1852,  3999]])
```

The ideal perfect classifier would have true positives and true negatives only. In this case the confusion matrix would have zero values outside the main diagonal.

```
In [46]: 1 y_train_perfect_pred = y_train_8
2 confusion_matrix(
3     y_train_8,
4     y_train_perfect_pred
5 )
```

```
Out[46]: array([[54149,    0],
               [    0,  5851]])
```

#### 4.1.3 Precision, Recall and Harmonic mean (F1 score)

- Precision or Positive Predicted Value:

$$PPV = \frac{TP}{TP + FP}$$

- Recall or TPR or Sensitivity

$$TPR = \frac{TP}{TP + FN}$$

- Harmonic Mean of Precision and Recall

$$F1 = \frac{2}{\frac{1}{Precision} + \frac{1}{Recall}} = \frac{2TP}{TP + \frac{FN+FP}{2}}$$

```
In [47]: 1 from sklearn.metrics import precision_score, recall_score, f1_score
          2 ps = precision_score(y_train_8, y_train_pred)
          3 rs = recall_score(y_train_8, y_train_pred)
          4 f1s = f1_score(y_train_8, y_train_pred)
          5 ps, rs, f1s
```

```
Out [47]: (0.5078740157480315, 0.6834729106135703, 0.5827322404371583)
```

Now our classifier looks way worse than before! It has a 50% precision and 68% recall.

Notice that F1 tends to favour models that have similar precision and recall. But in some context you might prefer a higher precision, while in others a higher recall, depending on the task.

There is however a trade-off between precision and recall.

A classifier such as our `SGDClassifier` performs the classification task by computing a score based on a "decision function". If a score is greater than a given threshold value, the instance is labeled with the positive class, otherwise with the negative class (from the theory of Logistic regression, if you remember, an estimated probability of class "1" greater than 0.5 means that we assign the value to class "1"). Raising this threshold will reduce the number of FP, thus increasing the precision. However, it will also increase the number of FN thus reducing the recall score.

Let's try to manipulate the `SGDClassifier`'s threshold manually, using the classifier's `.decision_function()` method

```
In [48]: 1 y_scores = sgd_cl.decision_function(X_train[:20])
          2 y_scores
```

```
Out [48]: array([ -6004.06437255, -4627.92299992, -9386.52192365, -565.379327,
        -4448.21399887, -810.18121937, -1870.60495763, -398.465083,
        -1754.65161801, -1952.75837913, -2652.48419663, 998.62692833,
        -13189.1250276, -5654.82281211, -814.11597091, -2811.97829008,
        -3299.87374082, 2336.10363749, -2216.20131501, -2398.72241464])
```

```
In [49]: 1 threshold = 0
          2 y_pred_on_scores = y_scores > threshold
          3 y_pred_on_scores
```

```
Out[49]: array([False, False, False, False, False, False, False, False, False, False,
                False, False, True, False, False, False, False, False, True,
                False, False])
```

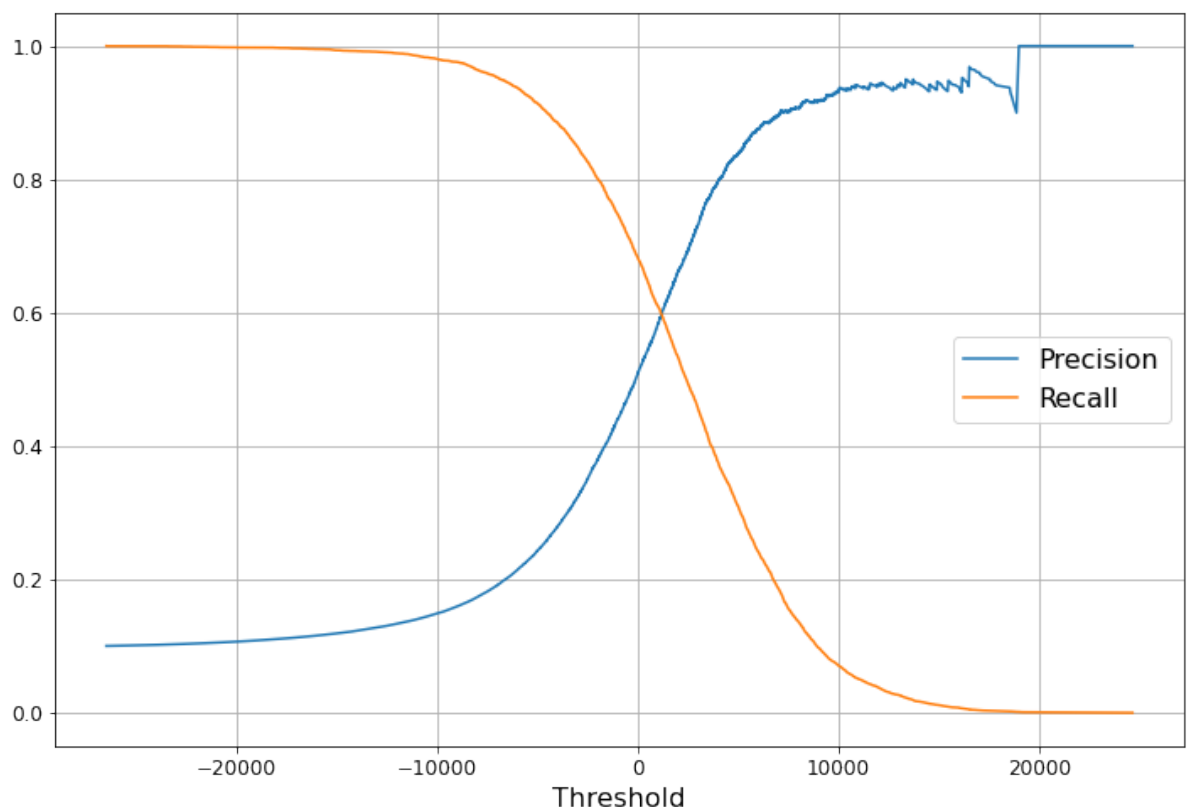
```
In [50]: 1 threshold = 500
          2 y_pred_on_scores = y_scores > threshold
          3 y_pred_on_scores
```

```
Out[50]: array([False, False, False, False, False, False, False, False, False, False,
                False, False, True, False, False, False, False, False, True,
                False, False])
```

Raising the threshold increases the number of FN, decreasing the recall. How can we then determine the right threshold value for our task? First let's use `cross_val_predict()` using the 'decision\_function' method on our entire training set, and then let's use the computed score together with the `precision_recall_curve()` to compute precision and recall for all the possible threshold values.

```
In [51]: 1 from sklearn.metrics import precision_recall_curve
          2 y_scores = cross_val_predict(
          3     sgd_cl,
          4     X_train,
          5     y_train_8,
          6     cv=3,
          7     method='decision_function'
          8 )
          9 precisions, recalls, thresholds = precision_recall_curve(
         10     y_train_8, y_scores
         11 )
```

```
In [60]: 1 import seaborn as sns
2
3 def plot_precision_and_recall_vs_threshold(
4     precisions, recalls, thresholds
5 ):
6     sns.lineplot(x=thresholds, y=precisions[:-1])
7
8     fig, ax = plt.subplots(figsize=(12, 8))
9     sns.lineplot(x=thresholds, y=precisions[:-1], ax=ax)
10    sns.lineplot(x=thresholds, y=recalls[:-1], ax=ax)
11    plt.legend(['Precision', 'Recall'], loc="center right", fontsize=12)
12    plt.xlabel("Threshold", fontsize=16)
13    plt.grid(True)
14    plt.show()
```



```
In [61]: 1 # Let's find the threshold for which we can achieve a 90% preci
2 threshold_90_prec = thresholds[np.argmax(precisions >= 0.90)]
3 threshold_90_prec
```

Out[61]: 7105.009210982933

We can now compute the predictions from the scores using this new threshold.

```
In [62]: 1 y_train_pred_90 = (y_scores >= threshold_90_prec)
```

```
In [63]: 1 precision_score(y_train_8, y_train_pred_90)
```

```
Out[63]: 0.9005059021922428
```

```
In [64]: 1 recall_score(y_train_8, y_train_pred_90)
```

```
Out[64]: 0.182532900358913
```

Now we have reached a 90% precision, at the expense of recall, which is now 18%!

## Week 5

### 4.1.3 Performance Measures: The ROC curve

Another tool that can be used to evaluate a classifier performance is the receiver-operating characteristic (ROC) curve. The ROC curves plots the true positive rate (TPR, i.e. recall) vs the false positive rate (FPR).

Specificity or TNR:

$$TNR = \frac{TN}{TN + FP}$$

True positive rate (TPR) or RECALL or SENSITIVITY:

$$TPR = \frac{TP}{TP + FN}$$

False positive rate (FPR):

$$FPR = 1 - TNR$$

The ROC curve plots sensitivity (TPR) against 1-specificity (FPR)

### Reminder:

TP: "this is 8" and it is actually an 8

FP: "this is 8" when it is **not actually** an 8 => **TYPE 1 ERROR**

FN: "this is not an 8" when **it is actually** an 8 => **TYPE 2 ERROR**

TN: "this is not an 8" when it is not actually an 8

```
In [67]: 1 from sklearn.metrics import roc_curve
          2
          3 import time
          4 start_time = time.time()
          5
          6 y_scores_sgd = cross_val_predict(
          7     sgd_cl, # our SGD classifier trained to fit a Logistic Regr
          8     X_train,
          9     y_train_8,
         10     cv=3,
         11     method='decision_function'
         12 )
         13
         14 fpr, tpr, thresholds = roc_curve(y_train_8, y_scores_sgd)
         15
         16 y_scores_sgd
         17
         18 #print("--- %s seconds ---" % (time.time() - start_time))
         19
```

```
Out[67]: array([-23493.52509479, -14499.86059748, -24419.59924926, ...,
               -1196.86161131, -23098.46876637,  5773.08810945])
```

```
In [80]: 1 tpr.mean()
```

```
Out[80]: 0.6652825240561266
```

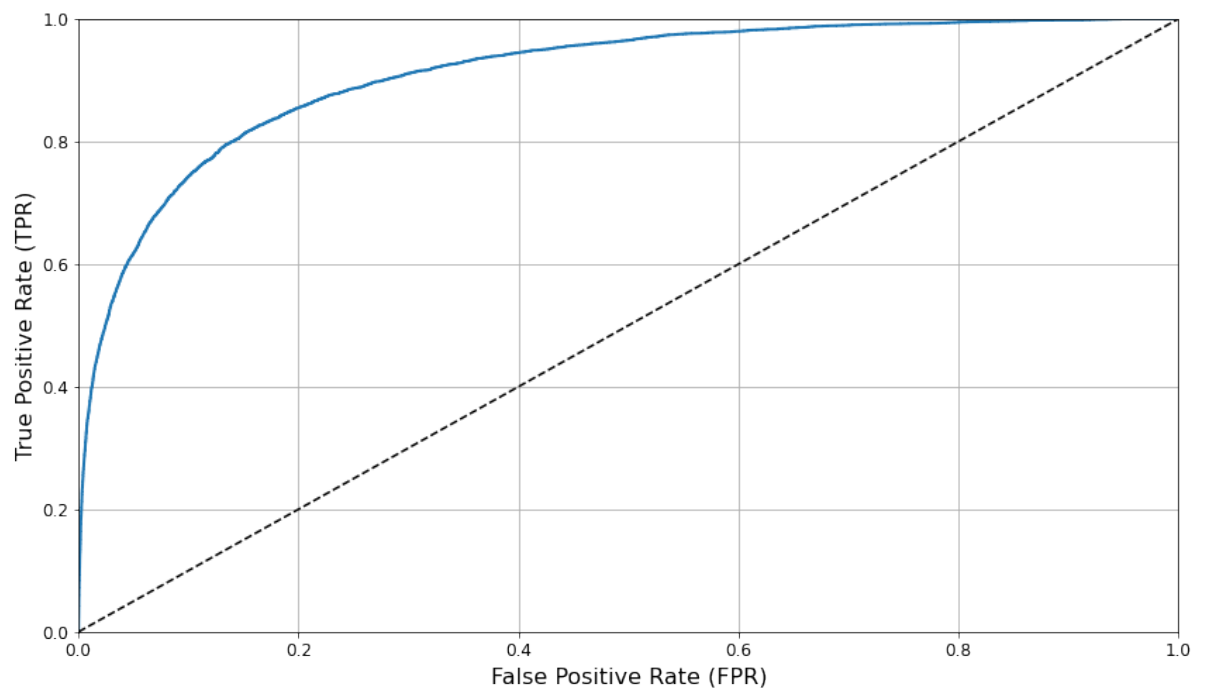
```
In [81]: 1 fpr.mean()
```

```
Out[81]: 0.14865795120755948
```

```
In [77]: 1 fpr.shape
```

```
Out[77]: (6466,)
```

```
In [53]: 1 plt.figure(figsize=(14, 8))
2 def plot_roc_curve(fpr, tpr, label=None):
3     plt.plot(fpr, tpr, linewidth=2, label=label)
4     plt.plot([0, 1], [0, 1], 'k--')
5     plt.axis([0, 1, 0, 1])
6     plt.xlabel('False Positive Rate (FPR)', fontsize=16)
7     plt.ylabel('True Positive Rate (TPR)', fontsize=16)
8     plt.grid(True)
9 plot_roc_curve(fpr, tpr)
10 plt.show()
11
```



You can measure the *area under the curve* (AUC) if you want to compare the performance of different classifiers.

```
In [54]: 1 from sklearn.metrics import roc_auc_score
2 roc_auc_score(y_train_8, y_scores_sgd)
```

Out [54]: 0.9102127159789787

As a general rule, prefer the Precision/Recall curve if the positive class is uncommon or if you worry more about the false positives rather than the false negatives. In the other scenarios, prefer the ROC curve.

## ROC AUC for KNNs and Random Forests

Let's try two different classifiers: a K-Nearest Neighbours classifier and a Random Forest classifier.

The K-Nearest Neighbours algorithm checks the K closest (i.e. most similar instances) in the training set and assigns as predicted class for the new instance the most represented class in the neighbourhood.

The Random Forest algorithm is an ensemble method which trains a number of decision tree classifiers on various sub-samples of the training set and uses averaging techniques to improve the predictive accuracy and control over-fitting.

We will see more on Decision Trees and Ensemble methods in the next weeks.

**NOTE:** K-Nearest Neighbour and Random Forest classifiers do not have a `decision_function()` method that returns the predicted scores for each instance in cross validation. They do have, however, a `predict_proba()` method that returns an array containing a row per instance and a column per class. This array contains the predicted probability that each instance belongs to a class. This can be used to draw ROC curves in lieu of `decision_function()`. Scikit-learn classifiers usually implement either one or the other method so you need to check their API to find out the one you need to use.

```
In [55]: 1 from sklearn.neighbors import KNeighborsClassifier
          2
          3 import time
          4 start_time = time.time()
          5
          6 kn_cl = KNeighborsClassifier(n_neighbors=9)
          7 y_probs_kn = cross_val_predict(
          8     kn_cl,
          9     X_train,
         10     y_train_8,
         11     cv=3,
         12     method='predict_proba'
         13 )
         14 print("--- %s seconds ---" % (time.time() - start_time))

--- 77.38686609268188 seconds ---
```



```

In [56]: 1 from sklearn.ensemble import RandomForestClassifier
          2
          3 import time
          4 start_time = time.time()
          5
          6 forest_cl = RandomForestClassifier(
          7     n_estimators=100, # a "forest" of 100 decision trees
          8     random_state=77
          9 )
         10 y_probs_forest = cross_val_predict(
         11     forest_cl,
         12     X_train,
         13     y_train_8,
         14     cv=3,
         15     method="predict_proba"
         16 )
         17 print("---- %s seconds ----" % (time.time() - start_time))

---- 34.85027623176575 seconds ----

```

```

In [57]: 1 y_probs_forest

```

```

Out[57]: array([[0.98, 0.02],
               [1.  , 0.  ],
               [0.98, 0.02],
               ...,
               [0.98, 0.02],
               [1.  , 0.  ],
               [0.31, 0.69]])

```

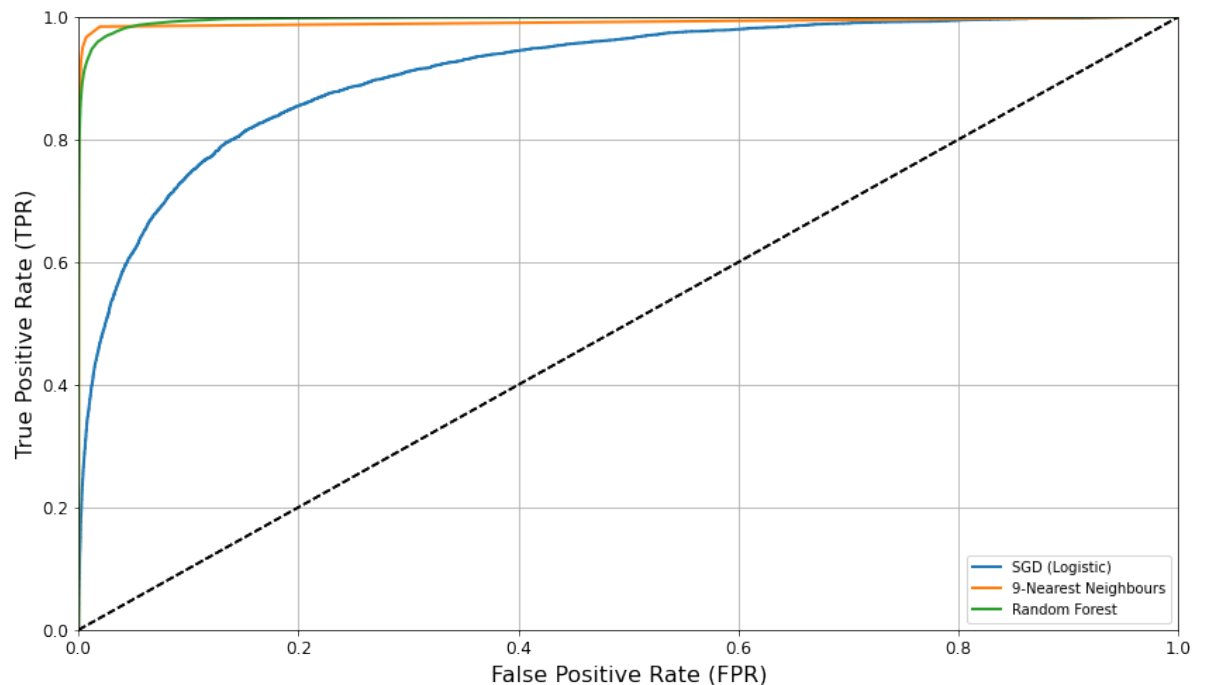
```

In [58]: 1 # Nearest neighbours scores
          2 y_scores_kn = y_probs_kn[:, 1] # score = proba of positive c
          3 fpr_kn, tpr_kn, thresholds_kn = roc_curve(y_train_8, y_scores_k
          4 # Random Forests scores
          5 y_scores_forest = y_probs_forest[:, 1] # score = proba of posit
          6 fpr_forest, tpr_forest, thresholds_forest = roc_curve(y_train_8

```

Use the scores computed above for `sgd_cl`, `kn_cl`, `forest_cl` to plot out the three ROC curves on the same plot. You can use, if you wish the `plot_roc_curve()` function defined above. Afterward compute the area under the curve for `mnb_cl` and `forest_cl`. Which is the best and the worst classifier?

```
In [59]: 1 plt.figure(figsize=(14, 8))
2 plot_roc_curve(fpr, tpr, "SGD (Logistic)")
3 plot_roc_curve(fpr_kn, tpr_kn, "9-Nearest Neighbours")
4 plot_roc_curve(fpr_forest, tpr_forest, "Random Forest")
5 plt.legend(loc="lower right")
6 plt.show()
```



```
In [61]: 1 print("SGD (Logistic):", roc_auc_score(y_train_8, y_scores_sgd))
2 print("9 Nearest Neighbours:", roc_auc_score(y_train_8, y_score
3 print("Random Forest:", roc_auc_score(y_train_8, y_scores_forest))
```

```
SGD (Logistic): 0.9102127159789787
9 Nearest Neighbours: 0.9913318627817932
Random Forest: 0.9965127855639053
```

```
In [83]: 1 roc_auc_score(y_train_8, y_scores_forest)
```

```
Out[83]: 0.9965127855639053
```

**Exercise 4:** Try and train a few more classifiers and plot their ROC curves. Which of the them has the better area under the curve? Which one has the "steepest" ROC curve?

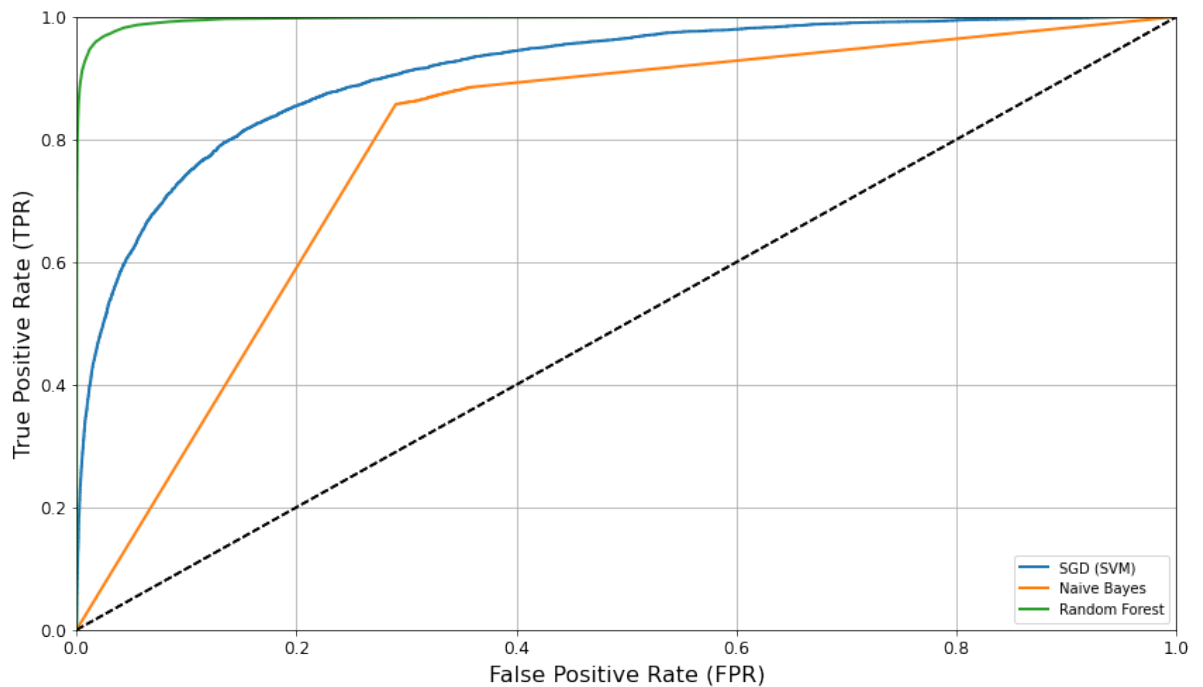
Hint 1: you can also train similar models but exploring the hyperparameter space.

Hint2: You could use the yellowbrick library (<https://www.scikit-yb.org/en/latest/index.html>) to get better visualization of the ROC curves. See <https://www.scikit-yb.org/en/latest/api/classifier/rocauc.html>. You can install yellowbrick using Anaconda Navigator (make sure to use the "conda-forge" channel) or by running `!conda install -c conda-forge -y yellowbrick` in a notebook cell.

```
In [62]: 1 # Write your solution here
          2
          3 # naive bayes
          4
          5 from sklearn.naive_bayes import MultinomialNB
          6 mnb_cl = MultinomialNB()
          7 y_probs_mnb = cross_val_predict(
          8     mnb_cl,
          9     X_train,
         10     y_train_8,
         11     cv=3,
         12     method='predict_proba'
         13 )
         14
```

```
In [63]: 1 # Naive Bayes scores
          2 y_scores_mnb = y_probs_mnb[:, 1] # score = proba of positive
          3 fpr_mnb, tpr_mnb, thresholds_mnb = roc_curve(y_train_8, y_score
          4 # Random Forests scores
          5 y_scores_forest = y_probs_forest[:, 1] # score = proba of posit
          6 fpr_forest, tpr_forest, thresholds_forest = roc_curve(y_train_8
          7 # plt.plot(fpr, tpr, "b:", label="SGD")
```

```
In [87]: 1 # Write your solution here:
2 plt.figure(figsize=(14, 8))
3 plot_roc_curve(fpr, tpr, "SGD (SVM)")
4 plot_roc_curve(fpr_mnb, tpr_mnb, "Naive Bayes")
5 plot_roc_curve(fpr_forest, tpr_forest, "Random Forest")
6 plt.legend(loc="lower right")
7 plt.show()
```



```
In [92]: 1 print("SGD (SVM):", roc_auc_score(y_train_8, y_scores_sgd))
2 print("Niave bayes:", roc_auc_score(y_train_8, y_scores_mnb))
3 print("Random Forest:", roc_auc_score(y_train_8, y_scores_forest))
```

```
SGD (SVM): 0.9102127159789787
Niave bayes: 0.7885469295383991
Random Forest: 0.9965127855639053
```

## 4.2 Multiclass Classification

When you need to distinguish more than two classes, you have a multiclass classification problem. ML learning algorithms, solve the multiclass problem mainly in three ways:

- they support multi-class classification natively (e.g. SGD Classifier, K-Nearest Neighbour, Naive Bayes, Random Forests)
- they are binary classifiers in one-vs-the-rest strategy (OvR) ( $N$  classifiers are needed for  $N$  classes)
- they are binary classifiers in one-vs-one-strategy (OvO) ( $N \times (N - 1)/2$  classifiers are needed for  $N$  classes). Each classifier has to be trained only on the subset of the training set for the two classes that it must distinguish.

(Offline) Support Vector Machines are generally trained using the OvO approach. Let's see an example (we will only use a subset of the dataset)

```
In [65]: 1 from time import time
          2 from sklearn.svm import SVC
          3 svm_cl = SVC(gamma='auto')
          4 start = time()
          5 svm_cl.fit(X_train[:1000], y_train[:1000])
          6 print('Duration: {} s'.format(time() - start))
          7 svm_cl.predict(X_train[:10])
```

Duration: 0.3419520854949951 s

```
Out [65]: array([5, 0, 4, 1, 9, 2, 1, 3, 1, 4], dtype=uint8)
```

**NOTE: do not run this in the class, it takes several hours to run, depending on your hardware specs. (on a MacBook Pro 2019 took 6 hours 55 min)**

```
In [ ]: 1 # Try and train it on the whole training set
          2 svm_cl = SVC(gamma='auto')
          3 start = time()
          4 svm_cl.fit(X_train, y_train)
          5 print('Duration: {} s'.format(time() - start))
          6 svm_cl.predict(X_train[:10])
```

```
In [66]: 1 scores = svm_cl.decision_function(X_train[:10])
          2 scores
```

```
Out[66]: array([[ 2.81585438,  7.09167958,  3.82972099,  0.79365551,  5.888
5703 ,
                9.29718395,  1.79862509,  8.10392157, -0.228207 ,  4.837
53243],
               [ 9.29838234,  7.09167958,  3.82972099,  1.79572006,  5.888
5703 ,
                0.7913911 ,  2.80027801,  8.10392157, -0.22656281,  4.837
53243],
               [ 3.82111996,  7.09167958,  4.83444983,  1.79943469,  9.299
32174,
                0.79485736,  2.80437474,  8.10392157, -0.22417259,  5.841
82891],
               [ 3.82760119,  9.29995923,  4.84239684,  1.80408497,  6.903
66992,
                0.79908447,  2.80938404,  8.10392157, -0.22130643,  5.850
60746],
               [ 3.81790353,  7.09167958,  4.83058232,  1.79710089,  5.888
5703 ,
                0.79272692,  2.80181215,  8.10392157, -0.22565337,  9.298
85985],
               [ 3.81723639,  7.09167958,  9.29871755,  1.79670336,  5.888
5703 ,
                0.79230931,  2.80133228,  8.10392157, -0.22592444,  4.837
53243],
               [ 3.82760119,  9.29995923,  4.84239684,  1.80408497,  6.903
66992,
                0.79908447,  2.80938404,  8.10392157, -0.22130643,  5.850
60746],
               [ 2.81585438,  7.09167958,  3.82972099,  9.29748313,  5.888
5703 ,
                0.78955012,  1.79862509,  8.10392157, -0.22784964,  4.837
53243],
               [ 3.82760119,  9.29995923,  4.84239684,  1.80408497,  6.903
66992,
                0.79908447,  2.80938404,  8.10392157, -0.22130643,  5.850
60746],
               [ 3.82111996,  7.09167958,  4.83444983,  1.79943469,  9.299
32174,
                0.79485736,  2.80437474,  8.10392157, -0.22417259,  5.841
82891]])
```

```
In [68]: 1 np.argmax(scores, axis=1)
```

```
Out[68]: array([5, 0, 4, 1, 9, 2, 1, 3, 1, 4])
```

```
In [69]: 1 svm_cl.classes_
```

```
Out[69]: array([0, 1, 2, 3, 4, 5, 6, 7, 8, 9], dtype=uint8)
```

```
In [71]: 1 import time
2 start_time = time.time()
3
4 forest_cl = RandomForestClassifier(n_estimators=100, random_state=77, cv=3)
5 scores = cross_val_score(forest_cl, X_train, y_train, cv=3, scoring='accuracy')
6
7 #print('Execution time {} s'.format(time() - start_time))
8 print("---- %s seconds ----" % (time.time() - start_time))

---- 46.61604690551758 seconds ----
```

```
In [72]: 1 import time
2 start_time = time.time()
3
4 sgd_cl = SGDClassifier(random_state=77, tol=1e-3, max_iter=2000, cv=3)
5 scores = cross_val_score(sgd_cl, X_train, y_train, cv=3, scoring='accuracy')
6
7 print("---- %s seconds ----" % (time.time() - start_time))

---- 54.88285803794861 seconds ----
```

Our `sgd_cl` should be initialized with a higher `max_iter` (eg 2000), otherwise the algorithm might not converge. This is not done in the live demo because the execution time becomes considerable.

```
In [73]: 1 from sklearn.preprocessing import StandardScaler
2 scaler = StandardScaler()
3 X_train_scaled = scaler.fit_transform(X_train.astype(np.float64))
```

**NOTE: do not run this in the class, it takes way too much time.**

```
In [ ]: 1 start = time()
2 scores = cross_val_score(sgd_cl, X_train_scaled, y_train, cv=3, scoring='accuracy')
3 print('Execution time {} s'.format(time() - start))
4 scores
```

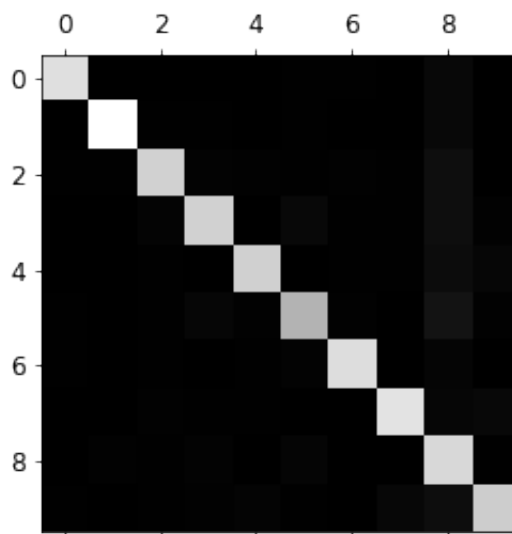
## 4.2.1 Model Evaluation

In a proper machine learning problem solving approach you would pre-process the input data (as shown in Class 2) and then try a few ML algorithms, also exploring the parameter space using `GridSearchCV` or `RandomizedSearchCV`. Hypothesising that you have done so, and that you have a good model, you will then want to evaluate its performance. Let's look at the confusion matrix, as a first step.

```
In [74]: 1 # NOTE: this cell should take about 3 minutes to run on a 2019 .
2 # depending on hardware
3 sgd_cl = SGDClassifier(random_state=77, tol=1e-3, max_iter=2000
4
5 import time
6 start_time = time.time()
7
8 y_train_pred = cross_val_predict(sgd_cl, X_train_scaled, y_train
9 c_mat = confusion_matrix(y_train, y_train_pred)
10 print("--- %s seconds ---" % (time.time() - start_time))

--- 151.326397895813 seconds ---
```

```
In [75]: 1 plt.matshow(c_mat, cmap=plt.cm.gray)
2 plt.show()
```



We can evaluate the performance of the predictor (precision, recall and F1 score) for each class, using the summary function `classification_report()`

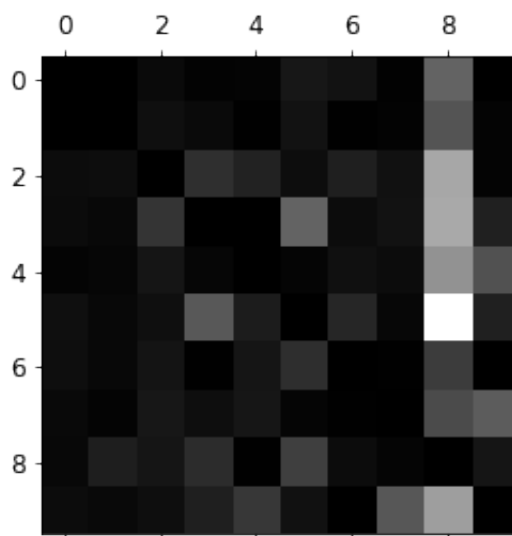


```
In [76]: 1 from sklearn.metrics import classification_report
2 report = classification_report(y_train, y_train_pred)
3 print(report)
```

	precision	recall	f1-score	support
0	0.97	0.94	0.96	5923
1	0.97	0.95	0.96	6742
2	0.93	0.88	0.90	5958
3	0.91	0.86	0.88	6131
4	0.94	0.89	0.91	5842
5	0.88	0.83	0.85	5421
6	0.95	0.94	0.95	5918
7	0.95	0.91	0.93	6265
8	0.68	0.93	0.78	5851
9	0.90	0.86	0.88	5949
accuracy			0.90	60000
macro avg	0.91	0.90	0.90	60000
weighted avg	0.91	0.90	0.90	60000

Let's plot the errors. Firstly, we divide each value by the number of images in the corresponding class so that we compare error rates rather than absolute numbers of errors. Then we fill the values along the diagonal with zeros to keep only the classification errors.

```
In [77]: 1 row_sums = c_mat.sum(axis=1, keepdims=True)
2 norm_c_mat = c_mat / row_sums
3 np.fill_diagonal(norm_c_mat, 0)
4 plt.matshow(norm_c_mat, cmap=plt.cm.gray)
5 plt.show()
```



Now, we can better see which type of errors our classifier makes. Which are these?

## 4.3 Multilabel and Multioutput Classification

```
In [78]: 1 from sklearn.neighbors import KNeighborsClassifier
          2
          3 y_train_prime = np.in1d(y_train, [2, 3, 5, 7])
          4 y_train_odd = y_train % 2 == 1
          5 y_multilabel = np.c_[y_train_prime, y_train_odd]
          6
          7 knn_clf = KNeighborsClassifier()
          8 knn_clf.fit(X_train, y_multilabel)
```

Out[78]: KNeighborsClassifier()

```
In [79]: 1 knn_clf.predict(X_train[:10])
```

Out[79]: array([[ True, True],  
[False, False],  
[False, False],  
[False, True],  
[False, True],  
[ True, False],  
[False, True],  
[ True, True],  
[False, True],  
[False, False]])

```
In [80]: 1 y_train[:10], y_multilabel[:10, 0], y_multilabel[:10, 1]
```

Out[80]: (array([5, 0, 4, 1, 9, 2, 1, 3, 1, 4], dtype=uint8),  
array([ True, False, False, False, False, True, False, True, False,  
False]),  
array([ True, False, False, True, True, False, True, True, True,  
False]))

**NOTE: do not run this in the class, it takes way too much time.**

```
In [81]: 1 y_train_knn_pred = cross_val_predict(knn_clf, X_train[:10], y_m
          2 f1_score(y_multilabel[:10, 1], y_train_knn_pred, average="macro
```

Out[81]: 0.37499999999999994

```
In [ ]: 1
```