# Measuring model performance

## Autumn 2020

```
###################################################
##                                               ##
##   Machine Learning Fundamentals with Python   ##
##                                               ##
###################################################
```

§1 Supervised Learning with scikit-learn

§1.1 Classification

§1.1.4 Measuring model performance

**1. What is the accuracy?**

- In classification, accuracy is a commonly used metric.

- Accuracy = *fraction of correct predictions*

**2. How to measuring model performance?**

- Split data into training and test set.

- Fit/train the classifier on the training set.

- Make predictions on the test set.

- Compare predictions with the known labels.

**3. Code of train/test split:**

```python
[1]: from sklearn import datasets
     from sklearn.neighbors import KNeighborsClassifier
     from sklearn.model_selection import train_test_split

     iris = datasets.load_iris()
     X = iris.data
     y = iris.target
     X_train, X_test, y_train, y_test = train_test_split(X,
                                                         y,
                                                         test_size=0.3,
                                                         random_state=21,
                                                         stratify=y)
     knn = KNeighborsClassifier(n_neighbors=8)
```

```
knn.fit(X_train, y_train)
y_pred = knn.predict(X_test)
print("Test set predictions:\n {}".format(y_pred))
```

```
Test set predictions:
 [2 1 2 2 1 0 1 0 0 1 0 2 0 2 2 0 0 0 1 0 2 2 2 0 1 1 1 0 0 1 2 2 0 0 1 2 2
 1 1 2 1 1 0 2 1]
```

```
[2]: knn.score(X_test, y_test)
```

```
[2]: 0.9555555555555556
```

### 4. What is model complexity?

- Larger k:
  - smoother decision boundary
  - less complex model
- Smaller k:
  - more complex model
  - can lead to overfitting

### 5. Practice exercise for measuring model performance:

▶ **Package pre-loading:**

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

▶ **The digits recognition dataset practice:**

```
[4]: # Import necessary modules
from sklearn import datasets
import matplotlib.pyplot as plt

# Load the digits dataset: digits
digits = datasets.load_digits()

# Print the keys and DESCR of the dataset
print(digits.DESCR)

# Print the shape of the images and data keys
print(digits.images.shape)
print(digits.data.shape)

# Display digit 1010
plt.imshow(digits.images[1010], cmap=plt.cm.gray_r, interpolation='nearest')
plt.show()
```

```
.. _digits_dataset:
```

Optical recognition of handwritten digits dataset
--------------------------------------------------

**Data Set Characteristics:**

    :Number of Instances: 5620
    :Number of Attributes: 64
    :Attribute Information: 8x8 image of integer pixels in the range 0..16.
    :Missing Attribute Values: None
    :Creator: E. Alpaydin (alpaydin '@' boun.edu.tr)
    :Date: July; 1998

This is a copy of the test set of the UCI ML hand-written digits datasets
https://archive.ics.uci.edu/ml/datasets/Optical+Recognition+of+Handwritten+Digits

The data set contains images of hand-written digits: 10 classes where
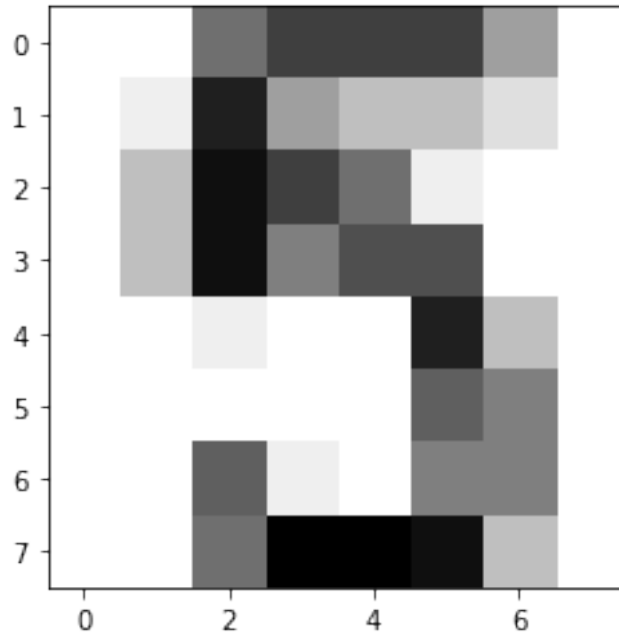each class refers to a digit.

Preprocessing programs made available by NIST were used to extract
normalized bitmaps of handwritten digits from a preprinted form. From a
total of 43 people, 30 contributed to the training set and different 13
to the test set. 32x32 bitmaps are divided into nonoverlapping blocks of
4x4 and the number of on pixels are counted in each block. This generates
an input matrix of 8x8 where each element is an integer in the range
0..16. This reduces dimensionality and gives invariance to small
distortions.

For info on NIST preprocessing routines, see M. D. Garris, J. L. Blue, G.
T. Candela, D. L. Dimmick, J. Geist, P. J. Grother, S. A. Janet, and C.
L. Wilson, NIST Form-Based Handprint Recognition System, NISTIR 5469,
1994.

.. topic:: References

  - C. Kaynak (1995) Methods of Combining Multiple Classifiers and Their
    Applications to Handwritten Digit Recognition, MSc Thesis, Institute of
    Graduate Studies in Science and Engineering, Bogazici University.
  - E. Alpaydin, C. Kaynak (1998) Cascading Classifiers, Kybernetika.
  - Ken Tang and Ponnuthurai N. Suganthan and Xi Yao and A. Kai Qin.
    Linear dimensionalityreduction using relevance weighted LDA. School of
    Electrical and Electronic Engineering Nanyang Technological University.
    2005.
  - Claudio Gentile. A New Approximate Maximal Margin Classification
    Algorithm. NIPS. 2000.
(1797, 8, 8)

```
(1797, 64)
```



▶ **Train/test split and fit/predict/accuracy practice:**

```
[5]: # Import necessary modules
     from sklearn.neighbors import KNeighborsClassifier
     from sklearn.model_selection import train_test_split

     # Create feature and target arrays
     X = digits.data
     y = digits.target

     # Split into training and test set
     X_train, X_test, y_train, y_test = train_test_split(X,
                                                         y,
                                                         test_size=0.2,
                                                         random_state=42,
                                                         stratify=y)

     # Create a k-NN classifier with 7 neighbors: knn
     knn = KNeighborsClassifier(n_neighbors=7)

     # Fit the classifier to the training data
     knn.fit(X_train, y_train)

     # Print the accuracy
```

```
print(knn.score(X_test, y_test))
```

0.9833333333333333

▶ **Overfitting and underfitting practice:**

```python
[6]: # Setup arrays to store train and test accuracies
neighbors = np.arange(1, 9)
train_accuracy = np.empty(len(neighbors))
test_accuracy = np.empty(len(neighbors))

# Loop over different values of k
for i, k in enumerate(neighbors):
    # Setup a k-NN Classifier with k neighbors: knn
    knn = KNeighborsClassifier(n_neighbors=k)

    # Fit the classifier to the training data
    knn.fit(X_train, y_train)

    #Compute accuracy on the training set
    train_accuracy[i] = knn.score(X_train, y_train)

    #Compute accuracy on the testing set
    test_accuracy[i] = knn.score(X_test, y_test)

# Generate plot
plt.title('k-NN: Varying Number of Neighbors')
plt.plot(neighbors, test_accuracy, label='Testing Accuracy')
plt.plot(neighbors, train_accuracy, label='Training Accuracy')
plt.legend()
plt.xlabel('Number of Neighbors')
plt.ylabel('Accuracy')
plt.show()
```

k-NN: Varying Number of Neighbors