

Laboratory 1

Data analysis and visualization with machine learning

An introduction to machine learning with scikit-learn Scikit-learn is a Python module integrating classic machine learning algorithms in the tightly-knit world of scientific Python packages (NumPy, SciPy, matplotlib). Analyze the following tutorial by executing the given examples <http://scikit-learn.org/stable/tutorial/basics/tutorial.html>

1 - In Python, usually, the functions are included in libraries which could be imported. For example to import the scikit-learn library :

```
In [1985]: from sklearn import *
```

2 - Import the libraries numpy (scientific computation) and matplotlib.pyplot (visualization).

```
In [1986]: import numpy as np
import matplotlib.pyplot as np
```

3-Load the Iris dataset using:

```
In [1987]: iris = datasets.load_iris()
```

The variable iris is an object which contains the dataset matrix iris.data, a vector containing the labels/classes (target), the name of variables (feature_names) and the name of classes (target_names).

4- Print the number of data, names of variables and the name of classes using print.

```
In [1988]: print(iris.data)

[[5.1 3.5 1.4 0.2]
 [4.9 3. 1.4 0.2]
 [4.7 3.2 1.3 0.2]
 [4.6 3.1 1.3 0.2]
 [5. 3.6 1.4 0.2]
 [5.4 3.9 1.7 0.4]
 [4.6 3.4 1.4 0.3]
 [5. 3.4 1.5 0.2]
 [4.4 2.9 1.4 0.2]
 [4.9 3.1 1.5 0.1]
 [5.4 3.7 1.5 0.2]
 [4.8 3.4 1.6 0.2]
 [4.8 3. 1.4 0.1]
 [4.3 3. 1.1 0.1]
 [5.8 4.1 1.2 0.2]
 [5.7 4.4 1.5 0.4]
 [5.4 3.9 1.3 0.4]
 [5.1 3.5 1.4 0.3]
 [5.7 3.8 1.7 0.3]
 [5.1 3.8 1.5 0.3]
 [5.4 3.1 1.7 0.2]
 [5.1 3.7 1.5 0.4]
 [4.6 3.6 1. 0.2]
 [5.1 3.7 1.7 0.5]
 [5.8 3.4 1.9 0.2]
 [9. 3. 1.6 0.2]
 [4.9 3.4 1.6 0.2]
 [5.2 3.5 1.5 0.2]
 [5.2 3.4 1.4 0.2]
 [4.7 3.2 1.6 0.2]
 [4.8 3.1 1.4 0.2]
 [4.9 3.4 1.5 0.4]
 [5.2 3.5 1.5 0.2]
 [4.9 3.1 1.5 0.2]
 [5.5 4.2 1.4 0.1]
 [5.5 3.5 1.3 0.3]
 [4.9 3.1 1.5 0.2]
 [4.9 3.6 1.4 0.1]
 [4.4 3.1 1.3 0.2]
 [5.1 3.4 1.5 0.3]
 [9. 3.5 1.3 0.3]
 [4.9 3.1 1.3 0.3]
 [9. 3.5 1.6 0.6]
 [5.5 3.8 1.9 0.4]
 [5.2 3.5 1.5 0.2]
 [4.8 3. 1.4 0.3]
 [5.1 3.8 1.6 0.2]
 [4.8 3.2 1.4 0.2]
 [5.3 3.7 1.5 0.2]
 [9. 3.3 1.4 0.2]
 [6. 3. 1.4 0.2]
 [6.2 3.2 1.7 0.4]
 [6.4 3.2 1.4 0.1]
 [6.3 3.7 1.5 0.1]
 [6.9 3.1 4.9 1.5]
 [5.5 3.2 4. 1.3]
 [6.5 2.8 4.6 1.3]
 [7.9 2.8 4.5 1.3]
 [6.3 3.1 4.7 1.6]
 [4.9 2.8 4.6 1.3]
 [6.6 2.9 4.6 1.3]
 [5.2 2.7 3.9 1.4]
 [6.7 2.9 4.5 1.3]
 [5.9 3. 4.2 1.5]
 [6.2 3.2 4. 1. ]
 [6.5 2.9 4.5 1.8]
 [5.6 2.9 3.6 1.3]
 [6.7 3.1 4.4 1.4]
 [6.2 3.4 4.3 1.3]
 [5.8 2.7 4.1 1. ]
 [6.2 3.2 4.5 1.5]
 [6.4 3.2 4.5 1.8]
 [5.9 3.2 4.8 1.8]
 [6.1 2.8 4.7 1.3]
 [6.4 3.1 4.7 1.3]
 [7.3 2.9 4.6 1.8]
 [6.8 3. 4.9 1.8]
 [6.4 2.8 5.6 2.1]
 [7.2 3. 5.8 1.6]
 [7.4 2.8 6.1 1.9]
 [7.9 3.8 6.4 2. ]
 [6.4 2.8 5.6 2.2]
 [6.3 2.8 5.1 1.5]
 [6.1 2.6 5.6 2.4]
 [6.4 3.1 5.5 1.8]
 [6.9 3.1 5.4 2.1]
 [6.7 3.1 5.5 2.4]
 [6.9 3.1 5.1 2.3]
 [5.8 2.7 5.1 1.9]
 [6.8 3.2 5.9 2.3]
 [6.7 3. 5.2 2.3]
 [6.3 3. 5.2 2.3]
 [6.5 3. 5.2 2. ]
 [6.2 3.4 5.4 2.3]
 [5.9 3. 5.1 1.8]]
```

```
In [1989]: print('The names of the dataset variables:\n',iris['feature_names'])

The names of the dataset variables:
['sepal length (cm)', 'sepal width (cm)', 'petal length (cm)', 'petal width (cm)']
```

```
In [1990]: print('Name of classes: "%s".' % list(iris.target_names))

Name of classes:
('setosa', 'versicolour', 'virginica')
```

B. Data normalization

The sklearn.preprocessing package provides several common utility functions and transformer classes to change raw feature vectors into a representation that is more suitable for the downstream estimators. Standardization of datasets is a common requirement for many machine learning estimators implemented in the scikit: they might behave badly if the individual feature do not more or less look like standard normally distributed data. Gaussian with zero mean and unit variance. In practice we often ignore the shape of the distribution and just transform the data to center it by removing the mean value of each feature. Then, scale it by dividing non-constant features by their standard deviation. For instance, many elements used in the objective function of a learning algorithm (such as the RBF kernel of Support Vector Machines or the l1 and l2 regularizers of linear models) assume that all features are centered around zero and have variance in the same order. If a feature has a variance that is orders of magnitude larger than others, it might dominate the objective function and make the estimator unable to learn from other features correctly as expected.

```
In [1991]: from sklearn.preprocessing import scale
```

1- Create the following matrix X: 1, -1, 2, 2, 0, 0, 0, 1, -1

```
In [1992]: X=[1,-1,2],[2,0,0],[0,1,-1]]

2- Print the matrix and compute the mean of the variables.
```

```
In [1993]: print(X)
M = np.mean(X)
print('The mean of the variables is: (%.2f)' % M)

[[1,-1,2],[2,0,0],[0,1,-1]]
The mean of the variables is: 0.44
```

3- Use the scale function to normalize X. Analyze the result.

```
In [1994]: scaled = scale(X)
scaled

array([[ 0.22744487, -1.22474487,  1.33630621],
       [-1.22474487,  1.22474487, -0.26262244],
       [-1.22474487,  1.22474487,  1.06904497]])
```

4- Compute the mean and the variance of the scaled X. What can you conclude?

```
In [1995]: mean = np.mean(scaled)
variance = np.var(scaled)
print('The mean of the scaled matrix is: (%.2f)' % mean)
print('The variance of the scaled matrix is: (%.2f)' % variance)
print('Mean = (%)' % np.mean(scaled, axis=0))
print('variance = (%)' % np.mean(scaled, axis=0))

The mean of the scaled matrix is: 0.00
The variance of the scaled matrix is: 1.00
Mean=[ 0.  0.  0.]
Variance = [ 1.  1.  1.]
```

We have scaled our matrix so that all 3 features could be in the same scaling. We verified this by checking the mean of the matrix and variance and we can observe that the mean is 0 and variance is 1, so they make our data unitless. This is because we have used the Standardization type for scaling. But there also is Normalization that is used when we want to bound our values between two numbers, typically, between 0.0 or 1.0. Machine learning algorithm just sees number and if there is a vast difference in the range say few ranging in thousands and few ranging in the tens, and makes the underlying assumption that higher ranging numbers have superiority over other sort. So these more significant number starts playing a more decisive role while training the model.

C. MinMax Normalization

An alternative standardization is scaling features to lie between a given minimum and maximum value, often between zero and one. This can be achieved using MinMaxScaler.

1- Create the following matrix X2: 1, -1, 2, 2, 0, 0, 0, 1, -1

```
In [1996]: X2=[1,-1,2],[2,0,0],[0,1,-1]]

2- Print the matrix and compute the mean of the variables.
```

```
In [1997]: print(X2)
M = np.mean(X2)
print('The mean of the variables is: (%.2f)' % M)

[[1,-1,2],[2,0,0],[0,1,-1]]
The mean of the variables is: 0.444
```

3- Normalize the data using a MinMaxScaler. Print the scaled matrix and compute the mean and the variance. What can you conclude?

```
In [1998]: from sklearn.preprocessing import MinMaxScaler
print(X2)
# define min max scaler
mm = MinMaxScaler()
# transform data
scaled = scaler.fit_transform(X2)
print(scaled)

[[[1,-1,2],[2,0,0],[0,1,-1]]
 [0.5  0.  0.
  0.5  0.33333333
  1.  0.
  0.  0.  ]]
```

```
In [1999]: mean = np.mean(scaled)
variance = np.var(scaled)
print('The mean of the scaled matrix is: (%.2f)' % mean)
print('The variance of the scaled matrix is: (%.2f)' % variance)

The mean of the scaled matrix is: 0.48
The variance of the scaled matrix is: 0.17
```

MinMaxScaler transforms features by scaling each feature to a given range. This estimator scales and translates each feature individually such that it is in the given range on the training set, e.g., between zero and one. This scaler shrinks the data within the range of -1 to 1 if there are negative values. This is why the variance of the scaled matrix is so small and mean is not 0, but 0.48.

D. Data visualization

1- Import the Iris dataset using : iris = datasets.load_iris()

This dataset consists of 3 different types of irises (Setosa, Versicolour, and Virginical) petal and sepal length, stored in a 150x4 numpy ndarray

The rows being the samples and the columns being: Sepal Length, Sepal Width, Petal Length and Petal Width. More information about this dataset can be found here:

Length and Petal Width. More information about this dataset can be found here: <https://archive.ics.uci.edu/ml/datasets/Iris>

(iris.data), the corresponding label (target), the names of the variables (feature_names) and the name of classes (target_names).

```
In [100]: iris = datasets.load_iris()
```

2- Plot the data points into 2D dimension with all the possible combination between variables and use the label for the color points. Vizually, which is the better combination of variables? Justify the answer.

```
In [100]: import pandas as pd
import matplotlib.pyplot as plt
df = pd.DataFrame(data=iris.data, columns=iris.feature_names)
features = df.columns.values.tolist()

# map target to species class
target_map = {i: iris.target_names[i] for i in range(0, len(iris.target_names))}

df['species'] = pd.Series(iris.target).map(target_map)
fig = px.scatter_matrix(
    dimensions=features,
    color='species')
fig.show()
```

Better combination of variables are: petal length & petal width, as we can observe a visible behavior as if the petal width and length are small, then it means these are of type of setosa class, and if the petal length is between 3 and 5 cm and petal width is between 1.2 and 1 cm then the flower is versicolour class. If the flower has petal length between 5 and 7 cm and petal width between 2.8 and 2 cm then it is mostly of class virginica.

3- Compute the correlations between each pair of variables by using the corcoef function of numpy package. Can you validate the answer of the question C.2 by using the obtained correlations? Justify your response.

```
In [100]: variables=iris.feature_names

['sepal length (cm)',
 'sepal width (cm)',
 'petal length (cm)',
 'petal width (cm)']
```

```
Out[100]: array([[ 0.22744487, -1.22474487,  1.33630621],
       [-1.22474487,  1.22474487, -0.26262244],
       [-1.22474487,  1.22474487,  1.06904497]])
```

```
In [100]: print('Correlation between (0) and (1) is (%.2f)' % (str(iris.feature_names[0]), str(iris.feature_names[1])))
print('Correlation between (0) and (1) is (%.2f)' % (str(iris.feature_names[0]), str(iris.feature_names[1])))
print('Correlation between (0) and (1) is (%.2f)' % (str(iris.feature_names[0]), str(iris.feature_names[2])))
print('Correlation between (0) and (1) is (%.2f)' % (str(iris.feature_names[0]), str(iris.feature_names[3])))
print('Correlation between (0) and (1) is (%.2f)' % (str(iris.feature_names[2]), str(iris.feature_names[3])))

Correlation between sepal length (cm) and sepal width (cm) is -0.117569784330018
Correlation between sepal length (cm) and petal length (cm) is 0.817941262715757
Correlation between sepal width (cm) and petal length (cm) is 0.817941262715757
Correlation between sepal width (cm) and petal width (cm) is -0.428440104303597
Correlation between petal length (cm) and petal width (cm) is 0.9628654314027965
```

As we can observe the biggest correlation is of variables petal length and petal width: 0.9629%

4- Subplots in matplotlib Test & Analyze the following code:

```
In [100]: import matplotlib.pyplot as plt

fig = plt.figure()

ax1 = fig.add_subplot(231)
ax1.scatter(iris.data[:, 0], iris.data[:, 1], c=iris.target, cmap=plt.cm.Set1)
ax1.set_xlabel(iris.feature_names[0])
ax1.set_ylabel(iris.feature_names[1])

ax2 = fig.add_subplot(232)
ax2.scatter(iris.data[:, 0], iris.data[:, 2], c=iris.target, cmap=plt.cm.Set1)
ax2.set_xlabel(iris.feature_names[0])
ax2.set_ylabel(iris.feature_names[2])

ax3 = fig.add_subplot(233)
ax3.scatter(iris.data[:, 1], iris.data[:, 3], c=iris.target, cmap=plt.cm.Set1)
ax3.set_xlabel(iris.feature_names[1])
ax3.set_ylabel(iris.feature_names[3])

ax4 = fig.add_subplot(234)
ax4.scatter(iris.data[:, 1], iris.data[:, 2], c=iris.target, cmap=plt.cm.Set1)
ax4.set_xlabel(iris.feature_names[1])
ax4.set_ylabel(iris.feature_names[2])

ax5 = fig.add_subplot(235)
ax5.scatter(iris.data[:, 2], iris.data[:, 3], c=iris.target, cmap=plt.cm.Set1)
ax5.set_xlabel(iris.feature_names[2])
ax5.set_ylabel(iris.feature_names[3])

plt.tight_layout()
fig = plt.gcf()

pltioy_fig = plt.mpl.to_ploiy(fig)
pltioy_fig['layout']['title'] = 'Simple Subplot Example Title'
pltioy_fig['layout']['margin'].update({'t':40})
pltioy_fig.show()
```

```
In [100]: import matplotlib.pyplot as plt

fig = plt.figure()

ax1 = fig.add_subplot(231)
ax1.scatter(iris.data[:, 0], iris.data[:, 1], c=iris.target, cmap=plt.cm.Set1)
ax1.set_xlabel(iris.feature_names[0])
ax1.set_ylabel(iris.feature_names[1])

ax2 = fig.add_subplot(232)
ax2.scatter(iris.data[:, 0], iris.data[:, 2], c=iris.target, cmap=plt.cm.Set1)
ax2.set_xlabel(iris.feature_names[0])
ax2.set_ylabel(iris.feature_names[2])

ax3 = fig.add_subplot(233)
ax3.scatter(iris.data[:, 1], iris.data[:, 3], c=iris.target, cmap=plt.cm.Set1)
ax3.set_xlabel(iris.feature_names[1])
ax3.set_ylabel(iris.feature_names[3])

ax4 = fig.add_subplot(234)
ax4.scatter(iris.data[:, 1], iris.data[:, 2], c=iris.target, cmap=plt.cm.Set1)
ax4.set_xlabel(iris.feature_names[1])
ax4.set_ylabel(iris.feature_names[2])

ax5 = fig.add_subplot(235)
ax5.scatter(iris.data[:, 2], iris.data[:, 3], c=iris.target, cmap=plt.cm.Set1)
ax5.set_xlabel(iris.feature_names[2])
ax5.set_ylabel(iris.feature_names[3])

plt.tight_layout()
fig = plt.gcf()

pltioy_fig = plt.mpl.to_ploiy(fig)
pltioy_fig['layout']['title'] = 'Simple Subplot Example Title'
pltioy_fig['layout']['margin'].update({'t':40})
pltioy_fig.show()
```

```
In [100]: iris.feature_names = ['sepal length (cm)', 'petal length (cm)', 'petal width (cm)']
iris.data = iris.data[:, [0, 2, 3]]

print('Correlation between (0) and (1) is (%.2f)' % (str(iris.feature_names[0]), str(iris.feature_names[1])))
print('Correlation between (0) and (1) is (%.2f)' % (str(iris.feature_names[0]), str(iris.feature_names[2])))
print('Correlation between (0) and (1) is (%.2f)' % (str(iris.feature_names[1]), str(iris.feature_names[2])))

Correlation between sepal length (cm) and petal length (cm) is 0.817941262715757
Correlation between sepal length (cm) and petal width (cm) is 0.817941262715757
Correlation between petal length (cm) and petal width (cm) is 0.9628654314027965
```

```
In [100]: iris.feature_names = ['petal length (cm)', 'petal width (cm)']
iris.data = iris.data[:, [1, 2]]

2- The PCA and LDA methods can be imported from the following packages : from sklearn.decomposition import PCA from sklearn.l
```

```
In [101]: from sklearn.decomposition import PCA
from sklearn.discriminant_analysis import LinearDiscriminantAnalysis as LDA

3- Analyze the help of these functions (pca and lda) and apply them on the iris dataset. You have to use here pca.fit(iris).transform(iris) and save the results in iris_pca for the PCA and iris_lda for the LDA.
```

```
In [101]: X = iris.data
y = iris.target
target_names = iris.target_names

pca = PCA(n_components=2)
iris_pca = pca.fit(X).transform(X)

lda = LDA(n_components=2)
iris_lda = lda.fit(X, y).transform(X)
```

4- Plot the data points on the new obtained projections : one image for the PCA and another for the LDA and use the label as color for the points. You should use the following function from Plyton: figure, scatter, title, xlim, ylim, xlabel, ylabel et show. Which difference you can see between the both results? Explain?

```
In [101]: plt.figure()
for c, i, target_name in zip('rgb', (0, 1, 2), target_names):
    plt.scatter(iris_pca[:, 0], iris_pca[:, 1], c=i, label=target_name)
plt.legend()
plt.ylim(-1.5, 1.5)
plt.xlim(-3, 3)
plt.title('PCA of IRIS dataset')
```

```
In [101]: plt.figure()
for c, i, target_name in zip('rgb', (0, 1, 2), target_names):
    plt.scatter(iris_lda[:, 0], iris_lda[:, 1], c=i, label=target_name)
plt.legend()
plt.ylim(-10, 10)
plt.xlim(-10, 10)
plt.title('LDA of IRIS dataset')
```

```
In [101]: plt.show()

PCA of IRIS dataset
```

```
LDA of IRIS dataset
```

As we know, to use the PCA we must first scale the data, but the LDA can handle the data that is not scaled. We can observe that the points have almost the same distribution, but the ranges are different: for PCA the X axis is between -3.2 and 4, and Y axis is between -1.5 and 1.5. And for LDA X axis is between range -10 and 10 and for Y axis is between -3.2 and 3. We can picture PCA as a technique that finds the directions of maximal variance. In contrast to PCA, LDA attempts to find a feature space that maximizes class separability.

```
In [101]: from sklearn.manifold import TSNE
tsne = TSNE(n_components=2)
tsne.fit(iris.data)
iris_tsne = tsne.fit_transform(iris.data)

plt.scatter(iris_tsne[:, 0], iris_tsne[:, 1], c=iris.target, cmap=plt.cm.Set1)
plt.xlabel(iris.feature_names[0])
plt.ylabel(iris.feature_names[1])
plt.title('TSNE')
```

```
C:\Users\Bef\AppData\Local\Programs\Python\Python39\lib\site-packages\sklearn\manifold\_t_sne.py:780: FutureWarning
The default initialization in TSNE will change from 'random' to 'pca' in 1.2.

C:\Users\Bef\AppData\Local\Programs\Python\Python39\lib\site-packages\sklearn\manifold\_t_sne.py:790: FutureWarning
The default learning rate in TSNE will change from 200.0 to 'auto' in 1.2.

C:\Users\Bef\AppData\Local\Programs\Python\Python39\lib\site-packages\sklearn\manifold\_t_sne.py:780: FutureWarning
The default initialization in TSNE will change from 'random' to 'pca' in 1.2.

C:\Users\Bef\AppData\Local\Programs\Python\Python39\lib\site-packages\sklearn\manifold\_t_sne.py:790: FutureWarning
The default learning rate in TSNE will change from 200.0 to 'auto' in 1.2.
```

```
TSNE
```

```
F. MNIST dataset
```

It is possible to load a dataset directly from mldata.org which contains a lot of available datasets using the function datasets.fetch_mldata.

1- Import the MNIST original.

```
In [101]: from sklearn.model_selection import train_test_split
from sklearn.datasets import load_digits
digits = load_digits()
print(digits.data.shape)
```

```
In [101]: (1797, 64)
```

```
Out[101]: digits
```

```
digits
```

```
dataset.matrix
```

```
print('Number of data')
print(iris.data.shape)
```

It means we have 1797 observations and 64 variables.

```
In [101]: print('Number of classes')
print(iris.target_names)
print('Number of variables: (%.2f)' % (format(len(digits['feature_names'])))
```

```
Variables
['pixel_0_0', 'pixel_0_1', 'pixel_0_2', 'pixel_0_3', 'pixel_0_4', 'pixel_0_5', 'pixel_0_6', 'pixel_0_7', 'pixel_0_8', 'pixel_0_9', 'pixel_1_0', 'pixel_1_1', 'pixel_1_2', 'pixel_1_3', 'pixel_1_4', 'pixel_1_5', 'pixel_1_6', 'pixel_1_7', 'pixel_1_8', 'pixel_1_9', 'pixel_2_0', 'pixel_2_1', 'pixel_2_2', 'pixel_2_3', 'pixel_2_4', 'pixel_2_5', 'pixel_2_6', 'pixel_2_7', 'pixel_2_8', 'pixel_2_9', 'pixel_3_0', 'pixel_3_1', 'pixel_3_2', 'pixel_3_3', 'pixel_3_4', 'pixel_3_5', 'pixel_3_6', 'pixel_3_7', 'pixel_3_8', 'pixel_3_9', 'pixel_4_0', 'pixel_4_1', 'pixel_4_2', 'pixel_4_3', 'pixel_4_4', 'pixel_4_5', 'pixel_4_6', 'pixel_4_7', 'pixel_4_8', 'pixel_4_9', 'pixel_5_0', 'pixel_5_1', 'pixel_5_2', 'pixel_5_3', 'pixel_5_4', 'pixel_5_5', 'pixel_5_6', 'pixel_5_7', 'pixel_5_8', 'pixel_5_9', 'pixel_6_0', 'pixel_6_1', 'pixel_6_2', 'pixel_6_3', 'pixel_6_4', 'pixel_6_5', 'pixel_6_6', 'pixel_6_7', 'pixel_6_8', 'pixel_6_9', 'pixel_7_0', 'pixel_7_1', 'pixel_7_2', 'pixel_7_3', 'pixel_7_4', 'pixel_7_5', 'pixel_7_6', 'pixel_7_7']
Number of variables: 64
```

```
In [101]: print('Classes')
print(digits['target_names'])
print('Number of classes: (%.2f)' % (format(len(digits['target_names'])))
```

```
Classes
['setosa', 'versicolour', 'virginica']
Number of classes: 10
```

3- Use the MNIST visualization tutorial to visualize the digits.

```
In [102]: import matplotlib.pyplot as plt
for i in range(0,10):
    plt.imshow(digits.images[i])
    plt.show()
```


