

Analyse et Traitement de l'Information

Extra TP1: Hierarchical Clustering (20 pts).

The MNIST dataset of handwritten digits ¹ consists of 28×28 grayscale images. To download the dataset use this code:

```
1 from sklearn.datasets import fetch_openml
2 from sklearn.model_selection import train_test_split
3
4 images, labels = fetch_openml('mnist_784', \
5                               return_X_y=True, as_frame=False, parser='auto')
6 train_images, test_images, train_labels, test_labels = \
7     train_test_split(images, labels, random_state=42)
```

To filter images of the specific class you can use the following method:

```
1 import numpy as np
2
3 def select_with_label(images, labels, desired_labels):
4     mask = np.isin(labels, desired_labels)
5     return images[mask], labels[mask]
6
7 images_of_two, labels_of_two = \
8     select_with_label(train_images, train_labels, desired_labels=['2'])
9 images_of_odd, labels_of_odd = \
10    select_with_label(train_images, train_labels, \
11                      desired_labels=['1', '3', '5', '7', '9'])
```

You can copy-paste from *this link*.

¹<http://yann.lecun.com/exdb/mnist/>

Clustering

Classification is a *supervised* problem, where labeled (training) data is used to label the unlabeled (testing) data.

Clustering is an *unsupervised* problem, where the objective is to group a given (unlabeled) dataset $\{\mathbf{x}\}_{i=1}^n$ into k clusters, so that similar samples appear in the same cluster, and dissimilar samples are in different clusters.

Hierarchical Clustering

We will take a look at another useful type of Clustering – Agglomerative Clustering. The idea is to unite the clusters with the shortest distance between them. This method has a high computational complexity so it is applied only **for small datasets**.

Linkage is a method which will be used to compute the distance between the clusters.

- *average* - the distance between two clusters is the average distance between all possible pairs of points where one belongs to one cluster and another one to another
- *single* - the distance between two clusters is the minimal distance between two points where one belongs to one cluster and another one to another
- *complete* - the distance between two clusters is the maximal distance between two points where one belongs to one cluster and another one to another

```
1 class AgglomerativeClustering:
2     def __init__(self, n_clusters=16, linkage="complete"):
3         raise NotImplementedError()
4
5     def fit_predict(self, X, y=None):
6         #this function returns a cluster number to each element of X
7         raise NotImplementedError()
8
```

Questions

1. Implement AgglomerativeClustering with *complete* distance using *numpy* package. Explain the algorithm in the report.
2. Take 200 instances of “3”s and “7”s from the MNIST dataset. For the case $k = 2$, give the confusion matrix for both Agglomerative and *k*-means clustering. (you can use sklearn implementation of KMeans for this task)
3. Sample 200 instances of “3”s and “5”s. Perform both clustering approaches with $k = 2$. Build the confusion matrix. Based on your results, explain the difference between these two clustering problems:
 - “3” vs. “7”
 - “3” vs. “5”.

Submission

Please archive your report and codes in “Prénom Nom.zip” (replace “Prénom” and “Nom” with your real name), and upload to “Upload Extra TP1 - Hierarchical Clustering.” on <https://moodle.unige.ch> before **Monday, December 1 2023, 21:59 PM**. Note, that the assessment is mainly based on your report, which should include your answers to all questions and the experimental results. *Importance is given on the mathematical explanations of your works and your codes should be commented*

Please use the template from [the first TP](#).