



SIA 3611

Année universitaire 2022-2023

---

# TP Machine Learning

## TP 4 : AutoML

---

Vassilis.Christophides@ensea.fr  
guillaume.renton@ensea.fr

# Table des matières

<b>1</b>	<b>Getting started with new datasets</b>	<b>1</b>
1.1	To do	1
1.2	To code	2
1.3	To code	2
1.4	To do	2
1.5	To code	3
1.6	To code	3
<b>2</b>	<b>AutoML</b>	<b>3</b>
2.1	To do	4
2.2	To code	4
<b>3</b>	<b>Bonus step</b>	<b>4</b>

## Introduction

In previous TP, you have learned to use machine learning for different kind of tasks, from regression to clustering through classification. In this TP, you are going to use the earned knowledge on new datasets for regression and classification.

You are going to use 2 new datasets in this TP. First one is california housing, whose target variable is the value of houses in california, expressed in hundred of thousand of dollars. For each house, a set of 9 features is available. There is a total of 20 060 data.

Second one is MNIST, a very popular dataset for handwritten recognition and image classification. The original dataset is made of 60 000 training images of shape 28x28 of handwritten digits from 0 to 9, and 10 000 images for test dataset. For computational time, you will work on a given random subset of MNIST made of 6000 images in train and 1000 images in test.

Objectives :

- Apply your knowledge on new datasets
- Tune models hyperparameters and explore metrics
- Apply principal components analysis and understand its effects on both dataset
- Understand and use Cross-Validation
- Use AutoML to find interesting models

The notebook corresponding to this TP is available on Moodle. Answer the questions directly on this notebook. The notebook has to be submitted before next week. The work can be done by pair, but one notebook is awaited per student.

## 1 Getting started with new datasets

### Regression

In first part of step 1, you will work on the regression problem with the dataset california housing.

## 1.1 To do

Execute the following cell to load the california housing dataset and normalize it.

```
1 from sklearn.datasets import fetch_california_housing
2 from sklearn.preprocessing import normalize
3 import numpy as np
4
5 X, y = fetch_california_housing(return_X_y = True)
6 X = normalize(X)
```

## 1.2 To code

Apply [Stochastic Gradient Descent](#) and [SVR](#) methods and cross validate your results using 5 folders. For this, you can either use the function [cross\\_val\\_score](#) (or any other method for cross validation in sklearn) or either compute yourself the cross validation. According to a relevant metric optimize both methods. For SGD you will optimize the value of alpha for both L2 and L1 penalty score. For SVR, you will optimize the kernel. Be careful with the metric if you use `cross_val_score`, the returned values are often negative.

### Question 1

According to your metric, which method obtain the best result ?

### Question 2

What is the interest of using cross validation in general ? Is it relevant in this particular case ?

## 1.3 To code

Transform your data according to [principal component analysis](#), and optimize the number of components according to the same metric than previously for both models.

```
1 from sklearn.decomposition import PCA
2 pca = PCA(n_components = 1)
3 X_pca = pca.fit_transform(X)
```

### Question 3

What is the interest of Principal Component Analysis in general ? Is it relevant here ?

## Classification

## 1.4 To do

Execute the following cells to load a subset of MNIST dataset. Since the dataset is already divided into training/test, we won't use cross validation this time.

```
1 import pickle
2 with open("data/mnist.pkl", "rb") as f:
3     ((X_train, y_train), (X_test, y_test)) = pickle.load(f)
```

## 1.5 To code

Compute classification on those images using a [KNN](#) classifier and an [Adaboost](#) classifier. For each classifier, optimize the parameters according to a relevant metric. For the KNN classifier, you will optimize the number of neighbor while for the Adaboost classifier, you will optimize the base estimator along with the number of estimators (for the basis estimator, limit yourself to different depth of decision tree classifier).

Also, for each model, compute the [confusion matrix](#).

### Question 4

According to your metric, which method obtain the best results ?

### Question 5

According to the confusion matrix, which class is the easiest to classify ? Which ones are the most difficult ? Which ones are the most confused with each other ?

### Bonus

For the Adaboost classifier, explore other classifier as base estimators. What are the limitations about those estimators ?

## 1.6 To code

Transform your data according to principal component analysis, and optimize the number of components according to the same metric than previously for each classifier.

Once again, compute the confusion matrix for each model.

### Question 6

Is the use of PCA relevant here ?

### Question 7

Did your answers from question 5 changed with PCA ?

## 2 AutoML

In this second section, we discuss on the utilisation of AutoML tools, such as auto-sklearn. If you are using colab or don't have auto-sklearn installed (much likely), you may need to run the following cell at first in order to install auto-sklearn. This will require you to restart the runtime (a prompt will invite you to).

Restarting the runtime will clear all your variables and imported libraries, so you will need to import them again.

```

1 !pip install --force-reinstall scipy==1.6
2 !pip install --force-reinstall auto-sklearn==0.15

```

## 2.1 To do

Execute the following cells.

```

1 from sklearn.model_selection import train_test_split
2 X, y = fetch_california_housing(return_X_y = True)
3 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

```

```

1 import autosklearn.regression
2 import sklearn.model_selection
3 import sklearn.datasets
4 import os, shutil
5 from sklearn.metrics import mean_squared_error, mean_absolute_error
6
7 automl = autosklearn.regression.AutoSklearnRegressor(
8     include = {'regressor': ["libsvm_svr", "sgd"]},
9     time_left_for_this_task=120,
10    per_run_time_limit=30,
11    tmp_folder='/tmp/california_housing_tmp',
12 )
13 automl.fit(X_train, y_train, dataset_name='California_Housing')
14
15 print(automl.leaderboard())
16
17 y_pred = automl.predict(X_test, y_test)
18 print("MSE = ", mean_squared_error(y_test, y_pred))
19 print("MRE = ", mean_absolute_error(y_test, y_pred))

```

```

1 from pprint import pprint
2 pprint(automl.show_models(), indent=4)

```

## Question 8

What are the evaluated models by autoML? Which model obtain the best performance? What are the parameters of the best model?

## 2.2 To code

With the help of the previous code, use autoML for the classification task on MNIST, by limiting the exploration to KNN and Adaboost.

## Question 9

What are the evaluated models by autoML? Which model obtain the best performance? What are the parameters of the best model?

## 3 Bonus step

As a bonus step, have fun and remove a maximum of constraints of your autoML model. Which model obtain the best performances? Describe the parameters of this model. You can do it for either for regression or classification or both.