

FTML practical session 1: 2023/03/02

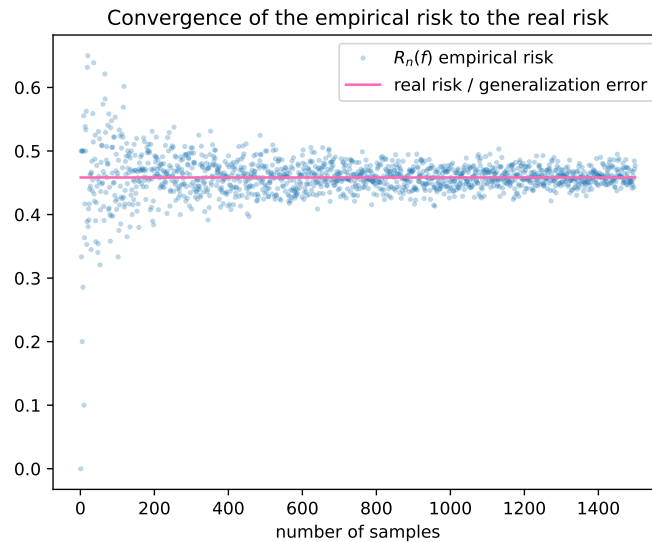


TABLE DES MATIÈRES

1	Miscellaneous python	2
1.1	Environments	2
1.2	Good python habits	2
1.2.1	Format your code	2
1.2.2	Sort your imports	2
1.2.3	Code style	2
1.3	Demos	2
1.4	Operations in Python	2
2	Experimenting with the law of large numbers	2
2.0.1	More Python profiling	3
3	Bayes risks	3
3.1	Setting	3
3.1.1	Risks	3
3.1.2	Estimating the generalization error	4
3.1.3	Bayes estimator	4
3.2	Two examples	4
3.2.1	Problem 1 : penalty shootout	4
3.2.2	Problem 2 : prediction of the number of spotify streams	4

INTRODUCTION

The goal of this practical session and of the next one is to experiment with some concepts that are specific and central to machine learning : the law of large numbers and risks (empirical risk and generalization error). During the lectures that will follow these two practical sessions, we will study and formalize the definitions of the different types of risks in more depth. You can do the different 3 parts in whatever order but the natural order is rather 2, 3, ?? (1 is not really an exercise). You do not have to finish everything during the session.

1 MISCELLANEOUS PYTHON

1.1 Environments

To install libraries, you can use virtual environments.

<https://docs.python.org/3/library/venv.html>

The list of libraries used will be in the `practical_sessions/requirements.txt` file, which will be updated periodically. You can use it to install all libraries directly with pip, e.g. with `pip install -r requirements.txt`.

1.2 Good python habits

You can explore these tools later.

1.2.1 Format your code

<https://github.com/psf/black>

1.2.2 Sort your imports

<https://github.com/PyCQA/isort>

1.2.3 Code style

<https://realpython.com/python-pep8/>

1.3 Demos

In `practical_sessions/tp1/demos/`, you can find a couple of simples demo files to use matplotlib, numpy (if needed).

1.4 Operations in Python

Time complexity of elementary operations in python :

<https://wiki.python.org/moin/TimeComplexity>

2 EXPERIMENTING WITH THE LAW OF LARGE NUMBERS

Let us consider the same variable as in exercice P3 : $Z_2 = Z_1$ and is Z_1^2 is a uniform law in $[1, 2]$. We have seen that $E[Z_2] = 7/3$. Hence, according to the law of large numbers, the empirical average of n draws of this variable converges in probability to this expected value.

In `exercice_1/law_of_large_numbers.py`, the function `empirical_average_loop` computes the empirical average with a for loop.

- Edit the function `empirical_average_array` in order to use numpy and array operations to perform the same computation in an optimized way, only using array operations and without a loop.
- Compare the speed of the methods by monitoring the `profile.prof` profiling file, for instance using `snakeviz profile.prof`.

https://en.wikipedia.org/wiki/Array_programming
<https://jiffyclub.github.io/snakeviz/>

2.0.1 More Python profiling

Profile individual lines :

<https://pypi.org/project/line-profiler/>

Profile memory usage :

<https://pypi.org/project/memory-profiler/>

3 BAYES RISKS

3.1 Setting

The goal of this exercise is to introduce the notion of empirical risk (risque empirique), generalization error (risque réel), and Bayes risk for some simple problems. We consider a supervised learning problem,

- an input space \mathcal{X}
- an output space \mathcal{Y}
- a loss function l
- and a dataset $D_n = \{(x_1, y_1), \dots, (x_n, y_n)\}$ of n samples.

3.1.1 Risks

An **estimator** f is a mapping from the input space to the output space.

Definition 1. Risks

Let l be a loss. The **risk** (or **statistical risk**, **generalization error**, **test error**, **risque réel in french**) of estimator f writes

$$R(f) = E_{(X,Y) \sim \rho} [l(Y, f(X))] \quad (1)$$

Here, X is the random variable that represents the inputs, and Y the variable that represents the output. ρ is the joint law.

The **empirical risk (ER)** (risque empirique) of an estimator f writes

$$R_n(f) = \frac{1}{n} \sum_{i=1}^n l(y_i, f(x_i)) \quad (2)$$

We emphasize that the risks depends on the loss l .

In supervised learning, we only have access to the empirical risk R_n but we actually want to find an estimator f which as a small generalization error! The problem is that in real situations, we do not have access to ρ , which allows its computation.

3.1.2 Estimating the generalization error

However, thanks to the law of large numbers, if we use a **fixed estimator** f , $R_n(f) \rightarrow R(f)$ when $n \rightarrow \infty$. Hence, if we have a large enough number of samples n , $R_n(f)$ is a good estimation of $R(f)$. The problem will then be : how large is sufficient ? the answer will depend on the context.

3.1.3 Bayes estimator

Under some simple hypotheses, for a given supervised learning problem, there exists an optimal estimator f^* called the Bayes estimator, which minimized the generalization error, given a distribution ρ . Its generalization error $R(f^*)$ is called the **Bayes risk**.

3.2 Two examples

For the two following problems, estimate the generalization error of various estimators of your choice by simulating the random variables, and try to find the Bayes estimator !

3.2.1 Problem 1 : penalty shootout

We represent a football penalty shootout. $X \in \{0, 1\}$ is the variable representing that team 1 shoots first. $Y \in \{0, 1\}$ is the variable representing the fact that team 1 wins. We assume that :

- X is uniformly distributed.
- If $X = 1$, Y follows a Bernoulli law of parameter 0.6. If $X = 0$, Y follows a law of parameter 0.4.
- l is the 0 – 1 loss (1 if there is a mistake, 0 otherwise)

For this setting, the Bayes risk is 0.4.

3.2.2 Problem 2 : prediction of the number of spotify streams

A music label is interested in predicting the number of streams of an artist, as a function of the investment. We will consider that the investment is represented by the number of persons who work with the artist during the production, which is a proxy to the investment. This variable is noted X . More precisely, we predict the number of streams of the song on a spotify, noted Y , during the first week after release, as a function of X . We assume that :

- $X - 1 \in \mathbb{N}$ follows a binomial law of parameters $n_X = 20$ and $p_X = 0.2$. Hence, $X > 0$.
- Given a value x of X , Y follows a binomial law of parameters $n_Y(x) = 3^x$ and $p_Y(x) = 0.5$.
- l is the squared loss.

For this setting, the Bayes risk is around 627.