

# A CHAINING ALGORITHM FOR ONLINE NONPARAMETRIC REGRESSION

PIERRE GAILLARD

PIERRE-P.GAILLARD@EDF.FR

SÉBASTIEN GERCHINOVITZ

SEBASTIEN.GERCHINOVITZ@MATH.UNIV-TOULOUSE.FR

## Problem: online nonparametric regression

For each round  $t = 1, \dots, T$ ,

- Environment chooses  $x_t \in \mathcal{X}$
- Forecaster predicts  $\hat{y}_t \in \mathbb{R}$
- Forecaster observes  $y_t \in \mathbb{R}$  and suffers square loss  $(\hat{y}_t - y_t)^2$

**Goal:** minimize the regret

$$\text{Reg}_T(\mathcal{F}) \stackrel{\text{def}}{=} \sum_{t=1}^T (y_t - \hat{y}_t)^2 - \inf_{f \in \mathcal{F}} \sum_{t=1}^T (y_t - f(x_t))^2$$

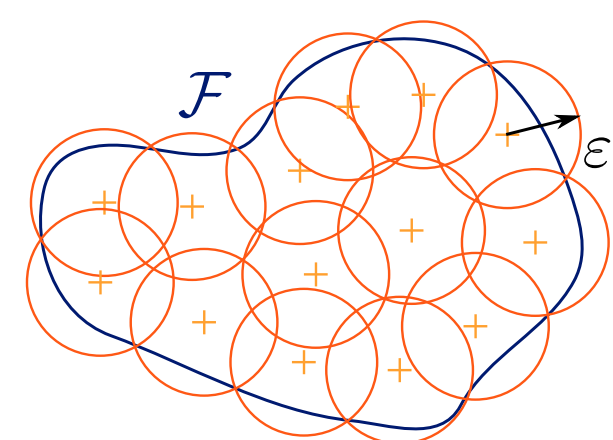
over nonparametric function classes  $\mathcal{F} \subset \mathbb{R}^{\mathcal{X}}$ .

## Main contributions

We design an algorithm that achieves a **Dudley-type** regret bound (in the same vein as in Rakhlin and Sridharan [3] but in a constructive fashion):

$$\text{Reg}_T(\mathcal{F}) \leq \square B^2(1 + \log \mathcal{N}_\infty(\mathcal{F}, \gamma)) + \square B \sqrt{T} \int_0^\gamma \sqrt{\log \mathcal{N}_\infty(\mathcal{F}, \varepsilon)} d\varepsilon,$$

where  $\gamma \in (\frac{B}{T}, B)$  is a parameter of the algorithm,  
 $B$  is an upper bound on  $\max_{1 \leq t \leq T} |y_t|$ ,  
 $\log \mathcal{N}_\infty(\mathcal{F}, \varepsilon)$  denotes the metric entropy of  $\mathcal{F}$ .



Function class	Metric entropy	Regret bound
	$\varepsilon^{-p}$ $p \in (0, 2)$	$T^{p/(p+2)}$
Lipschitz on $[0, 1]$	$\varepsilon^{-1}$	$T^{1/3}$
$\beta$ -Hölder on $[0, 1]$	$\varepsilon^{-1/\beta}$ $\beta > 1/2$	$T^{1/(2\beta+1)}$
Sparse lin. reg.	$\log \binom{d}{s} + s \log(1 + 1/(\varepsilon\sqrt{s}))$	$s \log(1 + dT/s)$

Second contribution: **efficient** version for Hölder classes (costs a log factor).

## Chaining algorithm: main ideas

**Chaining approximation:** we approximate any function  $f \in \mathcal{F}$  by a **sequence of refining approximations**  $\pi_0(f) \in \mathcal{F}^{(0)}, \pi_1(f) \in \mathcal{F}^{(1)}, \dots$ , such that for all  $k \geq 0$ ,

$$\sup_f \|\pi_k(f) - f\|_\infty \leq \gamma/2^k$$

and

$$\text{Card } \mathcal{F}^{(k)} = \mathcal{N}_\infty(\mathcal{F}, \gamma/2^k),$$

so that:

$$\inf_{f \in \mathcal{F}} \sum_{t=1}^T (y_t - f(x_t))^2 = \inf_{f \in \mathcal{F}} \sum_{t=1}^T \left( y_t - \pi_0(f)(x_t) - \underbrace{\sum_{k=0}^{\infty} [\pi_{k+1}(f) - \pi_k(f)](x_t)}_{|\text{small increments}| \leq 3\gamma/2^{k+1}} \right)^2.$$

**Our algorithm relies on a two-scale aggregation:**

- **high-scale aggregation:** run an Exponentially Weighted Average forecaster to be competitive against every function  $\pi_0(f)$  in the coarsest set  $\mathcal{F}^{(0)}$ ;

$$\text{Small Card } \mathcal{F}^{(0)} \implies \text{small regret}$$

- **low-scale aggregation:** run many instances of (an extension of) Exponentiated Gradient so as to be competitive against all increments  $\pi_{k+1}(f) - \pi_k(f)$ .

$$\text{Small } \|\pi_{k+1}(f) - \pi_k(f)\|_\infty \implies \text{small regret}$$

The Multi-variable Exponentiated Gradient algorithm (defined below) achieves this guarantee at all small scales  $\gamma/2^k$ ,  $k \geq 0$ , simultaneously.

## Key subroutine: Multi-variable Exponentiated Gradient

Let  $\Delta_N$  denote the simplex in  $\mathbb{R}^N$ . We design an extension of the Exponentiated Gradient algorithm to minimize a sequence of multi-variable loss functions  $(\mathbf{u}^{(1)}, \dots, \mathbf{u}^{(K)}) \mapsto \ell_t(\mathbf{u}^{(1)}, \dots, \mathbf{u}^{(K)})$  simultaneously over all the variables  $(\mathbf{u}^{(1)}, \dots, \mathbf{u}^{(K)}) \in \Delta_{N_1} \times \dots \times \Delta_{N_K}$ .

**input** : tuning parameters  $\eta^{(1)}, \dots, \eta^{(K)} > 0$ .

**initialization:** set  $\hat{\mathbf{u}}_1^{(k)} \stackrel{\text{def}}{=} (\frac{1}{N_k}, \dots, \frac{1}{N_k}) \in \Delta_{N_k}$  for all  $k = 1, \dots, K$ .

**for** each round  $t = 2, 3, \dots$  **do**

Compute the weight vectors  $(\hat{\mathbf{u}}_t^{(1)}, \dots, \hat{\mathbf{u}}_t^{(K)}) \in \Delta_{N_1} \times \dots \times \Delta_{N_K}$  as follows ( $Z_t^{(k)}$  is a normalization factor):

$$\hat{\mathbf{u}}_{t,i}^{(k)} \stackrel{\text{def}}{=} \frac{\exp\left(-\eta^{(k)} \sum_{s=1}^{t-1} \partial_{\hat{\mathbf{u}}_{s,i}^{(k)}} \ell_s(\hat{\mathbf{u}}_s^{(1)}, \dots, \hat{\mathbf{u}}_s^{(K)})\right)}{Z_t^{(k)}}, \quad i \in \{1, \dots, N_k\}.$$

**end**

**Regret bound:** Assume that the  $\ell_t$  are jointly convex and differentiable with bound  $\|\nabla_{\mathbf{u}^{(k)}} \ell_t\|_\infty \leq G^{(k)}$ . Then, the Multi-variable Exponentiated Gradient algorithm tuned with the parameters  $\eta^{(k)} = \sqrt{2 \log(N_k)/T} / G^{(k)}$  satisfies:

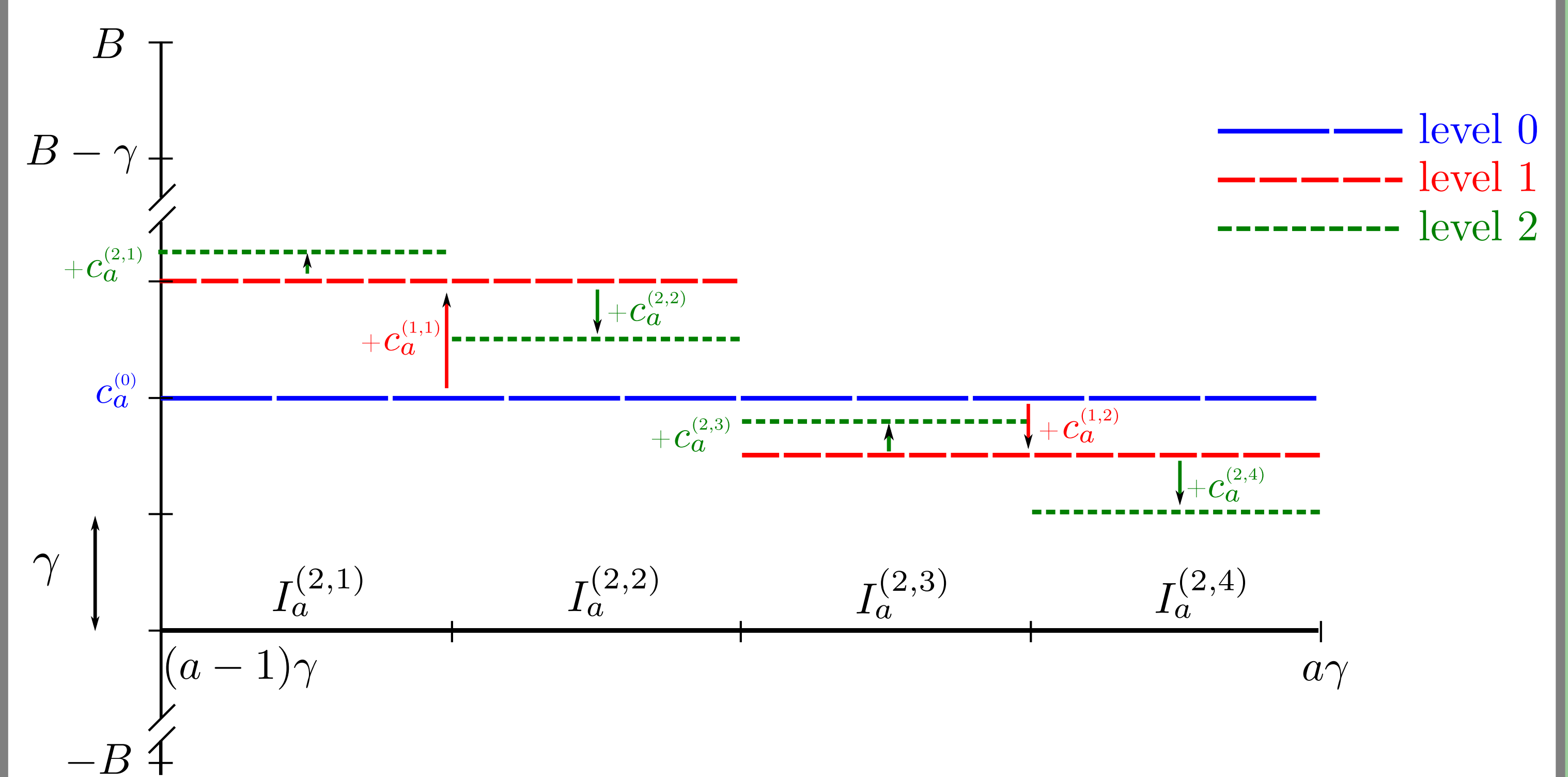
$$\sum_{t=1}^T \ell_t(\hat{\mathbf{u}}_t^{(1)}, \dots, \hat{\mathbf{u}}_t^{(K)}) - \min_{\mathbf{u}^{(1)}, \dots, \mathbf{u}^{(K)}} \sum_{t=1}^T \ell_t(\mathbf{u}^{(1)}, \dots, \mathbf{u}^{(K)}) \leq \sqrt{2T} \sum_{k=1}^K G^{(k)} \sqrt{\log N_k},$$

where the minimum is taken over all  $(\mathbf{u}^{(1)}, \dots, \mathbf{u}^{(K)}) \in \Delta_{N_1} \times \dots \times \Delta_{N_K}$ .

## An efficient version for Hölder classes

The idea is to design **computationally manageable coverings**  $\mathcal{F}^{(k)}$ ,  $k \geq 0$ :

- approximate any Lipschitz function  $f \in [0, 1] \rightarrow [-B, B]$  with **piecewise constant** functions (level  $k = 0$ );
- refine the approximation via a **dyadic discretization** (levels  $k \geq 1$ ).



At each round  $t$ , the point  $x_t$  falls into only one subinterval for each level  $k$   $\implies$  No need to update all coefficients  $\implies$  **manageable complexity**.

**For Hölder functions:** replace piecewise constant functions with **piecewise polynomials**.

Function class	Time complexity	Space complexity
Lipschitz on $[0, 1]$	$T^{4/3} \log T$	$T^{4/3} \log T$
$\beta$ -Hölder on $[0, 1]$	$\text{poly}(T)$	$\text{poly}(T)$

## References

- [1] N. Cesa-Bianchi and G. Lugosi. On prediction of individual sequences. *Ann. Statist.*, 27:1865–1895, 1999.
- [2] N. Cesa-Bianchi and G. Lugosi. Worst-case bounds for the logarithmic loss of predictors. *Mach. Learn.*, 43:247–264, 2001.
- [3] A. Rakhlin and K. Sridharan. Online nonparametric regression. In *Proceedings of COLT'14*, volume 35, pages 1232–1264. JMLR W&CP, 2014.