

# A chaining algorithm for online nonparametric regression

Sébastien Gerchinovitz

Institut de Mathématiques de Toulouse, Université Toulouse III – Paul Sabatier

This is a joint work with Pierre Gaillard.

# Introduction

- This talk is about regret minimization in a specific online learning problem known as **online nonparametric regression**.
- We consider the **adversarial** setting (almost no assumptions on the data sequence).
- We present an algorithm based on the **chaining** technique.

## Outline of the talk:

- 1 The chaining technique in the stochastic setting
- 2 Our setting: online regression with individual sequences
- 3 Large (nonparametric) function sets
- 4 An algorithm based on the chaining technique

1 The chaining technique in the stochastic setting

2 Online regression with individual sequences

3 Large (nonparametric) function sets

4 An algorithm based on the chaining technique

# Bounding the expected supremum of a stochastic process

Technique introduced by Dudley (1967). Let  $(X_f)_{f \in \mathcal{F}}$  be a centered stochastic process (indexed by a finite metric space  $(\mathcal{F}, d)$ ) with subgaussian increments:

$$\forall f, g \in \mathcal{F}, \quad \forall \lambda > 0, \quad \log \mathbb{E} e^{\lambda(X_f - X_g)} \leq \frac{\lambda^2}{2} d(f, g)^2.$$

**Goal:** upper bound the quantity  $\mathbb{E}[\sup_{f \in \mathcal{F}} X_f] = \mathbb{E}[\sup_{f \in \mathcal{F}} (X_f - X_{f_0})]$  for any  $f_0 \in \mathcal{F}$ .

Lemma (see, e.g., Boucheron et al. 2013)

Let  $Z_1, \dots, Z_N$  be such that  $\log \mathbb{E} e^{\lambda Z_i} \leq \lambda^2 v / 2$  for all  $\lambda > 0$  and  $i \in [N]$ .  
Then,  $\mathbb{E} \max_{i=1, \dots, N} Z_i \leq \sqrt{2v \log N}$ .

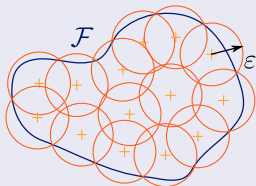
This lemma entails the pessimistic bound (**correlations are not used**):  
 $\mathbb{E}[\sup_{f \in \mathcal{F}} (X_f - X_{f_0})] \leq B \sqrt{2 \log(\text{card } \mathcal{F})}$  with  $B = \sup_{f \in \mathcal{F}} d(f, f_0)$ .

# Discretizing the space $(\mathcal{F}, d)$ into small balls

## Definition (metric entropy)

- Let  $(\mathcal{F}, d)$  be a metric space of finite cardinality.
- $\varepsilon$ -net: any subset  $\mathcal{G} \subseteq \mathcal{F}$  such that

$$\forall f \in \mathcal{F}, \exists g \in \mathcal{G} : d(f, g) \leq \varepsilon \iff \bigcup_{g \in \mathcal{G}} \bar{B}(g, \varepsilon) = \mathcal{F}$$



- $\mathcal{N}_d(\mathcal{F}, \varepsilon)$ : smallest cardinality of an  $\varepsilon$ -net.
- **metric entropy of  $\mathcal{F}$  at scale  $\varepsilon$** :  $\log \mathcal{N}_d(\mathcal{F}, \varepsilon)$ .  
It measures the complexity (richness) of the space  $(\mathcal{F}, d)$ .

# Multi-scale discretization to exploit the correlations

## Successive refining discretizations:

Let  $\mathcal{F}^{(0)} = \{f_0\}$ ,  $\mathcal{F}^{(1)}, \dots, \mathcal{F}^{(K-1)}$ ,  $\mathcal{F}^{(K)} = \mathcal{F}$  be minimal  $B/2^k$ -nets of  $\mathcal{F}$ :

$$\forall f \in \mathcal{F}, \exists \pi_k(f) \in \mathcal{F}^{(k)}, d(f, \pi_k(f)) \leq B/2^k .$$

**Chaining argument:** using the lemma at multiple scales, we get:

$$\begin{aligned} \mathbb{E} \left[ \sup_{f \in \mathcal{F}} (X_f - X_{f_0}) \right] &= \mathbb{E} \left[ \sup_{f \in \mathcal{F}} \sum_{k=1}^K \left( X_{\pi_k(f)} - X_{\pi_{k-1}(f)} \right) \right] \\ &\leq \sum_{k=1}^K \mathbb{E} \left[ \sup_{f \in \mathcal{F}} \underbrace{\left( X_{\pi_k(f)} - X_{\pi_{k-1}(f)} \right)}_{\text{small increments}} \right] \\ &\leq 6 \sum_{k=1}^K B 2^{-k} \sqrt{\log \mathcal{N}_d(\mathcal{F}, B/2^k)} \\ &\leq 12 \int_0^{B/2} \sqrt{\log \mathcal{N}_d(\mathcal{F}, \varepsilon)} d\varepsilon . \end{aligned}$$

# Multi-scale discretization to exploit the correlations

## Successive refining discretizations:

Let  $\mathcal{F}^{(0)} = \{f_0\}$ ,  $\mathcal{F}^{(1)}, \dots, \mathcal{F}^{(K-1)}$ ,  $\mathcal{F}^{(K)} = \mathcal{F}$  be minimal  $B/2^k$ -nets of  $\mathcal{F}$ :

$$\forall f \in \mathcal{F}, \exists \pi_k(f) \in \mathcal{F}^{(k)}, d(f, \pi_k(f)) \leq B/2^k .$$

**Chaining argument:** using the lemma at multiple scales, we get:

$$\begin{aligned} \mathbb{E} \left[ \sup_{f \in \mathcal{F}} (X_f - X_{f_0}) \right] &= \mathbb{E} \left[ \sup_{f \in \mathcal{F}} \sum_{k=1}^K (X_{\pi_k(f)} - X_{\pi_{k-1}(f)}) \right] \\ &\leq \sum_{k=1}^K \mathbb{E} \left[ \sup_{f \in \mathcal{F}} \underbrace{(X_{\pi_k(f)} - X_{\pi_{k-1}(f)})}_{\text{small increments}} \right] \\ &\leq 6 \sum_{k=1}^K B 2^{-k} \sqrt{\log \mathcal{N}_d(\mathcal{F}, B/2^k)} \\ &\leq 12 \underbrace{\int_0^{B/2} \sqrt{\log \mathcal{N}_d(\mathcal{F}, \varepsilon)} d\varepsilon}_{\text{Dudley's entropy integral}} . \end{aligned}$$

- 1 The chaining technique in the stochastic setting
- 2 Online regression with individual sequences
- 3 Large (nonparametric) function sets
- 4 An algorithm based on the chaining technique



# Setting: online regression with individual sequences

**Prediction task:** at each time  $t \in \mathbb{N}^*$ , predict the observation  $y_t \in \mathbb{R}$  from the input  $x_t \in \mathcal{X}$ , on the basis of the past data  $(x_1, y_1), \dots, (x_{t-1}, y_{t-1})$ .

**Initial step:** the environment chooses **arbitrary deterministic sequences**  $(y_t)_{t \geq 1}$  in  $\mathbb{R}$  and  $(x_t)_{t \geq 1}$  in  $\mathcal{X}$  but the forecaster has not access to them.

**At each time round**  $t \in \mathbb{N}^*$ ,

- 1 The environment reveals the input  $x_t \in \mathcal{X}$ .
- 2 The forecaster chooses a prediction  $\hat{y}_t \in \mathbb{R}$ .
- 3 The environment reveals the observation  $y_t \in \mathbb{R}$  and the forecaster incurs the loss  $(y_t - \hat{y}_t)^2$ .

# Goal: minimizing regret

Let  $\mathcal{F} \subseteq \mathbb{R}^{\mathcal{X}}$  be a set of functions.

**Goal of the forecaster:** on the long run, to predict almost as well as the best function  $f \in \mathcal{F}$  in hindsight, that is, to minimize the **regret**:

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

**Individual sequence setting:** our goal is to minimize the regret  $\text{Reg}_T(\mathcal{F})$  **uniformly** over all sequences  $(y_t)_{t \geq 1}$  in  $[-B, B]$  and  $(x_t)_{t \geq 1}$  in  $\mathcal{X}$ ; typically:

$$\sup_{\substack{|y_t| \leq B \\ x_t \in \mathcal{X}}} \left\{ \frac{1}{T} \sum_{t=1}^T (y_t - \hat{y}_t)^2 - \inf_{f \in \mathcal{F}} \frac{1}{T} \sum_{t=1}^T (y_t - f(x_t))^2 \right\} \leq o(1) \quad \text{when } T \rightarrow +\infty .$$

## Particular case: finite $\mathcal{F}$

Assume that  $\mathcal{F} = \{f_1, f_2, \dots, f_N\} \subseteq \mathbb{R}^X$  is **finite**. We can use a well-known algorithm studied, e.g., by Kivinen and Warmuth (1999) and Vovk (2001):

### Algorithm (Exponentially Weighted Average forecaster (EWA))

Parameter:  $\eta > 0$

At each round  $t \geq 1$ ,

- Using past data, compute the weight vector  $\hat{\mathbf{w}}_t = (\hat{w}_{t,1}, \dots, \hat{w}_{t,N})$  as

$$\hat{w}_{t,j} \triangleq \frac{\exp\left(-\eta \sum_{s=1}^{t-1} (y_s - f_j(x_s))^2\right)}{\sum_{j'=1}^N \exp\left(-\eta \sum_{s=1}^{t-1} (y_s - f_{j'}(x_s))^2\right)}, \quad 1 \leq j \leq N;$$

- Compute the convex combination (**convex aggregate**):

$$\hat{y}_t \triangleq \sum_{j=1}^N \hat{w}_{t,j} f_j(x_t).$$

## Regret guarantee when $\mathcal{F}$ is finite

If  $\mathcal{F}$  contains  $N$  functions, then we have a  $\mathcal{O}(\log N)$  upper bound on the regret under the boundedness assumption:

$$|y_1|, \dots, |y_T| \leq B \quad \text{and} \quad \|f_1\|_\infty, \dots, \|f_N\|_\infty \leq B .$$

**Theorem (Kivinen and Warmuth 1999)**

Assume that  $\mathcal{F} = \{f_1, f_2, \dots, f_N\} \subseteq [-B, B]^{\mathcal{X}}$ .

Then, the EWA algorithm tuned with  $\eta = 1/(8B^2)$  satisfies: for all sequences  $(y_t)_{t \geq 1}$  in  $[-B, B]$  and  $(x_t)_{t \geq 1}$  in  $\mathcal{X}$ , for all  $T \geq 1$ ,

$$\sum_{t=1}^T (y_t - \hat{y}_t)^2 - \min_{1 \leq j \leq N} \sum_{t=1}^T (y_t - f_j(x_t))^2 \leq 8B^2 \log N .$$

Remark 1: the requirement  $\forall j, \|f_j\|_\infty \leq B$  can be removed via clipping.

## Regret guarantee when $\mathcal{F}$ is finite

If  $\mathcal{F}$  contains  $N$  functions, then we have a  $\mathcal{O}(\log N)$  upper bound on the regret under the boundedness assumption:

$$|y_1|, \dots, |y_T| \leq B \quad \text{and} \quad \|f_1\|_\infty, \dots, \|f_N\|_\infty \leq B .$$

**Theorem (Kivinen and Warmuth 1999)**

Assume that  $\mathcal{F} = \{f_1, f_2, \dots, f_N\} \subseteq [-B, B]^{\mathcal{X}}$ .

Then, the EWA algorithm tuned with  $\eta = 1/(8B^2)$  satisfies: for all sequences  $(y_t)_{t \geq 1}$  in  $[-B, B]$  and  $(x_t)_{t \geq 1}$  in  $\mathcal{X}$ , for all  $T \geq 1$ ,

$$\sum_{t=1}^T (y_t - \hat{y}_t)^2 - \min_{1 \leq j \leq N} \sum_{t=1}^T (y_t - f_j(x_t))^2 \leq 8B^2 \log N .$$

Remark 1: the requirement  $\forall j, \|f_j\|_\infty \leq B$  can be removed via clipping.

Remark 2: we can obtain a similar bound if  $B = \max_{1 \leq t \leq T} |y_t|$  is **unknown**.

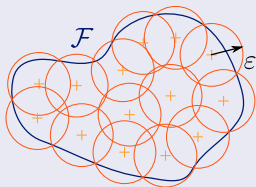
- 1 The chaining technique in the stochastic setting
- 2 Online regression with individual sequences
- 3 Large (nonparametric) function sets**
- 4 An algorithm based on the chaining technique

# Large function sets $\mathcal{F}$ : finite approximation

## Definition (metric entropy for sup norm)

- Let  $\mathcal{F} \subseteq \mathbb{R}^{\mathcal{X}}$  be a set of bounded functions endowed with the sup norm  $\|f\|_{\infty} \triangleq \sup_{x \in \mathcal{X}} |f(x)|$ .
- $\varepsilon$ -net: any subset  $\mathcal{G} \subseteq \mathcal{F}$  such that

$$\forall f \in \mathcal{F}, \exists g \in \mathcal{G} : \|f - g\|_{\infty} \leq \varepsilon .$$



- $\mathcal{N}_{\infty}(\mathcal{F}, \varepsilon)$ : smallest cardinality of an  $\varepsilon$ -net.
- metric entropy of  $\mathcal{F}$  at scale  $\varepsilon$ :  $\log \mathcal{N}_{\infty}(\mathcal{F}, \varepsilon)$ .

## Large function sets $\mathcal{F}$ : finite approximation (2)

Assume that  $\mathcal{F}$  is infinite (the EWA algorithm cannot be used). Small regret is still achievable if  $\mathcal{F}$  can be well approximated by a finite set.

**Discretizing**  $\mathcal{F}$  (Vovk, 2006): approximate  $\mathcal{F}$  with a minimal  $\varepsilon$ -net and run the EWA algorithm on this finite subset:

$$\begin{aligned} \sum_{t=1}^T (y_t - \hat{y}_t)^2 &\leq \min_{1 \leq j \leq \mathcal{N}_\infty(\mathcal{F}, \varepsilon)} \sum_{t=1}^T (y_t - f_j(x_t))^2 + 8B^2 \log \mathcal{N}_\infty(\mathcal{F}, \varepsilon) \\ &\leq \inf_{f \in \mathcal{F}} \sum_{t=1}^T (y_t - f(x_t))^2 + T\varepsilon^2 + 4TB\varepsilon + 8B^2 \log \mathcal{N}_\infty(\mathcal{F}, \varepsilon) \end{aligned}$$

**Finite-dimensional case:** given  $\varphi_j : \mathcal{X} \rightarrow [-B, B]$  and a compact set  $\Theta \subseteq \mathbb{R}^d$ , define

$$\mathcal{F} = \left\{ \sum_{j=1}^d \theta_j \varphi_j : \theta \in \Theta \right\} \subseteq \mathbb{R}^{\mathcal{X}}.$$

Note that  $\mathcal{N}_\infty(\mathcal{F}, \varepsilon) \lesssim (1/\varepsilon)^d$ . Choosing  $\varepsilon \approx 1/T$  yields a regret at most of the order of  $d \log(T)$ , which is optimal (**parametric** rate).



# What if $\mathcal{F}$ is very large (nonparametric)?

**Nonparametric set:** assume that  $\mathcal{F}$  is much larger than in the finite-dimensional case:

$$\log \mathcal{N}_\infty(\mathcal{F}, \varepsilon) \approx (1/\varepsilon)^p \quad \text{as} \quad \varepsilon \rightarrow 0 .$$

**Example:** Hölder class  $\mathcal{F} \subseteq \mathbb{R}^{[0,1]}$  of regularity  $\beta = q + \alpha$ :

$$|f^{(q)}(x) - f^{(q)}(y)| \leq \lambda |x - y|^\alpha \quad \text{and} \quad \forall k \leq q, \|f^{(k)}\|_\infty \leq B$$

In this case,  $\log \mathcal{N}_\infty(\mathcal{F}, \varepsilon) \approx \varepsilon^{-1/\beta}$  so that  $p = 1/\beta$ .

**EWA is suboptimal:** the regret bound  $T\varepsilon^2 + 4TB\varepsilon + 8B^2 \log \mathcal{N}_\infty(\mathcal{F}, \varepsilon)$  becomes roughly  $T\varepsilon + (1/\varepsilon)^p$ . Optimizing in  $\varepsilon$  only yields:

$$\sum_{t=1}^T (y_t - \hat{y}_t)^2 \leq \inf_{f \in \mathcal{F}} \sum_{t=1}^T (y_t - f(x_t))^2 + \mathcal{O}(T^{p/(p+1)}) ,$$

which is worse than the optimal rate  $\mathcal{O}(T^{p/(p+2)})$ .

# Optimal rates by Rakhlin and Sridharan (2014)

We still assume that  $\log \mathcal{N}_\infty(\mathcal{F}, \varepsilon) \approx (1/\varepsilon)^p$  as  $\varepsilon \rightarrow 0$ .

**Optimal regret:** through a **non-constructive approach** (reduction to a stochastic problem via von Neumann minimax theorem), Rakhlin and Sridharan (2014) proved that, if  $p \in (0, 2)$ , then

$$\begin{aligned} \text{Reg}_T(\mathcal{F}) &\leq c_1 B^2 (1 + \log \mathcal{N}_\infty(\mathcal{F}, \gamma)) + c_2 B \sqrt{T} \int_0^\gamma \sqrt{\log \mathcal{N}_\infty(\mathcal{F}, \varepsilon)} d\varepsilon \\ &\lesssim \gamma^{-p} + \sqrt{T} \int_0^\gamma \varepsilon^{-p/2} d\varepsilon \\ &\lesssim T^{p/(p+2)} \quad \text{for } \gamma = T^{-1/(p+2)}. \end{aligned}$$

The rate  $T^{p/(p+2)}$  is better than  $T^{p/(p+1)}$  obtained previously with EWA, and it is (in a sense) **optimal**.

# Optimal rates by Rakhlin and Sridharan (2014)

We still assume that  $\log \mathcal{N}_\infty(\mathcal{F}, \varepsilon) \approx (1/\varepsilon)^p$  as  $\varepsilon \rightarrow 0$ .

**Optimal regret:** through a non-constructive approach (reduction to a stochastic problem via von Neumann minimax theorem), Rakhlin and Sridharan (2014) proved that, if  $p \in (0, 2)$ , then

$$\begin{aligned} \text{Reg}_T(\mathcal{F}) &\leq c_1 B^2 (1 + \log \mathcal{N}_\infty(\mathcal{F}, \gamma)) + c_2 B \sqrt{T} \int_0^\gamma \sqrt{\log \mathcal{N}_\infty(\mathcal{F}, \varepsilon)} d\varepsilon \\ &\lesssim \gamma^{-p} + \sqrt{T} \int_0^\gamma \varepsilon^{-p/2} d\varepsilon \\ &\lesssim T^{p/(p+2)} \quad \text{for } \gamma = T^{-1/(p+2)}. \end{aligned}$$

**Example (Hölder class with regularity  $\beta$ ):**

Since  $p = 1/\beta$ , we get  $\text{Reg}_T(\mathcal{F})/T = \mathcal{O}(T^{-2\beta/(2\beta+1)})$  if  $\beta > 1/2$ .

Therefore, same rate as in the statistical setting (for  $\beta > 1/2$ ).

# Optimal rates by Rakhlin and Sridharan (2014)

We still assume that  $\log \mathcal{N}_\infty(\mathcal{F}, \varepsilon) \approx (1/\varepsilon)^p$  as  $\varepsilon \rightarrow 0$ .

**Optimal regret:** through a non-constructive approach (reduction to a stochastic problem via von Neumann minimax theorem), Rakhlin and Sridharan (2014) proved that, if  $p \in (0, 2)$ , then

$$\begin{aligned} \text{Reg}_T(\mathcal{F}) &\leq c_1 B^2 (1 + \log \mathcal{N}_\infty(\mathcal{F}, \gamma)) + c_2 B \sqrt{T} \int_0^\gamma \sqrt{\log \mathcal{N}_\infty(\mathcal{F}, \varepsilon)} d\varepsilon \\ &\lesssim \gamma^{-p} + \sqrt{T} \int_0^\gamma \varepsilon^{-p/2} d\varepsilon \\ &\lesssim T^{p/(p+2)} \quad \text{for } \gamma = T^{-1/(p+2)}. \end{aligned}$$

The above integral is a **Dudley entropy integral**.

- In statistical learning with i.i.d. data, useful to derive risk bounds for empirical risk minimizers (e.g., Massart 2007; Rakhlin et al. 2013).
- Also appears in online learning with individual sequences. Earlier appearances: Opper and Haussler (1997); Cesa-Bianchi and Lugosi (1999, 2001).

# Our contributions

- 1 We provide an **explicit algorithm** that achieves the Dudley-type regret bound (when  $p \in (0, 2)$ ):

$$\text{Reg}_T(\mathcal{F}) \leq c_1 B^2 (1 + \log \mathcal{N}_\infty(\mathcal{F}, \gamma)) + c_2 B \sqrt{T} \int_0^\gamma \sqrt{\log \mathcal{N}_\infty(\mathcal{F}, \varepsilon)} d\varepsilon.$$

Nota: contrary to Rakhlin and Sridharan (2014), our bounds are not in terms of the stronger notion of *sequential entropy*.

- 2 This algorithm uses ideas from the **chaining technique**, and relies on a new subroutine (Multi-variable Exponentiated Gradient algorithm) to perform optimization at different scales simultaneously.
- 3 We address computational issues by showing how to construct **more efficient** and quasi-optimal  **$\varepsilon$ -nets** (for Hölder classes).

- 1 The chaining technique in the stochastic setting
- 2 Online regression with individual sequences
- 3 Large (nonparametric) function sets
- 4 An algorithm based on the chaining technique**

# Linearizing the square loss can help locally (1)

Suppose we play with loss functions  $\mathbf{u} \mapsto \ell_t(\mathbf{u})$ ,  $t \geq 1$ , that are convex and differentiable over the simplex  $\Delta_N = \{\mathbf{u} \in \mathbb{R}_+^N : \sum_{i=1}^N u_i = 1\}$ .

## Algorithm (Exponentiated Gradient—EG)

*Parameter:*  $\eta > 0$

*At each round*  $t \geq 1$ , *compute the weight vector*  $\hat{\mathbf{u}}_t \in \Delta_N$  *by*

$$\hat{u}_{t,j} \triangleq \frac{1}{Z_t} \exp\left(-\eta \sum_{s=1}^{t-1} \partial_{\hat{u}_{s,j}} \ell_s(\hat{\mathbf{u}}_s)\right), \quad 1 \leq j \leq N.$$

## Theorem (Kivinen and Warmuth 1999 and Cesa-Bianchi 1999)

*Assume*  $\ell_t$  *convex, diff, and*  $\|\nabla \ell_t\|_\infty \leq G$ . *For*  $\eta = G^{-1} \sqrt{2 \log(N)/T}$ ,

$$\sum_{t=1}^T \ell_t(\hat{\mathbf{u}}_t) \leq \min_{\mathbf{u} \in \Delta_N} \sum_{t=1}^T \ell_t(\mathbf{u}) + G \sqrt{2T \log N}.$$

# Linearizing the square loss can help locally (1)

Suppose we play with loss functions  $\mathbf{u} \mapsto \ell_t(\mathbf{u})$ ,  $t \geq 1$ , that are convex and differentiable over the simplex  $\Delta_N = \{\mathbf{u} \in \mathbb{R}_+^N : \sum_{i=1}^N u_i = 1\}$ .

## Algorithm (Exponentiated Gradient—EG)

Parameter:  $\eta > 0$

At each round  $t \geq 1$ , compute the weight vector  $\hat{\mathbf{u}}_t \in \Delta_N$  by

$$\hat{u}_{t,j} \triangleq \frac{1}{Z_t} \exp\left(-\eta \sum_{s=1}^{t-1} \partial_{\hat{u}_{s,j}} \ell_s(\hat{\mathbf{u}}_s)\right), \quad 1 \leq j \leq N.$$

## Theorem (Kivinen and Warmuth 1999 and Cesa-Bianchi 1999)

Assume  $\ell_t$  convex, diff, and  $\|\nabla \ell_t\|_\infty \leq G$ . For  $\eta = G^{-1} \sqrt{2 \log(N)/T}$ ,

$$\sum_{t=1}^T \ell_t(\hat{\mathbf{u}}_t) \leq \min_{\mathbf{u} \in \Delta_N} \sum_{t=1}^T \ell_t(\mathbf{u}) + G \sqrt{2T \log N}.$$



## Linearizing the square loss can help locally (2)

**Application:** we want to predict almost as well as the best function in  $\mathcal{F} = \{f_0 + g_j : j = 1, \dots, N\}$  with  $\|g_j\|_\infty$  small (neighbors of  $f_0$ ).

We use EG with  $\ell_t(\mathbf{u}) = \left(y_t - f_0(x_t) - \sum_{j=1}^N u_j g_j(x_t)\right)^2$ ,  $\mathbf{u} \in \Delta_N$ .

Since  $\|\nabla \ell_t\|_\infty \lesssim B \max_j \|g_j\|_\infty$ , the EG algorithm satisfies:

$$\sum_{t=1}^T \underbrace{\left(y_t - f_0(x_t) - \sum_{j=1}^N \hat{u}_{t,j} g_j(x_t)\right)^2}_{= \hat{y}_t} \leq \min_{1 \leq j \leq N} \sum_{t=1}^T (y_t - f_0(x_t) - g_j(x_t))^2 + \square B \max_{1 \leq j \leq N} \|g_j\|_\infty \sqrt{T \log N}$$

**Advantage:** the above regret bound  $B \max_j \|g_j\|_\infty \sqrt{T \log N}$  improves on  $B^2 \log N$  (obtained by EWA) when  $\max_j \|g_j\|_\infty \ll B \sqrt{\log(N)/T}$ .

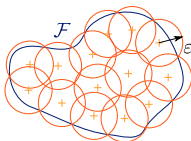
Thus, **linearizing the square loss** can help if the functions in  $\mathcal{F}$  are **close**.

# Turning the chaining technique into an online algorithm

We still assume that  $\log \mathcal{N}_\infty(\mathcal{F}, \varepsilon) \approx (1/\varepsilon)^p$  as  $\varepsilon \rightarrow 0$ . Recall that we want to prove a bound of the form:

$$\sum_{t=1}^T (y_t - \hat{y}_t)^2 \leq \inf_{f \in \mathcal{F}} \sum_{t=1}^T (y_t - f(x_t))^2 + [\text{small term}]$$

**Chaining principle:** as previously, we discretize  $\mathcal{F}$  and use projections  $\pi_k(f)$  such that  $\sup_f \|\pi_k(f) - f\|_\infty \leq \gamma/2^k$  for all  $k \geq 0$ .



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

# Aggregation at two different levels

$$\inf_{f \in \mathcal{F}} \sum_{t=1}^T (y_t - f(x_t))^2 = \inf_{f \in \mathcal{F}} \sum_{t=1}^T \left( y_t - \underbrace{\pi_0(f)}_{\in \mathcal{F}^{(0)}}(x_t) - \sum_{k=1}^{\infty} \underbrace{[\pi_k(f) - \pi_{k-1}(f)]}_{\in \mathcal{G}^{(k)}}(x_t) \right)^2$$

Sufficient goal:

$$\sum_{t=1}^T (y_t - \hat{y}_t)^2 \leq \inf_{f_0, g_1, \dots, g_K} \sum_{t=1}^T (y_t - (f_0 + g_1 + \dots + g_K)(x_t))^2 + [\text{small term}]$$

Two aggregation levels:

$$\left. \begin{array}{l} f_{0,1} \xrightarrow{\text{low scale multi-variable EG}} \hat{f}_{t,1} \\ f_{0,2} \xrightarrow{\quad \quad \quad \quad \quad \quad \quad} \hat{f}_{t,2} \\ \vdots \\ f_{0,N_0} \xrightarrow{\quad \quad \quad \quad \quad \quad \quad} \hat{f}_{t,N_0} \end{array} \right\} \xrightarrow{\text{high scale EWA}} \hat{y}_t = \sum_{j=1}^{N_0} \hat{w}_{t,j} \hat{f}_{t,j}(x_t)$$

# Combining two regret guarantees

**High-scale aggregation** Using an Exponentially Weighted Average (EWA) forecaster  $\hat{f}_t = \sum_{j=1}^{N_0} \hat{w}_{t,j} \hat{f}_{t,j}$  yields

$$\sum_{t=1}^T (y_t - \hat{y}_t)^2 \leq \min_{1 \leq j \leq N_0} \sum_{t=1}^T \left( y_t - \hat{f}_{t,j}(x_t) \right)^2 + \square B^2 \log N_0$$

**Low-scale aggregation** Recall that  $\mathcal{G}^{(k)} = \{ \pi_k(f) - \pi_{k-1}(f) : f \in \mathcal{F} \}$ . Denote  $\mathcal{G}^{(k)} = \{ g_1^{(k)}, \dots, g_{N_k}^{(k)} \}$ .

We designed a multi-variable extension of the Exponentiated Gradient algorithm:

$$\hat{f}_{t,j} \triangleq f_{0,j} + \sum_{k=1}^K \sum_{i=1}^{N_k} \hat{u}_{t,i}^{(j,k)} g_i^{(k)}$$

which yields, for all  $j = 1, \dots, N_0$ ,

$$\begin{aligned} \sum_{t=1}^T \left( y_t - \hat{f}_{t,j}(x_t) \right)^2 &\leq \min_{g_1, \dots, g_K} \sum_{t=1}^T \left( y_t - (f_{0,j} + g_1 + \dots + g_K)(x_t) \right)^2 \\ &\quad + 120B\sqrt{T} \int_0^{\gamma/2} \sqrt{\log \mathcal{N}_\infty(\mathcal{F}, \varepsilon)} d\varepsilon. \end{aligned}$$

# Main result

The next theorem indicates that the Chaining Exponentially Weighted Average forecaster satisfies a **Dudley-type regret bound**.

## Theorem (Gaillard and G., 2015)

Let  $B > 0$ ,  $T \geq 1$ , and  $\gamma \in (\frac{B}{T}, B)$ .

- Assume that  $\max_{1 \leq t \leq T} |y_t| \leq B$  and that  $\sup_{f \in \mathcal{F}} \|f\|_\infty \leq B$ .
- Assume that  $(\mathcal{F}, \|\cdot\|_\infty)$  is totally bounded and define  $\mathcal{F}^{(0)}$  and  $\mathcal{G}^{(k)}$  as above.

Then, the Chaining Exponentially Weighted Average forecaster (tuned with appropriate parameters) satisfies:

$$\text{Reg}_T(\mathcal{F}) \leq B^2(5 + 50 \log \mathcal{N}_\infty(\mathcal{F}, \gamma)) + 120B\sqrt{T} \int_0^{\gamma/2} \sqrt{\log \mathcal{N}_\infty(\mathcal{F}, \varepsilon)} d\varepsilon.$$

# Computational issues: dyadic discretization

We assume that  $\mathcal{F} = \{f : [0, 1] \rightarrow [-B, B] : f \text{ is 1-Lipschitz}\}$ .

## Regret bound:

We know that  $\log \mathcal{N}_\infty(\mathcal{F}, \varepsilon) = \mathcal{O}(\varepsilon^{-1})$ .

Therefore, our algorithm obtains  $\text{Reg}_T(\mathcal{F}) = \mathcal{O}(T^{1/3})$ , which is optimal.

## Computational issue:

Our algorithm updates  $\exp(\mathcal{O}(T))$  weights at every round  $t$ .

Hence **very poor time and space computational complexities**.

## Solution:

$\mathcal{F}$  has a sufficiently nice structure that can be exploited to construct **computationally manageable  $\varepsilon$ -nets** with quasi-optimal cardinality.

For example: piecewise-constant approximations on a dyadic discretization lead to  $\mathcal{O}(T^{1/3} \log T)$  regret and per-round time complexity.

# Conclusion

- We designed an **explicit algorithm** with a **Dudley-type** regret bound for online nonparametric regression.
- We provided an **efficient** implementation for **Hölder** classes.

Thank you for your attention!

---

## Appendix

---



5 Computational issues: dyadic discretization

# Lipschitz class $\mathcal{F}$ : a computationally efficient discretization

We assume that  $\mathcal{F} = \{f : [0, 1] \rightarrow [-B, B] : f \text{ is 1-Lipschitz}\}$ .

## Regret bound:

We know that  $\log \mathcal{N}_\infty(\mathcal{F}, \varepsilon) = \mathcal{O}(\varepsilon^{-1})$ .

Therefore, our algorithm obtains  $\text{Reg}_T(\mathcal{F}) = \mathcal{O}(T^{1/3})$ , which is optimal.

## Computational issue:

Our algorithm updates exponentially many weights at every round  $t$ .

Hence **poor time and space computational complexities**.

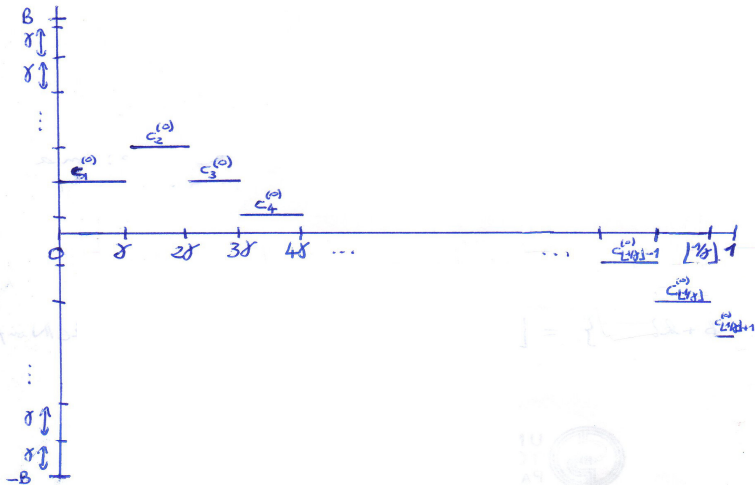
## Solution:

$\mathcal{F}$  has a sufficiently nice structure that can be exploited to construct **computationally manageable  $\varepsilon$ -nets** with quasi-optimal cardinality.

# High-level discretization (piecewise-constant approximation)

- Partition the  $x$ -axis  $[0, 1]$ :  $I_a \triangleq [(a-1)\gamma, a\gamma)$ ,  $a = 1, \dots, \frac{1}{\gamma}$ .
- Discretize the  $y$ -axis  $[-B, B]$ :  $\mathcal{C}^{(0)} = \{-B + j\gamma : j = 0, \dots, \frac{2B}{\gamma}\}$ .

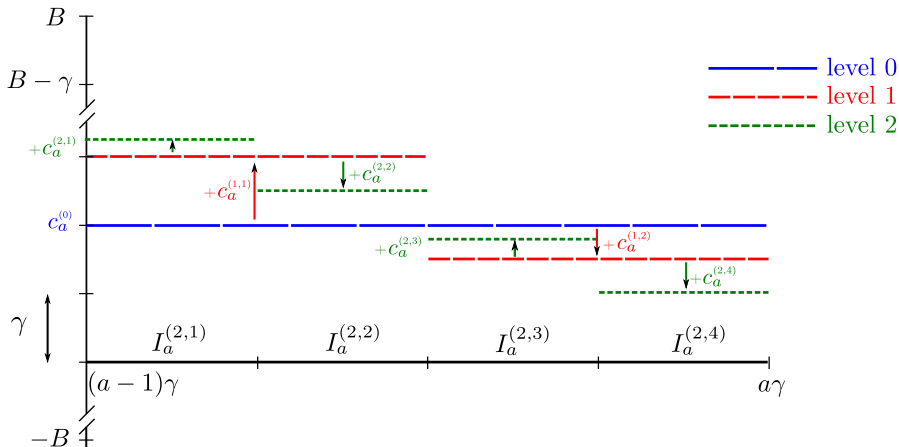
$\mathcal{F}^{(0)}$ : set of piecewise-constant functions  $f^{(0)}(x) = \sum_{a=1}^{1/\gamma} c_a^{(0)} \mathbb{I}_{x \in I_a}$ ,  $c_a^{(0)} \in \mathcal{C}^{(0)}$ .



# Low-level discretization (dyadic approximation)

$\mathcal{F}^{(M)}$ : set of all functions  $f_c : [0, 1] \rightarrow \mathbb{R}$  of the form

$$f_c(x) = \underbrace{\sum_{a=1}^{1/\gamma} c_a^{(0)} \mathbb{I}_{x \in I_a}}_{f^{(0)}(x)} + \sum_{m=1}^M \underbrace{\sum_{a=1}^{1/\gamma} \sum_{n=1}^{2^m} c_a^{(m,n)} \mathbb{I}_{x \in I_a^{(m,n)}}}_{g^{(m)}(x)} .$$



## Theorem (Gaillard and G., 2015)

Let  $B > 0$ ,  $T \geq 2$ , and  $\mathcal{F}$  be the set of all 1-Lipschitz functions from  $[0, 1]$  to  $[-B, B]$ . Assume that  $\max_{1 \leq t \leq T} |y_t| \leq B$ .

Then, the Dyadic Chaining Algorithm (see preprint) satisfies, for some absolute constant  $c > 0$ ,

$$\text{Reg}_T(\mathcal{F}) \leq c \max\{B, B^2\} T^{1/3} \log T .$$

Remark: additional log factor, but computationally **tractable**:

- per-round time complexity:  $\mathcal{O}(T^{1/3} \log T)$ ;
- space complexity:  $\mathcal{O}(T^{4/3} \log T)$ .

# Bibliographie I

- S. Boucheron, G. Lugosi, and P. Massart. *Concentration inequalities: a nonasymptotic theory of independence*. Oxford University Press, 2013.
- N. Cesa-Bianchi. Analysis of two gradient-based algorithms for on-line regression. *J. Comput. System Sci.*, 59(3):392–411, 1999.
- N. Cesa-Bianchi and G. Lugosi. On prediction of individual sequences. *Ann. Statist.*, 27: 1865–1895, 1999.
- N. Cesa-Bianchi and G. Lugosi. Worst-case bounds for the logarithmic loss of predictors. *Mach. Learn.*, 43:247–264, 2001.
- R.M. Dudley. The sizes of compact subsets of hilbert space and continuity of gaussian processes. *Journal of Functional Analysis*, 1(3):290 – 330, 1967.
- J. Kivinen and M. K. Warmuth. Averaging expert predictions. In *Proceedings of the 4th European Conference on Computational Learning Theory (EuroCOLT'99)*, pages 153–167, 1999.
- P. Massart. *Concentration Inequalities and Model Selection*, volume 1896 of *Lecture Notes in Mathematics*. Springer, Berlin, 2007.
- M. Opper and D. Haussler. Worst case prediction over sequences under log loss. In *The Mathematics of Information Coding, Extraction, and Distribution*. Spinger Verlag, 1997.
- A. Rakhlin and K. Sridharan. Online nonparametric regression. *JMLR W&CP*, 35 (Proceedings of COLT 2014):1232–1264, 2014.

- A. Rakhlin, K. Sridharan, and A.B. Tsybakov. Empirical entropy, minimax regret and minimax risk. *Bernoulli*, 2013. URL <http://arxiv.org/abs/1308.1147>. To appear.
- V. Vovk. Competitive on-line statistics. *Internat. Statist. Rev.*, 69:213–248, 2001.
- V. Vovk. Metric entropy in competitive on-line prediction. *arXiv*, 2006. URL <http://arxiv.org/abs/cs.LG/0609045>.