



# MINIMAX ADAPTIVE BOOSTING IN ONLINE NON-PARAMETRIC REGRESSION

Boosting & Optimisation Workshop, Besançon

---

**Paul Liautaud**

October 14, 2024

Sorbonne University, Paris

## Joint work with



Pierre Gaillard  
CR Inria/UGA



Olivier Wintenberger  
PR Sorbonne University

# Table of Contents

1. Online Learning & Non-Parametric Regression
2. Building Predictions with Online Gradient Boosting
3. Online Gradient Boosting in Chaining-Tree
4. Adaptive Boosting in Online NonParametric Regression

# Online Learning & Non-Parametric Regression

---

# Classical Machine Learning

The learner:



- 1 observes a **whole training dataset** with labels/targets:

$$(x_1, y_1), \dots, (x_T, y_T) \stackrel{\text{iid}}{\sim} (X, Y) \text{ with distribution } \mathbb{P} \text{ over } \mathcal{X} \times \mathcal{Y}.$$

# Classical Machine Learning

The learner:



→ Learning method

- 1 observes a **whole training dataset** with labels/targets:

$$(x_1, y_1), \dots, (x_T, y_T) \stackrel{\text{iid}}{\sim} (X, Y) \text{ with distribution } \mathbb{P} \text{ over } \mathcal{X} \times \mathcal{Y}.$$

- 2 learn a function  $\hat{f} : \mathcal{X} \rightarrow \mathcal{Y} \in \mathcal{F}$  with small risk  $\mathbb{E}_{\mathbb{P}}[\ell(\hat{f}(X), Y)]$  by minimizing:

$$R(\hat{f}) = \frac{1}{T} \sum_{t=1}^T \ell(\hat{f}(x_t), y_t),$$

where  $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$  is a prescribed loss function.

# Classical Machine Learning



- 1 observes a **whole training dataset** with labels/targets:

$$(x_1, y_1), \dots, (x_T, y_T) \stackrel{\text{iid}}{\sim} (X, Y) \text{ with distribution } \mathbb{P} \text{ over } \mathcal{X} \times \mathcal{Y}.$$

- 2 learn a function  $\hat{f} : \mathcal{X} \rightarrow \mathcal{Y} \in \mathcal{F}$  with small risk  $\mathbb{E}_{\mathbb{P}}[\ell(\hat{f}(X), Y)]$  by minimizing:

$$R(\hat{f}) = \frac{1}{T} \sum_{t=1}^T \ell(\hat{f}(x_t), y_t),$$

where  $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$  is a prescribed loss function.

- 3 controls the error of new data if they are similar to the training data.

# Classical Machine Learning

The learner:



→ Learning method → Prediction on test data

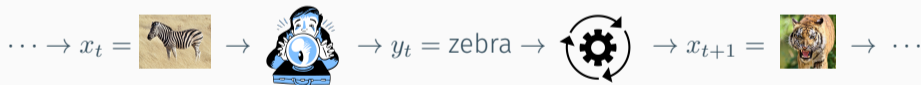
We won't deal with it!



# A dive into Sequential Learning

In **sequential learning**:

- Data are acquired and treated **on the fly**;
- Data are **not** necessarily **iid**, possibly **adversarial**;
- Feedbacks are received and algorithms updated **step by step**.



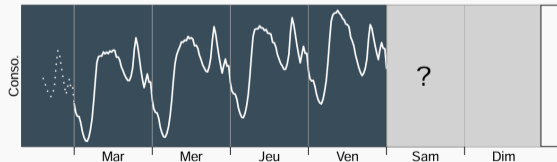
# A dive into Sequential Learning

## ? Why **Online** Learning?

In some applications, the environment may **evolve over time** and data may be available **sequentially**.

Examples:

- ads to display,
- electricity consumption forecast,
- spam detection,
- aggregation of expert knowledge.



## Setting of the talk

- Data arrives **sequentially** as a stream

$$(x_1, y_1), \dots, (x_{t-1}, y_{t-1}), (x_t, ??) \in \mathcal{X} \times \mathbb{R}$$

- At time  $t \geq 1$ : we want to **predict** each next response  $y_t$  as a function of  $x_t$

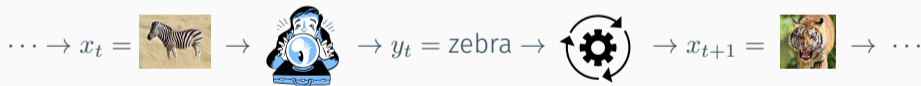
$$\hat{f}_t(x_t), \quad \text{with } \hat{f}_t \in \mathbb{R}^{\mathcal{X}} \text{ sequentially updated with } (x_s, y_s)_{1 \leq s \leq t-1}.$$

# Setting of the talk

The scenario is as follows:

At each round  $t = 1, \dots, T$ , the learner or algorithm

- 1 observes input  $x_t \in \mathcal{X}$
- 2 makes prediction  $\hat{f}_t(x_t) \in \mathbb{R}$
- 3 incurs loss  $l_t(\hat{f}_t(x_t))$  and discover gradients  $g_t$
- 4 updates prediction function  $\hat{f}_t \rightarrow \hat{f}_{t+1}$

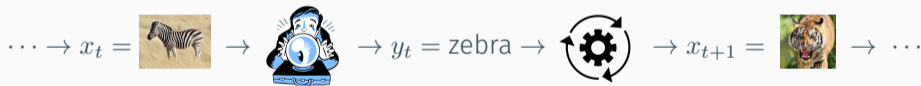


# Setting of the talk

The scenario is as follows:

At each round  $t = 1, \dots, T$ , the learner or algorithm

- 1 observes input  $x_t \in \mathcal{X}$
- 2 makes prediction  $\hat{f}_t(x_t) \in \mathbb{R}$  Choose  $\hat{f}_t$  before observing  $l_t$
- 3 incurs loss  $l_t(\hat{f}_t(x_t))$  and discover gradients  $g_t$  No assumptions on how  $l_t$  is generated!
- 4 updates prediction function  $\hat{f}_t \rightarrow \hat{f}_{t+1}$



- $l_1, \dots, l_T$  are convex, differentiable and  $G$ -Lipschitz, with  $G > 0$ ;
- $\mathcal{X} \subset \mathbb{R}^d$  bounded subset with  $|\mathcal{X}| = \sup_{x, x' \in \mathcal{X}} \|x - x'\|_\infty$ .

# Setting of the talk

**Goal:** find  $\hat{f}_1, \dots, \hat{f}_T$  that...  
minimize the cumulative loss

$$\sum_{t=1}^T \ell_t(\hat{f}_t(x_t))$$

# Setting of the talk

**Goal:** find  $\hat{f}_1, \dots, \hat{f}_T$  that...

minimize the cumulative loss  $\Leftrightarrow$  predict almost as well as the best function  $f$

$$\sum_{t=1}^T \ell_t(\hat{f}_t(x_t))$$

$$\underbrace{\sum_{t=1}^T \ell_t(\hat{f}_t(x_t)) - \sum_{t=1}^T \ell_t(f(x_t))}_{:= \text{Reg}_T(f)}$$

**"Non-Parametric regression"**  $\Leftrightarrow$  we compare  $(\hat{f}_t)$  to benchmark functions  $f \in \mathcal{F}$  (e.g., Lipschitz)

# Regret Analysis

→ We want  $\hat{f}_1, \dots, \hat{f}_T$  such that regret

$$\text{Reg}_T(f) = \underbrace{\sum_{t=1}^T \ell_t(\hat{f}_t(x_t))}_{\text{our performance}} - \underbrace{\sum_{t=1}^T \ell_t(f(x_t))}_{\text{reference performance}}$$

against  $f \in \mathcal{F}$  is ----- .



# Regret Analysis

→ We want  $\hat{f}_1, \dots, \hat{f}_T$  such that regret

$$\text{Reg}_T(f) = \underbrace{\sum_{t=1}^T \ell_t(\hat{f}_t(x_t))}_{\text{our performance}} - \underbrace{\sum_{t=1}^T \ell_t(f(x_t))}_{\text{reference performance}} = \underbrace{o(T)}_{\text{goal}}$$


against  $f \in \mathcal{F}$  is as **small** as possible.

# Regret Analysis

→ We want  $\hat{f}_1, \dots, \hat{f}_T$  such that regret

$$\text{Reg}_T(f) = \underbrace{\sum_{t=1}^T \ell_t(\hat{f}_t(x_t))}_{\text{our performance}} - \underbrace{\sum_{t=1}^T \ell_t(f(x_t))}_{\text{reference performance}} = \underbrace{o(T)}_{\text{goal}}$$

against  $f \in \mathcal{F}$  is as **small** as possible.

 **Difficulty: no stochastic assumption** on data  $(x_t, y_t, \ell_t)$ !

→  $(\hat{f}_t)$  have to perform **well** with all **arbitrary** time series i.e. approaching

$$\inf_{\hat{f}_1} \sup_{x_1, \ell_1} \inf_{\hat{f}_2} \sup_{x_2, \ell_2} \cdots \inf_{\hat{f}_T} \sup_{x_T, \ell_T} \sup_{f \in \mathcal{F}} \text{Reg}_T(f).$$

# Regret Analysis

→ We want  $\hat{f}_1, \dots, \hat{f}_T$  such that regret

$$\text{Reg}_T(f) = \underbrace{\sum_{t=1}^T \ell_t(\hat{f}_t(x_t))}_{\text{our performance}} - \underbrace{\sum_{t=1}^T \ell_t(f(x_t))}_{\text{reference performance}} = \underbrace{o(T)}_{\text{goal}}$$

against  $f \in \mathcal{F}$  is as **small** as possible.

 **Difficulty: no stochastic assumption** on data  $(x_t, y_t, \ell_t)$ !

→  $(\hat{f}_t)$  have to perform **well** with all **arbitrary** time series i.e. approaching

$$\inf_{\hat{f}_1} \sup_{x_1, \ell_1} \inf_{\hat{f}_2} \sup_{x_2, \ell_2} \cdots \inf_{\hat{f}_T} \sup_{x_T, \ell_T} \sup_{f \in \mathcal{F}} \text{Reg}_T(f).$$

? How to sequentially build predictors  $\hat{f}_1, \dots, \hat{f}_T$ ?

# Building Predictions with Online Gradient Boosting

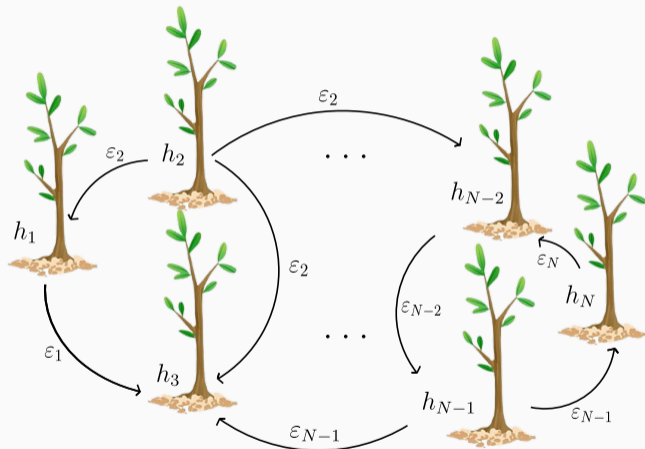
---

# Boosting uses "wisdom of the crowd"

- **Boosting**: ensemble method combining multiple **weak learners** to create a **strong learner**

- Each **model** corrects/learns from errors of its peers

→ Resulting in a **highly accurate** predictive model [1]



[1] e.g. AdaBoost and XGBoost

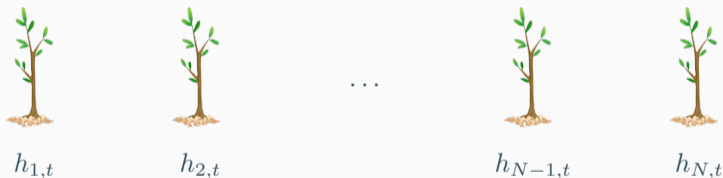
## How to deal with weak learners?

- $\mathcal{W} \subset \mathbb{R}^{\mathcal{X}}$  a set of real valued functions  $\mathcal{X} \rightarrow \mathbb{R}$ ;
- $\text{span}_N(\mathcal{W}) = \{\sum_{n=1}^N \beta_n h_n, h_n \in \mathcal{W}, \beta_n \in \mathbb{R}\}$  linear function space associated to  $\mathcal{W}$ .

## How to deal with weak learners?

- $\mathcal{W} \subset \mathbb{R}^{\mathcal{X}}$  a set of real valued functions  $\mathcal{X} \rightarrow \mathbb{R}$ ;
- $\text{span}_N(\mathcal{W}) = \{\sum_{n=1}^N \beta_n h_n, h_n \in \mathcal{W}, \beta_n \in \mathbb{R}\}$  linear function space associated to  $\mathcal{W}$ .

For each  $t = 1, \dots, T$ , we use  $N \geq 1$  *sequential* weak predictors from  $\mathcal{W}$



and we form *strong predictor* at any time  $t \geq 1$  as

$$\hat{f}_t = \sum_{n=1}^N \beta_{n,t} h_{n,t}, \quad \beta_{n,t} \in \mathbb{R}, n \in [N]$$

## A new Online Gradient Boosting procedure

→ **Goal:** We want to find a sequence of functions

$$\hat{f}_t = \sum_n \beta_{n,t} h_{n,t} \in \text{span}_N(\mathcal{W}), \quad 1 \leq t \leq T,$$

minimizing regret against  $\mathcal{F} = \text{span}_N(\mathcal{W})$ .



# A new Online Gradient Boosting procedure

→ **Goal:** We want to find a sequence of functions

$$\hat{f}_t = \sum_n \beta_{n,t} h_{n,t} \in \text{span}_N(\mathcal{W}), \quad 1 \leq t \leq T,$$

minimizing regret against  $\mathcal{F} = \text{span}_N(\mathcal{W})$ .

🍃 At  $t \geq 1$ , each  $n \in [N]$  is boosted with OGB as:

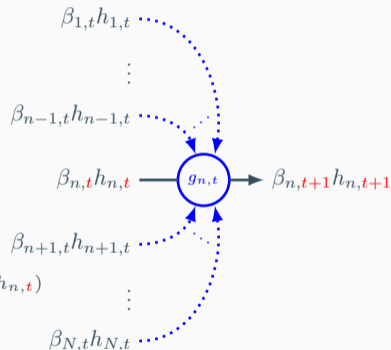
- 1 Predict  $\hat{f}_t(x_t)$
- 2  $(\beta_{n,t}, h_{n,t})$  receives its gradient

$$g_{n,t} = \left[ \nabla_{(\beta_n, h_n)} \ell_t(\hat{f}_{-n,t}(x_t) + \beta_n h_n(x_t)) \right]_{(\beta_n, h_n) = (\beta_{n,t}, h_{n,t})}$$

- 3 Find  $(\beta_{n,t+1}, h_{n,t+1}) \in \mathbb{R} \times \mathcal{W}$  to solve

$$\min_{\beta_n, h_n} \ell_t(\hat{f}_{-n,t}(x_t) + \beta_n h_n(x_t)) \quad (1)$$

using gradient  $g_{n,t}$ .



**Figure 1:** Boosting at time  $t$ .

# Online Gradient Boosting in Chaining-Tree

---

# Tree-based Method

Regular decision-tree  $(\mathcal{T}, \bar{\mathcal{X}}, \bar{W})$  over  $\mathcal{X}$  is made of:

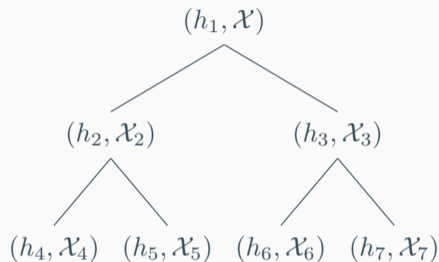
- a set of nodes  $\mathcal{N}(\mathcal{T})$  including leaves  $\mathcal{L}(\mathcal{T})$ ;
- a family of subregions

$$\bar{\mathcal{X}} = \{\mathcal{X}_n, n \in \mathcal{N}(\mathcal{T})\}$$

partitionning  $\mathcal{X}$  by level ;

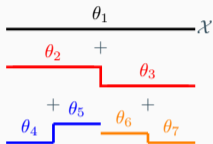
- a family of prediction functions

$$\bar{W} = \{h_n, n \in \mathcal{N}(\mathcal{T})\}.$$



**Figure 2:** Example of  $\mathcal{T}$  with depth  $d(\mathcal{T}) = 3$  over  $\mathcal{X} \subset \mathbb{R}$ .

# Chaining-Tree



**Figure 3:** Prediction of a CT  $\mathcal{T}$  of depth  $d(\mathcal{T}) = 3$  on  $\mathcal{X} \subset \mathbb{R}$ .

## Definition (Chaining-Tree)

A Chaining-Tree (CT) prediction function  $\hat{f}$  over  $\mathcal{X}$  is defined as

$$\hat{f}(x) = \sum_{n \in \mathcal{N}(\mathcal{T})} h_n(x), \quad x \in \mathcal{X},$$

where:

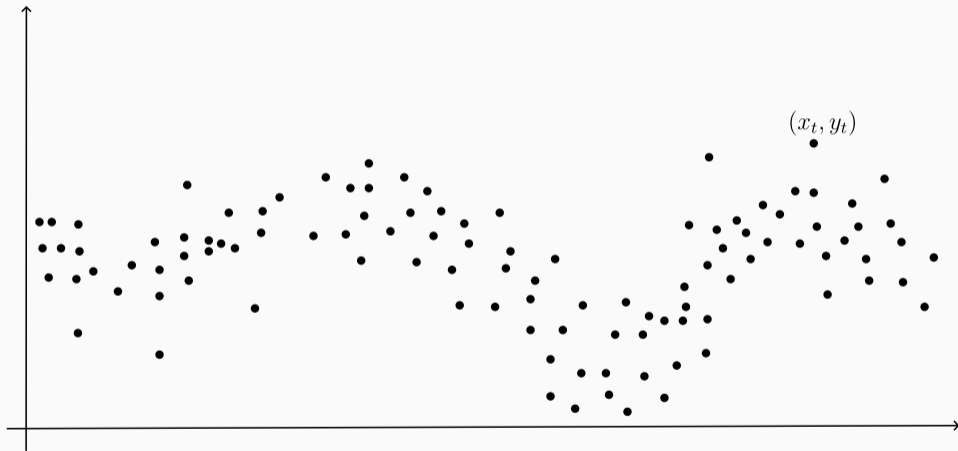
- $h_n(x) = \theta_n \mathbf{1}_{x \in \mathcal{X}_n}$  are constant functions;
- each interior node  $n \in \mathcal{N}(\mathcal{T}) \setminus \mathcal{L}(\mathcal{T})$  has  $2^d$  children forming a regular partition of  $\mathcal{X}_n$ .

 Remark: contrary to standard methods, we predict with *all* nodes  $n \in \mathcal{N}(\mathcal{T})$ .

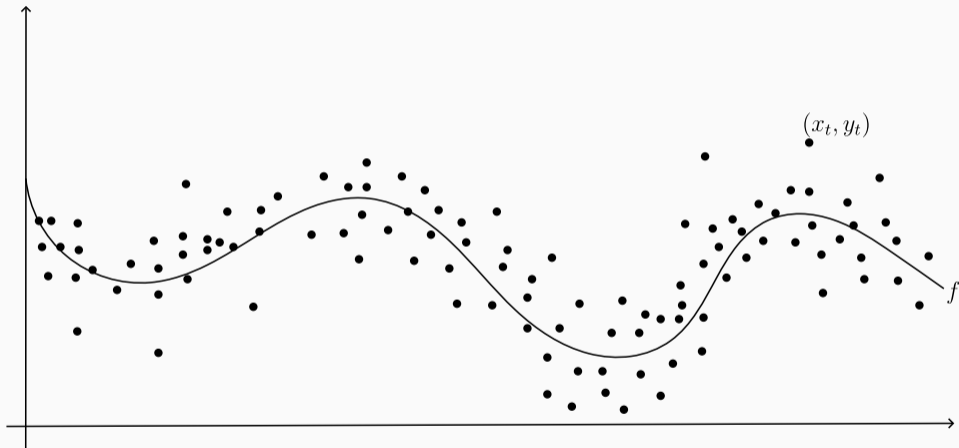
## Illustration of approximation by Chaining-Tree

→ Assume  $\ell_t(\hat{y}) = (\hat{y} - y_t)^2$  square loss function and we launch a CT  $\mathcal{T}$  with depth  $d(\mathcal{T}) = 1, 2, 3$ , over  $T$  data. We have the following illustration:

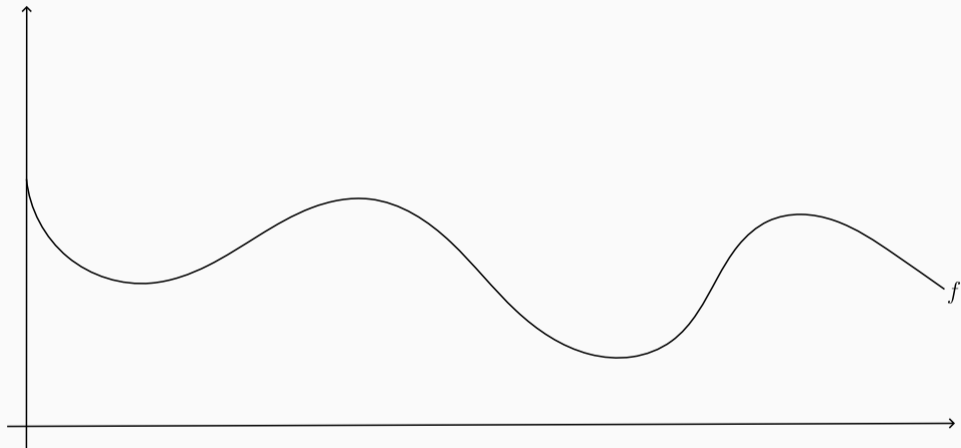
# Illustration of approximation by Chaining-Tree



# Illustration of approximation by Chaining-Tree

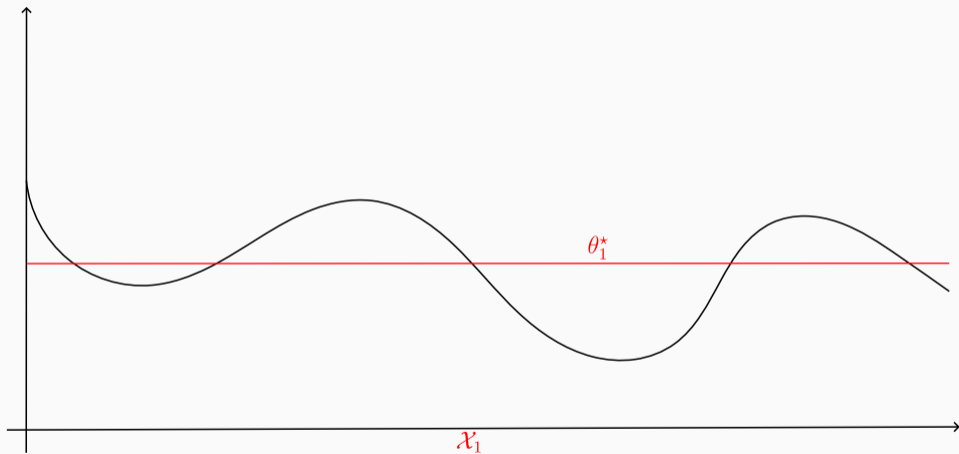


# Illustration of approximation by Chaining-Tree

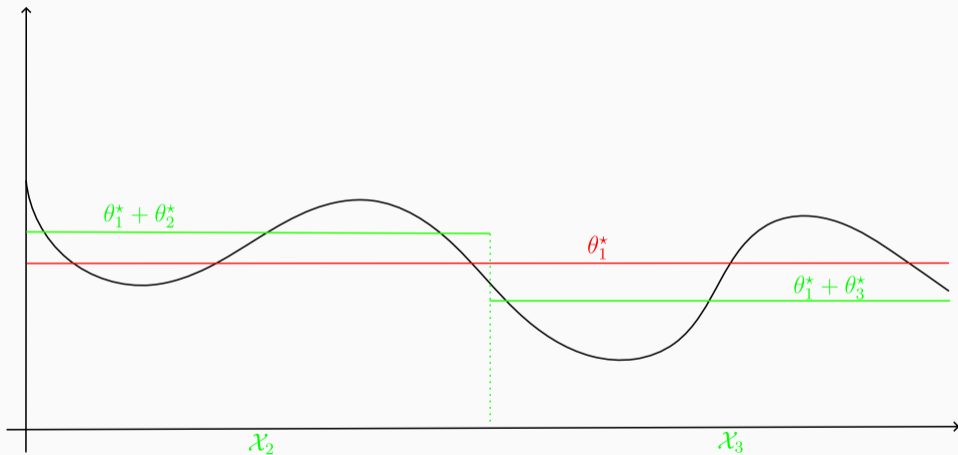




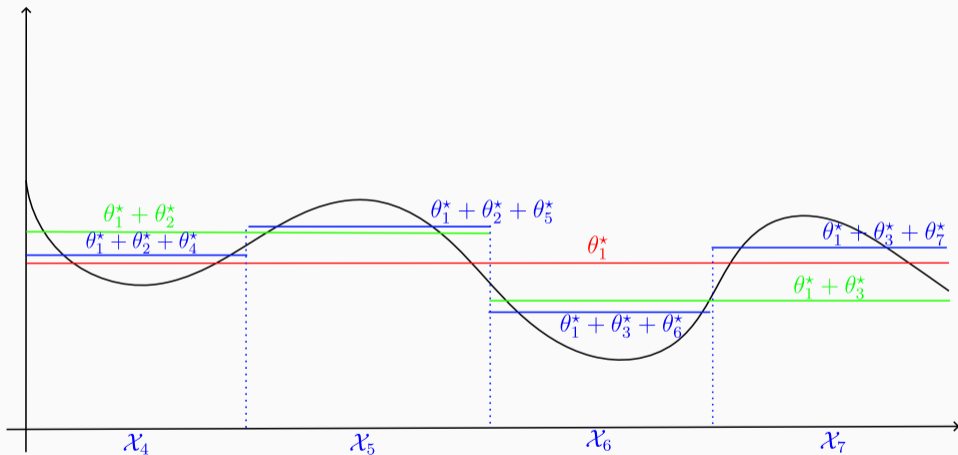
# Illustration of approximation by Chaining-Tree



# Illustration of approximation by Chaining-Tree



# Illustration of approximation by Chaining-Tree



# Online Boosting in a Chaining-Tree

→ **Goal:** Sequentially training CT  $\mathcal{T}$ , i.e. tuning over time the family

$$\bar{\mathcal{W}}_t = \{h_{n,t} = \theta_{n,t} \mathbb{1}_{\mathcal{X}_n}, n \in \mathcal{N}(\mathcal{T})\}.$$



We use **OGB** on  $\bar{\mathcal{W}}_t$ , with  $\beta_n = 1$ ,  $N = |\mathcal{N}(\mathcal{T})|$ . Gradient step becomes, for all  $n \in [N]$ :

$$\theta_{n,t+1} \leftarrow \text{grad-step}(\theta_{n,t}, g_{n,t}), \quad \text{where} \quad g_{n,t} = \ell'_t(\hat{f}_t(x_t)) \mathbb{1}_{x_t \in \mathcal{X}_n}.$$

## Assumption 1 (Parameter Free regret)

Let  $n \in \mathcal{N}(\mathcal{T})$ ,  $\forall g_{n,1}, \dots, g_{n,T} \in [-G, G]$ ,  $G > 0$ , **grad-step** produces  $(\theta_{n,t})$  such that:

$$\sum_{t \in T_n} g_{n,t} (\theta_{n,t} - \theta_n) \lesssim G |\theta_n| \sqrt{|T_n|}, \quad \text{with} \quad T_n = \{1 \leq t \leq T, g_{n,t} \neq 0\},$$

for every  $\theta_n \in \mathbb{R}$ .

e.g., parameter free algorithms in Orabona and Pál; Mhammedi and Koolen; Cutkosky and Orabona (2016; 2020; 2018).

# Optimal Regret and Adaptivity to Hölder functions

$\alpha$ -Hölder continuous functions over  $\mathcal{X} \subset \mathbb{R}^d$ :

$$\text{Lip}_L^\alpha(\mathcal{X}) = \{f : \mathcal{X} \rightarrow \mathbb{R} : |f(x) - f(x')| \leq L\|x - x'\|_\infty^\alpha, \forall x, x' \in \mathcal{X}\}.$$

## Theorem (Regret of OGB-CT vs Hölder functions - Liautaud et al. (2024))

Under Assumption 1, OGB on CT  $(\mathcal{T}, \bar{\mathcal{X}}, \bar{\mathcal{W}})$  with  $\mathcal{X}_{\text{root}} = \mathcal{X}$ ,  $\theta_{n,1} = 0, n \in \mathcal{N}(\mathcal{T})$  and  $d(\mathcal{T}) = \frac{1}{d} \log_2 T$  has regret:

$$\sup_{f \in \text{Lip}_L^\alpha(\mathcal{X})} \text{Reg}_T(f) \lesssim GLX^\alpha \begin{cases} \sqrt{T} & \text{if } d < 2\alpha, \\ \log_2 T \sqrt{T} & \text{if } d = 2\alpha, \\ T^{1-\frac{\alpha}{d}} & \text{if } d > 2\alpha, \end{cases}$$

for any  $L > 0, \alpha \in (0, 1]$ .


# Optimal Regret and Adaptivity to Hölder functions


## Theorem (Regret of OGB-CT vs Hölder functions - Liautaud et al. (2024))

Under Assumption 1, OGB on CT  $(\mathcal{T}, \bar{\mathcal{X}}, \bar{\mathcal{W}})$  with  $\mathcal{X}_{\text{root}} = \mathcal{X}$ ,  $\theta_{n,1} = 0$ ,  $n \in \mathcal{N}(\mathcal{T})$  and  $d(\mathcal{T}) = \frac{1}{d} \log_2 T$  has regret:

$$\sup_{f \in \text{Lip}_L^\alpha(\mathcal{X})} \text{Reg}_T(f) \lesssim GLX^\alpha \begin{cases} \sqrt{T} & \text{if } d < 2\alpha, \\ \log_2 T \sqrt{T} & \text{if } d = 2\alpha, \\ T^{1-\frac{\alpha}{d}} & \text{if } d > 2\alpha, \end{cases}$$

for any  $L > 0$ ,  $\alpha \in (0, 1]$ .

 Our rates are **minimax** over  $\text{Lip}_L^\alpha$  (Rakhlin et al. (2015)) + we **do not need** prior knowledge of neither  $L$  nor  $\alpha$ .

 Computationally tractable:  $x_t$  only falls into one subregion  $\mathcal{X}_n$  for each level  $1, \dots, d(\mathcal{T})$ : we update  $\mathcal{O}(\frac{T}{d} \log_2(T))$  for  $T$  rounds.

# **Adaptive Boosting in Online NonParametric Regression**

---

# Locally Adaptive Boosting - LocAdaBoost



We base our predictions on a core tree  $(\mathcal{T}_0, \bar{\mathcal{X}}, \bar{W})$  associated to:

$$\hat{f}_t(x_t) = \sum_{n \in \mathcal{N}(\mathcal{T}_0)} w_{n,t} \hat{f}_{n,t}(x_t), \quad \forall t \geq 1,$$

where for any  $n \in \mathcal{N}(\mathcal{T}_0)$ :

- $\hat{f}_n$  is a CT rooted at  $\mathcal{X}_n$ ;
- $w_{n,t}$  weight associated.

We use OGB on

- $\beta_{n,t} = w_{n,t}, n \in \mathcal{N}(\mathcal{T}_0)$ ;
- and gradient  $g_t = \nabla_{(w_n)} \ell_t(\hat{f}_t(x_t))|_{(w_n)=(w_{n,t})}$ .

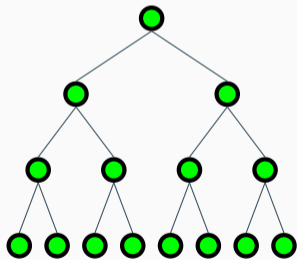


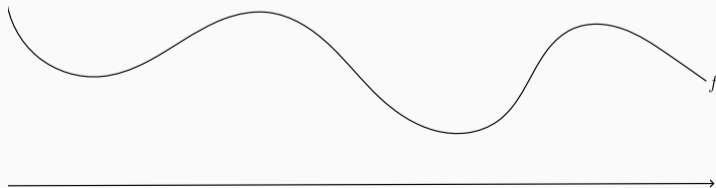
Figure 4:  $\mathcal{T}_0$



## Adaptivity to local profile of the competitor

→ **Goal:** Learn the best pruned tree from  $\mathcal{T}_0$  in  $\mathcal{P}(\mathcal{T}_0)$  to fit the competitor.

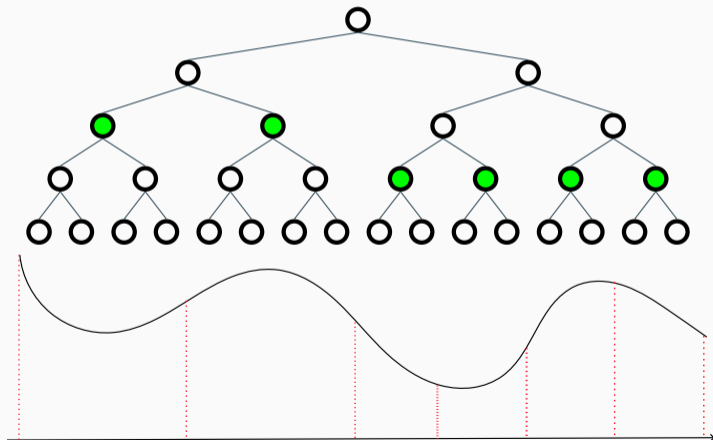
Example 1:



# Adaptivity to local profile of the competitor

→ **Goal:** Learn the best pruned tree from  $\mathcal{T}_0$  in  $\mathcal{P}(\mathcal{T}_0)$  to fit the competitor.

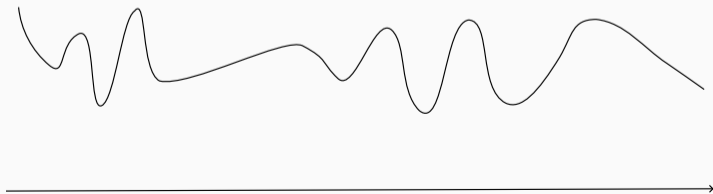
Example 1:



## Adaptivity to local profile of the competitor

→ **Goal:** Learn the best pruned tree from  $\mathcal{T}_0$  in  $\mathcal{P}(\mathcal{T}_0)$  to fit the competitor.

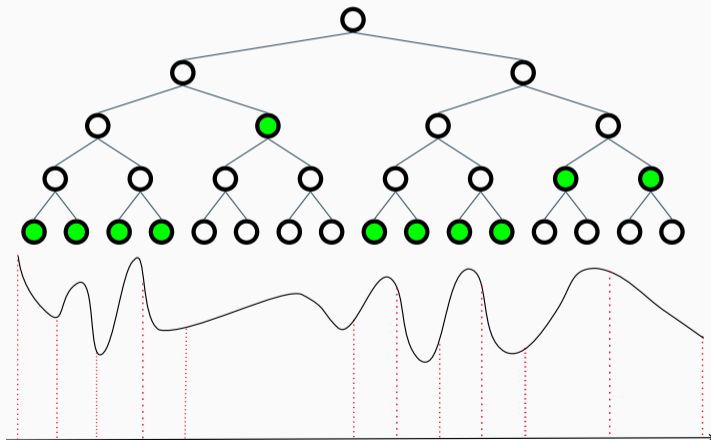
Example 2:



# Adaptivity to local profile of the competitor

→ **Goal:** Learn the best pruned tree from  $\mathcal{T}_0$  in  $\mathcal{P}(\mathcal{T}_0)$  to fit the competitor.

Example 2:



# Optimal and Locally Adaptive Regret (1/2)

**Theorem (Locally Adaptive Regret, case  $d = 1, \alpha > \frac{1}{2}$  - Liautaud et al. (2024))**

Under assumptions, for any  $f \in \text{Lip}_L^\alpha(\mathcal{X})$ , *LocAdaBoost* achieves

$$\text{Reg}_T(f) \lesssim \inf_{\mathcal{T} \in \mathcal{P}(\mathcal{T}_0)} \left\{ \sqrt{T|\mathcal{L}(\mathcal{T})|} + |\mathcal{L}(\mathcal{T})| + X^\alpha \sum_{n \in \mathcal{L}(\mathcal{T})} L_n(f) 2^{-\alpha d(n)} \sqrt{|T_n|} \right\},$$

with  $L_n(f)$  local Hölder constants.

If  $(\ell_t)$  are exp-concave (e.g. square loss)

$$\text{Reg}_T(f) \lesssim \inf_{\mathcal{T} \in \mathcal{P}(\mathcal{T}_0)} \left\{ |\mathcal{L}(\mathcal{T})| + X^\alpha \sum_{n \in \mathcal{L}(\mathcal{T})} L_n(f) 2^{-\alpha d(n)} \sqrt{|T_n|} \right\}$$



Remark: *LocAdaBoost* could also adapt to local regularities  $(\alpha_n)$

## Optimal and Locally Adaptive Regret (2/2)

### Corollary (Minimax Regret - Liautaud et al. (2024))

For any  $f \in \text{Lip}_L^\alpha(\mathcal{X})$ ,  $L > 0$ , *LocAdaBoost* achieves

$$\text{Reg}_T(f) \lesssim \begin{cases} (X^\alpha \bar{L}(f))^{\frac{2}{2\alpha+1}} T^{\frac{1}{2\alpha+1}} & \text{if } \ell_t \text{ are exp-concave,} \\ (X^\alpha \bar{L}(f))^{\frac{1}{2\alpha}} \sqrt{T}, & \end{cases}$$

where  $\bar{L}(f) = \left(\frac{1}{X} \sum_{n \in \mathcal{L}(T)} |\mathcal{X}_n| L_n(f)^{1/\alpha}\right)^\alpha$ .

- ✓ Minimax optimality
- ✓ Adaptivity to local regularities ( $L_n$ ) and  $\alpha$
- ✓ Adaptivity to the loss curvature

# Conclusion

- New generic Online Gradient Boosting procedure;
- Online Gradient Boosting coupled with Chaining-Tree achieve **minimax regret**;
- Our unique **LocAdaBoost** algorithm both adapts optimally to **local regularities** of the competitor and **curvature** of sequential losses;
- **First** constructive algorithm to achieve **optimal locally adaptive regret**;
- Future work: extend the boosting procedure to other learners to approach other classes of functions.

Thank you!

Questions?

---

[1] Link to the paper: <https://arxiv.org/abs/2410.03363>

## Comparison with the literature

Ref.	Assumptions	Upper bound
[2]	$(\ell_t)$ exp-concave, $L > 0$ unknown $(\ell_t)$ convex, $L > 0$ unknown	$\min \{ \sqrt{LT}, L^{\frac{2}{3}} T^{\frac{1}{3}} \}$ $\sqrt{LT}$
[3]	$(\ell_t)$ square loss, $L > 0$ unknown	$\sqrt{LT}$
[4]	$(\ell_t)$ absolute loss, $L > 0$ known $(\ell_t)$ square loss, $L > 0$ known	$L^{\frac{1}{3}} T^{\frac{2}{3}}$ $\sqrt{LT}$
[5]	$(\ell_t)$ square loss, $L = 1$ known	$T^{\frac{1}{3}}$
[6]	$(\ell_t)$ convex, $L = 1$ known	$\sqrt{T}$

[2] Liataud, Gaillard, and Wintenberger, “Minimax Adaptive Boosting for Online Nonparametric Regression”.

[3] Kuzborskij and Cesa-Bianchi, “Locally-adaptive nonparametric online learning”.

[4] Hazan, Agarwal, and Kale, “Logarithmic regret algorithms for online convex optimization”.

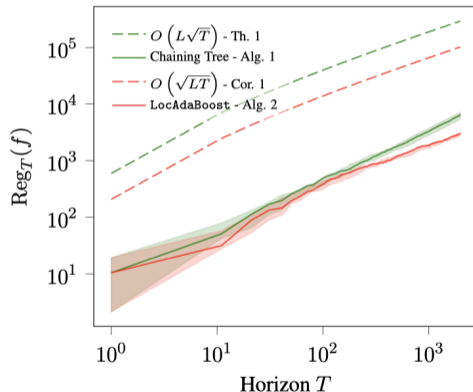
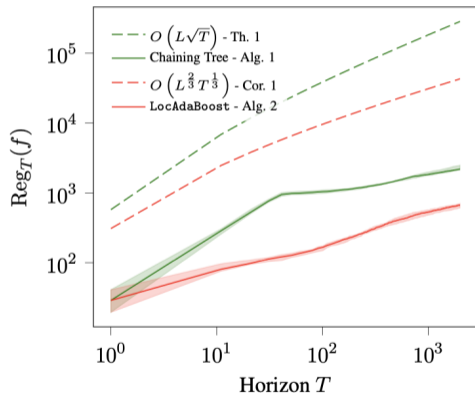
[5] Gaillard and Gerchinovitz, “A Chaining Algorithm for Online Nonparametric Regression”.

[6] Cesa-Bianchi et al., “Algorithmic chaining and the role of partial feedback in online nonparametric learning”.

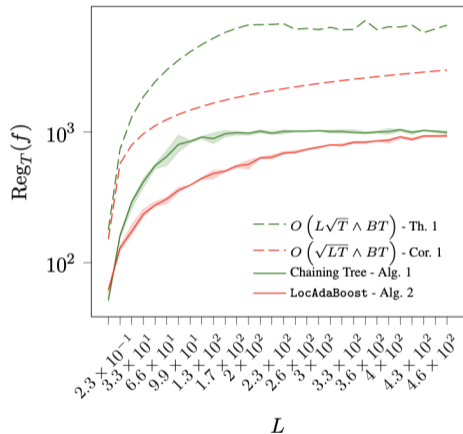
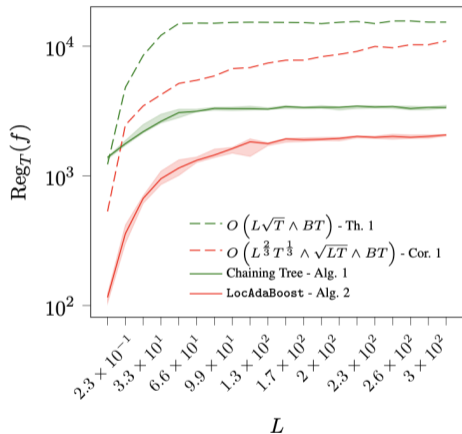


# Experiments (1/3)

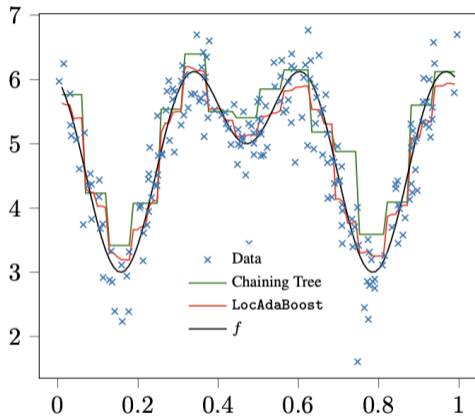
Regression setting:  $y_t = f(x_t) + \varepsilon_t$ , where  $\varepsilon_t \sim \mathcal{N}(0, \sigma^2)$  with  $\sigma = 0.5, f(x) = \sin(10x) + \cos(5x) + 5$ , for  $x \in \mathcal{X} = [0, 1]$  and  $\sup_x |f'(x)| \leq 15 =: L$ .



# Experiments (2/3)















## Experiments (3/3)



# References

---

-  Cesa-Bianchi, Nicolò and Gábor Lugosi (2006). *Prediction, Learning, and Games*. Cambridge University Press.
-  Cesa-Bianchi, Nicolò et al. (2017). “Algorithmic chaining and the role of partial feedback in online nonparametric learning”. In: *Conference on Learning Theory*. PMLR, pp. 465–481.
-  Cutkosky, Ashok and Francesco Orabona (2018). “Black-box reductions for parameter-free online learning in banach spaces”. In: *Conference On Learning Theory*. PMLR, pp. 1493–1529.
-  Gaillard, Pierre and Sebastien Gerchinovitz (2015). “A Chaining Algorithm for Online Nonparametric Regression”. In: *COLT*.
-  Hazan, Elad, Amit Agarwal, and Satyen Kale (2007). “Logarithmic regret algorithms for online convex optimization”. In: *Machine Learning* 69.2, pp. 169–192.
-  Kuzborskij, Ilja and Nicolo Cesa-Bianchi (2020). “Locally-adaptive nonparametric online learning”. In: *Advances in Neural Information Processing Systems* 33, pp. 1679–1689.
-  Liautaud, Paul, Pierre Gaillard, and Olivier Wintenberger (2024). “Minimax Adaptive Boosting for Online Nonparametric Regression”. In: *arXiv preprint arXiv:2410.03363*.

-  Mhammedi, Zakaria and Wouter M Koolen (2020). “Lipschitz and comparator-norm adaptivity in online learning”. In: *Conference on Learning Theory*. PMLR, pp. 2858–2887.
-  Orabona, Francesco and Dávid Pál (2016). “Coin betting and parameter-free online learning”. In: *Advances in Neural Information Processing Systems* 29.
-  Rakhlin, Alexander and Karthik Sridharan (2014). “Online non-parametric regression”. In: *Conference on Learning Theory*. PMLR, pp. 1232–1264.
-  — (2015). “Online nonparametric regression with general loss functions”. In: *arXiv preprint arXiv:1501.06598*.
-  Zinkevich, Martin (2003). “Online convex programming and generalized infinitesimal gradient ascent”. In: *Proceedings of the 20th international conference on machine learning (icml-03)*, pp. 928–936.