

A CHAINING ALGORITHM FOR ONLINE NON PARAMETRIC REGRESSION

Pierre Gaillard

December 2, 2015

University of Copenhagen

This is joint work with Sebastien Gerchinovitz

1. Online prediction of arbitrary sequences
2. Finite reference class: prediction with expert advice
3. Large reference class
4. Extensions, current (and future) work

ONLINE PREDICTION OF ARBITRARY
SEQUENCES

Sequential prediction of arbitrary time-series¹:

- a time-series $y_1, \dots, y_n \in \mathcal{Y} = [-B, B]$ is to be predicted step by step
- covariates $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathcal{X}$ are sequentially available

At each forecasting instance $t = 1, \dots, n$

- the environment reveals $\mathbf{x}_t \in \mathcal{X}$
- the player is asked to form a prediction \hat{y}_t of y_t based on
 - the **past** observations y_1, \dots, y_{t-1}
 - the **current** and **past** covariates $\mathbf{x}_1, \dots, \mathbf{x}_t$
- the environment reveals y_t

Goal: minimize the cumulative loss: $\hat{L}_n = \sum_{t=1}^n (\hat{y}_t - y_t)^2$.

Difficulty: **no stochastic assumption** on the time series

- neither on the observations (y_t)
- neither on the covariates (\mathbf{x}_t)

¹ N. Cesa-Bianchi and G. Lugosi. *Prediction, learning, and games*. 2006.

Sequential prediction of arbitrary time-series:

- a time-series $y_1, \dots, y_n \in \mathcal{Y} = [-B, B]$ is to be predicted step by step
- covariates $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathcal{X}$ are sequentially available

At each forecasting instance $t = 1, \dots, n$

- the environment reveals $\mathbf{x}_t \in \mathcal{X}$
- solution: produce the prediction as a function of \mathbf{x}_t

$$\hat{y}_t = \hat{f}_t(\mathbf{x}_t)$$

- the environment reveals y_t

Goal: minimize our regret against a reference function class $\mathcal{F} \in \mathcal{Y}^{\mathcal{X}}$

$$\text{Reg}_n(\mathcal{F}) \stackrel{\text{def}}{=} \underbrace{\sum_{t=1}^n (\hat{f}_t(\mathbf{x}_t) - y_t)^2}_{\text{our performance}} - \underbrace{\inf_{f \in \mathcal{F}} \sum_{t=1}^n (f(\mathbf{x}_t) - y_t)^2}_{\text{reference performance}}$$

Sequential prediction of arbitrary time-series:

- a time-series $y_1, \dots, y_n \in \mathcal{Y} = [-B, B]$ is to be predicted step by step
- covariates $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathcal{X}$ are sequentially available

At each forecasting instance $t = 1, \dots, n$

- the environment reveals $\mathbf{x}_t \in \mathcal{X}$
- solution: produce the prediction as a function of \mathbf{x}_t

$$\hat{y}_t = \hat{f}_t(\mathbf{x}_t)$$

- the environment reveals y_t

Goal: minimize our regret against a reference function class $\mathcal{F} \in \mathcal{Y}^{\mathcal{X}}$

$$\text{Reg}_n(\mathcal{F}) \stackrel{\text{def}}{=} \underbrace{\sum_{t=1}^n (\hat{f}_t(\mathbf{x}_t) - y_t)^2}_{\text{our performance}} - \underbrace{\inf_{f \in \mathcal{F}} \sum_{t=1}^n (f(\mathbf{x}_t) - y_t)^2}_{\text{reference performance}} = \underbrace{o(n)}_{\text{Goal}}$$

FINITE REFERENCE CLASS: PREDICTION
WITH EXPERT ADVICE

Assumption: $\mathcal{F} = \{f_1, \dots, f_K\} \subset \mathcal{Y}^{\mathcal{X}}$ is finite

The exponentially weighted average forecaster (EWA)¹

At each forecasting instance t ,

- assign to each function f_k the weight

$$\hat{p}_{k,t} = \frac{\exp\left(-\eta \sum_{s=1}^{t-1} (y_s - f_k(\mathbf{x}_s))^2\right)}{\sum_{j=1}^K \exp\left(-\eta \sum_{s=1}^{t-1} (y_s - f_j(\mathbf{x}_s))^2\right)}$$

- form function $\hat{f}_t = \sum_{k=1}^K \hat{p}_{k,t} f_k$ and predict $\hat{y}_t = \hat{f}_t(\mathbf{x}_t)$

Performance: if $\mathcal{Y} = [-B, B]$ and $\eta = 1/(8B^2)$

$$\text{Reg}_n(\mathcal{F}) \stackrel{\text{def}}{=} \sum_{t=1}^n (y_t - \hat{f}_t(\mathbf{x}_t))^2 - \inf_{f \in \mathcal{F}} \sum_{t=1}^n (y_t - f(\mathbf{x}_t))^2 \leq 8B^2 \log K$$

If B is not known in advance, η can be tuned online (doubling trick).

¹ Littlestone and M. K. Warmuth (1994) and Vovk (1990)

1. Upper bound the instantaneous loss

$$\begin{aligned}
 (y_t - \widehat{f}_t(\mathbf{x}_t))^2 &= (y_t - \sum_{k=1}^K \widehat{p}_{k,t} f_k(\mathbf{x}_t))^2 \\
 &\stackrel{\text{for } \eta \leq 1/(8B^2)}{\leq} -\frac{1}{\eta} \log \left(\sum_{k=1}^K \widehat{p}_{k,t} e^{-\eta (y_t - f_k(\mathbf{x}_t))^2} \right) \\
 &\stackrel{\text{by definition of } \widehat{p}_{k,t+1}}{=} -\frac{1}{\eta} \log \left(\frac{\widehat{p}_{k,t}}{\widehat{p}_{k,t+1}} e^{-\eta (y_t - f_k(\mathbf{x}_t))^2} \right) \\
 &= (y_t - f_k(\mathbf{x}_t))^2 + \frac{1}{\eta} \log \frac{\widehat{p}_{k,t+1}}{\widehat{p}_{k,t}}
 \end{aligned}$$

2. Sum over all t , the sum telescopes

$$\sum_{t=1}^n (y_t - \widehat{f}_t(\mathbf{x}_t))^2 - (y_t - f_k(\mathbf{x}_t))^2 \leq \frac{1}{\eta} \log \frac{\widehat{p}_{k,n+1}}{\widehat{p}_{k,1}} \leq \frac{\log K}{\eta} = 8B^2 \log K$$

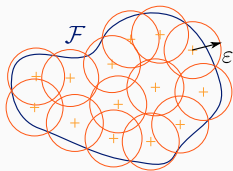
LARGE REFERENCE CLASS

1. Approximate \mathcal{F} by a finite set \mathcal{F}_ε such that

$$\forall f \in \mathcal{F} \quad \exists f_\varepsilon \in \mathcal{F}_\varepsilon \quad \|f - f_\varepsilon\|_\infty \leq \varepsilon. \quad (1)$$

Such set \mathcal{F}_ε is called an ε -net of \mathcal{F}

2. Run EWA on \mathcal{F}_ε



Definition (metric entropy)

The cardinal of the smallest ε -net \mathcal{F}_ε that satisfies (1) is denoted $\mathcal{N}_\infty(\mathcal{F}, \varepsilon)$. The **metric entropy** of \mathcal{F} is $\log \mathcal{N}_\infty(\mathcal{F}, \varepsilon)$.

Regret bound of order (forgetting constants):

$$\begin{aligned} \text{Reg}_n(\mathcal{F}) &= \text{Reg}_n(\mathcal{F}_\varepsilon) + \left[\inf_{f_\varepsilon \in \mathcal{F}_\varepsilon} \sum_{t=1}^n (y_t - f_\varepsilon(x_t))^2 - \inf_{f \in \mathcal{F}} \sum_{t=1}^n (y_t - f(x_t))^2 \right] \\ &\lesssim \underbrace{\log \mathcal{N}_\infty(\mathcal{F}, \varepsilon)}_{\text{Regret of EWA on } \mathcal{F}_\varepsilon} + \underbrace{\varepsilon n}_{\text{Approximation of } \mathcal{F} \text{ by } \mathcal{F}_\varepsilon} \end{aligned}$$

EXAMPLES OF REFERENCE CLASSES: THE PARAMETRIC CASE

If $\mathcal{N}_\infty(\mathcal{F}, \varepsilon) \lesssim \varepsilon^{-p}$ for $p > 0$ as $\varepsilon \rightarrow 0$,

$$\begin{aligned} \text{Reg}_n(\mathcal{F}) &\lesssim \log \mathcal{N}_\infty(\mathcal{F}, \varepsilon) + \varepsilon n \\ &\approx \log(\varepsilon^{-p}) + \varepsilon n \stackrel{\varepsilon \approx 1/n}{\approx} p \log(n) \end{aligned}$$

Example

Assume you have $d \geq 1$ black-box forecasters $\varphi_1, \dots, \varphi_d \in \mathcal{X}^{\mathcal{Y}}$

- linear regression in a compact ball

$$\mathcal{F} = \left\{ \sum_{j=1}^d u_j \varphi_j : \text{for } \mathbf{u} \in \Theta \underset{\text{comp.}}{\subset} \mathbb{R}^d \right\} \rightarrow \mathcal{N}_\infty(\mathcal{F}, \varepsilon) \lesssim \varepsilon^{-d}$$

- sparse linear regression

$$\mathcal{F} = \left\{ \sum_{j=1}^d u_j \varphi_j : \text{for } \mathbf{u} \in [0, 1]^d \text{ s.t. } \|\mathbf{u}\|_1 = 1 \text{ and } \|\mathbf{u}\|_0 = s \right\}$$

Then²,

$$\log \mathcal{N}_\infty(\mathcal{F}, \varepsilon) \lesssim \log \binom{d}{s} + s \log(1 + 1/(\varepsilon \sqrt{s})) \rightarrow \text{Reg}_n(\mathcal{F}) \lesssim s \log(1 + dn/s)$$

² F. Gao, C-K. Ing, and Y. Yang. "Metric entropy and sparse linear approximation of ℓ_q -hulls for $0 < q \leq 1$ ". In: *Journal of Approximation Theory* (2013).

EXAMPLES OF REFERENCE CLASSES: THE PARAMETRIC CASE

If $\mathcal{N}_\infty(\mathcal{F}, \varepsilon) \lesssim \varepsilon^{-p}$ for $p > 0$ as $\varepsilon \rightarrow 0$,

$$\begin{aligned} \text{Reg}_n(\mathcal{F}) &\lesssim \log \mathcal{N}_\infty(\mathcal{F}, \varepsilon) + \varepsilon n \\ &\approx \log(\varepsilon^{-p}) + \varepsilon n \stackrel{\varepsilon \approx 1/n}{\approx} p \log(n) \quad \rightarrow \text{optimal} \end{aligned}$$

Example

Assume you have $d \geq 1$ black-box forecasters $\varphi_1, \dots, \varphi_d \in \mathcal{X}^{\mathcal{Y}}$

- linear regression in a compact ball

$$\mathcal{F} = \left\{ \sum_{j=1}^d u_j \varphi_j : \text{for } \mathbf{u} \in \Theta \underset{\text{comp.}}{\subset} \mathbb{R}^d \right\} \quad \rightarrow \quad \mathcal{N}_\infty(\mathcal{F}, \varepsilon) \lesssim \varepsilon^{-d}$$

- sparse linear regression

$$\mathcal{F} = \left\{ \sum_{j=1}^d u_j \varphi_j : \text{for } \mathbf{u} \in [0, 1]^d \text{ s.t. } \|\mathbf{u}\|_1 = 1 \text{ and } \|\mathbf{u}\|_0 = s \right\}$$

Then²,

$$\log \mathcal{N}_\infty(\mathcal{F}, \varepsilon) \lesssim \log \binom{d}{s} + s \log(1 + 1/(\varepsilon \sqrt{s})) \quad \rightarrow \quad \text{Reg}_n(\mathcal{F}) \lesssim s \log(1 + dn/s)$$

² F. Gao, C-K. Ing, and Y. Yang. "Metric entropy and sparse linear approximation of ℓ_q -hulls for $0 < q \leq 1$ ". In: *Journal of Approximation Theory* (2013).

WHAT IF \mathcal{F} IS NON PARAMETRIC?

if $\log \mathcal{N}_\infty(\mathcal{F}, \varepsilon) \lesssim \varepsilon^{-p}$ for $p > 0$ as $\varepsilon \rightarrow 0$,

$$\begin{aligned} \text{Reg}_n(\mathcal{F}) &\lesssim \log \mathcal{N}_\infty(\mathcal{F}, \varepsilon) + \varepsilon n \\ &\lesssim \varepsilon^{-p} + \varepsilon n \quad \begin{array}{l} \varepsilon = n^{-1/(p+1)} \\ \approx n^{\frac{p}{p+1}} \end{array} \end{aligned}$$

Example

- 1-Lipschitz ball on $[0, 1]$

$$\mathcal{F} = \left\{ f \in \mathcal{Y}^{\mathcal{X}} : \forall x, y \in \mathcal{X} \subset [0, 1] \quad \|f(x) - f(y)\| \leq \|x - y\| \right\}$$

Then $\log \mathcal{N}_\infty(\mathcal{F}, \varepsilon) \approx \varepsilon^{-1} \rightarrow \text{Reg}_n(\mathcal{F}) \lesssim \sqrt{n}$

- Hölder ball on $\mathcal{X} \subset [0, 1]$ with regularity $\beta = q + \alpha > 1/2$

$$\mathcal{F} = \left\{ f \in \mathcal{Y}^{\mathcal{X}} : \forall x, y \in \mathcal{X} \quad |f^{(q)}(x) - f^{(q)}(y)| \leq |x - y|^\alpha \text{ and } \forall k \leq q, \|f^{(k)}\|_\infty \leq B \right\}$$

Then³ $\log \mathcal{N}_\infty(\mathcal{F}, \varepsilon) \approx \varepsilon^{-1/\beta} \rightarrow \text{Reg}_n(\mathcal{F}) \lesssim n^{\frac{1}{1+\beta}}$.

³ G.G. Lorentz. "Metric Entropy, Widths, and Superpositions of Functions". In: *Amer. Math. Monthly* 6 (1962).

WHAT IF \mathcal{F} IS NON PARAMETRIC?

if $\log \mathcal{N}_\infty(\mathcal{F}, \varepsilon) \lesssim \varepsilon^{-p}$ for $p > 0$ as $\varepsilon \rightarrow 0$,

$$\begin{aligned} \text{Reg}_n(\mathcal{F}) &\lesssim \log \mathcal{N}_\infty(\mathcal{F}, \varepsilon) + \varepsilon n \\ &\lesssim \varepsilon^{-p} + \varepsilon n \end{aligned} \quad \begin{array}{l} \varepsilon = n^{-1/(p+1)} \\ \approx n^{\frac{p}{p+1}} \end{array}$$

→ suboptimal:

$$\begin{array}{ll} n^{\frac{p}{p+2}} & \text{if } p < 2 \\ n^{\frac{p-1}{p}} & \text{if } p > 2 \end{array}$$

Example

- 1-Lipschitz ball on $[0, 1]$

$$\mathcal{F} = \left\{ f \in \mathcal{Y}^{\mathcal{X}} : \forall x, y \in \mathcal{X} \subset [0, 1] \quad \|f(x) - f(y)\| \leq \|x - y\| \right\}$$

Then $\log \mathcal{N}_\infty(\mathcal{F}, \varepsilon) \approx \varepsilon^{-1} \rightarrow \text{Reg}_n(\mathcal{F}) \lesssim \sqrt{n} \rightarrow$ suboptimal: $n^{\frac{1}{3}}$

- Hölder ball on $\mathcal{X} \subset [0, 1]$ with regularity $\beta = q + \alpha > 1/2$

$$\mathcal{F} = \left\{ f \in \mathcal{Y}^{\mathcal{X}} : \forall x, y \in \mathcal{X} \quad |f^{(q)}(x) - f^{(q)}(y)| \leq |x - y|^\alpha \text{ and } \forall k \leq q, \|f^{(k)}\|_\infty \leq B \right\}$$

Then³ $\log \mathcal{N}_\infty(\mathcal{F}, \varepsilon) \approx \varepsilon^{-1/\beta} \rightarrow \text{Reg}_n(\mathcal{F}) \lesssim n^{\frac{1}{1+\beta}} \rightarrow$ suboptimal: $n^{\frac{1}{1+2\beta}}$.

³ G.G. Lorentz. "Metric Entropy, Widths, and Superpositions of Functions". In: *Amer. Math. Monthly* 6 (1962).

Theorem (Rakhlin and Sridharan 2014⁴)

The minimax rate of the regret if of order

$$\inf_{\gamma \geq \epsilon \geq 0} \left\{ \log \mathcal{N}^{\text{seq}}(\mathcal{F}, \gamma) + \sqrt{n} \int_{\epsilon}^{\gamma} \sqrt{\log \mathcal{N}^{\text{seq}}(\tau, \mathcal{F})} d\tau + \epsilon n \right\}$$

where $\log \mathcal{N}^{\text{seq}}(\mathcal{F}, \epsilon) \leq \log \mathcal{N}_{\infty}(\mathcal{F}, \epsilon)$ is the sequential entropy of \mathcal{F} .

$\log \mathcal{N}_{\infty}(\mathcal{F}, \gamma)$: regret of EWA against γ -net → crude approximation

ϵn : approximation error of the ϵ -net → fine approximation

$\sqrt{n} \int_{\epsilon}^{\gamma} \sqrt{\log \mathcal{N}_{\infty}(\mathcal{F}, \tau)} d\tau$: from large scale γ to small scale ϵ .

This term is a **Dudley entropy integral** that appears in

- Chaining to bound the supremum of a stochastic process (Dudley 1967)
- Statistical learning with i.i.d. data to derive risk bounds (e.g., Massart 2007; Rakhlin et al. 2013)
- Online learning with arbitrary sequences (Opper and Haussler 1997; Cesa-Bianchi and Lugosi 1999)

⁴ A. Rakhlin and K. Sridharan. "Online Nonparametric Regression". In: *COLT* (2014).

Theorem (Rakhlin and Sridharan 2014⁴)

The minimax rate of the regret if of order

$$\inf_{\gamma \geq \varepsilon \geq 0} \left\{ \log \mathcal{N}_\infty(\mathcal{F}, \gamma) + \sqrt{n} \int_\varepsilon^\gamma \sqrt{\log \mathcal{N}_\infty(\tau, \mathcal{F})} d\tau + \varepsilon n \right\}$$

if $\log \mathcal{N}_\infty(\mathcal{F}, \varepsilon) \approx \log \mathcal{N}^{\text{seq}}(\mathcal{F}, \varepsilon)$.

$\log \mathcal{N}_\infty(\mathcal{F}, \gamma)$: regret of EWA against γ -net → crude approximation

εn : approximation error of the ε -net → fine approximation

$\sqrt{n} \int_\varepsilon^\gamma \sqrt{\log \mathcal{N}_\infty(\mathcal{F}, \tau)} d\tau$: from large scale γ to small scale ε .

This term is a **Dudley entropy integral** that appears in

- Chaining to bound the supremum of a stochastic process (Dudley 1967)
- Statistical learning with i.i.d. data to derive risk bounds (e.g., Massart 2007; Rakhlin et al. 2013)
- Online learning with arbitrary sequences (Opper and Haussler 1997; Cesa-Bianchi and Lugosi 1999)

⁴ A. Rakhlin and K. Sridharan. "Online Nonparametric Regression". In: *COLT* (2014).

Theorem (Rakhlin and Sridharan 2014⁴)

The minimax rate of the regret if of order

$$\inf_{\gamma \geq \varepsilon \geq 0} \left\{ \log \mathcal{N}_\infty(\mathcal{F}, \gamma) + \sqrt{n} \int_\varepsilon^\gamma \sqrt{\log \mathcal{N}_\infty(\tau, \mathcal{F})} d\tau + \varepsilon n \right\}$$

if $\log \mathcal{N}_\infty(\mathcal{F}, \varepsilon) \approx \log \mathcal{N}^{\text{seq}}(\mathcal{F}, \varepsilon)$.

Example: let $p \in (0, 2)$ and \mathcal{F} such that

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

The minimax regret is then of order

$$\gamma^{-p} + \sqrt{n} \int_\varepsilon^\gamma \tau^{-p/2} d\tau + \varepsilon n = \gamma^{-p} + \sqrt{n} \gamma^{1-p/2} + 0 \approx n^{\frac{p}{p+2}}$$

for the optimal choices $\varepsilon = 0$ and $\gamma \approx n^{-1/(p+2)}$.

⁴ A. Rakhlin and K. Sridharan. "Online Nonparametric Regression". In: COLT (2014).

Main algorithm which:

- achieves the Dudley-type regret bound

$$\text{Reg}_n \lesssim \log \mathcal{N}_\infty(\mathcal{F}, \gamma) + \sqrt{n} \int_\epsilon^\gamma \sqrt{\log \mathcal{N}_\infty(\mathcal{F}, \tau)} d\tau + \epsilon n$$

- **efficient** version for Hölder class in $[0, 1]$ (costs a log factor)

Key-subroutine (Multi-variable EG) to go from scale γ to scale ϵ .

Function class	Metric entropy	Regret of EWA	Our Regret
	$\epsilon^{-p} \quad p \in (0, 2)$	$n^{p/(p+1)}$	$n^{p/(p+2)}$
Lipschitz on $[0, 1]$	ϵ^{-1}	$n^{1/2}$	$n^{1/3}$
β -Hölder on $[0, 1]$	$\epsilon^{-1/\beta} \quad \beta > 1/2$	$n^{1/(\beta+1)}$	$n^{1/(2\beta+1)}$
Sparse lin. reg.	$\log \binom{d}{s} + s \log(1 + 1/(\epsilon\sqrt{s}))$	$s \log(1 + dn/s)$	$s \log(1 + dn/s)$

Instead of competing directly with \mathcal{F}_ε for small ε (too many functions)

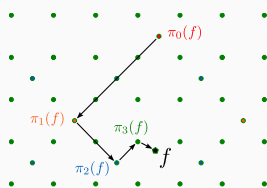
1. create a “chain” of refining approximations

$\pi_0(f) \in \mathcal{F}^{(0)}, \pi_1(f) \in \mathcal{F}^{(1)}, \dots$ of any function $f \in \mathcal{F}_\varepsilon$ such that

$$\forall k \geq 0 \quad \sup_{f \in \mathcal{F}} \|\pi_k(f) - f\|_\infty \leq \gamma/2^k$$

and

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



2. compete with the chains

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

|small increments| $\leq 3\gamma/2^{k+1}$

where $\mathcal{G}^{(k)} \stackrel{\text{def}}{=} \{\pi_k(f) - \pi_{k-1}(f) : f \in \mathcal{F}\}$

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

|small increments| $\leq 3\gamma/2^{k+1}$

It thus suffices to compete with

$$\inf_{f_0 \in \mathcal{F}^{(0)}} \left\{ \underbrace{\inf_{g_0 \in \mathcal{G}^{(0)}, \dots, g_{K_\varepsilon} \in \mathcal{G}^{(K_\varepsilon)}} \sum_{t=1}^n \left(y_t - (f_0 + g_0 + \dots + g_{K_\varepsilon})(x_t) \right)^2}_{\text{low-scale aggregation}} \right\}$$

low-scale aggregation: gradient descents

simultaneously at all g_k

Regret cost: $\sqrt{n} \int_\varepsilon^\gamma \sqrt{\log \mathcal{N}_\infty(\mathcal{F}, \tau)} d\tau$

high-scale aggregation: run EWA to be competitive against all $f_0 \in \mathcal{F}^{(0)}$

Regret cost: $\log(\text{Card} \mathcal{F}^{(0)}) = \log \mathcal{N}_\infty(\mathcal{F}, \gamma)$

THE EXPONENTIATED GRADIENT FORECASTER (EG)

Let $\Delta_N \stackrel{\text{def}}{=} \{u \in \mathbb{R}_+^N : \sum_{i=1}^N u_i = 1\}$.

Setting: at each step $t \geq 1$, player plays $\hat{u}_t \in \Delta_N$ and environment chooses **convex** and **differentiable** loss function $\ell_t : \Delta_N \rightarrow \mathbb{R}$.

The Exponentiated Gradient forecaster⁵

At each forecasting instance $t \geq 1$, define $\hat{u}_t \in \Delta_N$ component-wise by

$$\hat{u}_{k,t} \stackrel{\text{def}}{=} \frac{1}{Z_t} \exp \left(-\eta \sum_{s=1}^{t-1} \partial_{\hat{u}_{k,s}} \ell_s(\hat{u}_s) \right)$$

where Z_t is a normalization factor.

Regret bound: if $\|\nabla \ell_t\|_\infty \leq G$. For $\eta = G^{-1} \sqrt{2(\log N)/n}$

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

If G is small, $G \sqrt{n \log N}$ can be better than $B^2 \log N$.

⁵ Kivinen and M. Warmuth (1997) and Cesa-Bianchi (1999)

EXPONENTIATED GRADIENT SIMULTANEOUSLY ON ALL CHAIN LINKS

Let Δ_N denote the simplex in \mathbb{R}^N .

Goal: minimize a sequence of multi-variable losses $(\mathbf{u}^{(1)}, \dots, \mathbf{u}^{(K)}) \mapsto \ell_t(\mathbf{u}^{(1)}, \dots, \mathbf{u}^{(K)})$ simultaneously over all variables $(\mathbf{u}^{(1)}, \dots, \mathbf{u}^{(K)}) \in \Delta_{N_1} \times \dots \times \Delta_{N_K}$.

The Multi-variable Exponentiated Gradient forecaster

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):

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

end

Regret bound: If the ℓ_t are jointly convex and differentiable with $\|\nabla_{\mathbf{u}^{(k)}} \ell_t\|_\infty \leq G^{(k)}$, then Multi-variable EG tuned with $\eta^{(k)} = \sqrt{2 \log(N_k)}/n / G^{(k)}$ satisfies:

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

Remember the goal:

$$\inf_{f_0 \in \mathcal{F}^{(0)}} \left\{ \underbrace{\inf_{g_0 \in \mathcal{G}^{(0)}, \dots, g_{K_\varepsilon} \in \mathcal{G}^{(K_\varepsilon)}} \sum_{t=1}^n \left(y_t - (f_0 + g_0 + \dots + g_{K_\varepsilon})(x_t) \right)^2}_{\text{low-scale aggregation}} \right\}$$

low-scale aggregation: Multi-variable EG with loss functions:

$$\ell_t(\mathbf{u}^{(1)}, \dots, \mathbf{u}^{(K_\varepsilon)}) = (y_t - f_0 - \mathbf{u}^{(1)} \cdot g^{(0)} + \dots + \mathbf{u}^{(K_\varepsilon)} \cdot g^{(K_\varepsilon)})^2$$

$$\text{Regret cost: } \sqrt{n} \int_\varepsilon^\gamma \sqrt{\log \mathcal{N}_\infty(\mathcal{F}, \tau)} d\tau$$

high-scale aggregation: run EWA to be competitive against all $f_0 \in \mathcal{F}^{(0)}$

$$\text{Regret cost: } \log(\text{Card} \mathcal{F}^{(0)}) = \log \mathcal{N}_\infty(\mathcal{F}, \gamma)$$

Theorem (G. and Gerchinovitz 2015)

Let $B > 0$, $n \geq 1$, and $\gamma \in (B/n, B)$. Assume that $\max_{1 \leq t \leq n} |y_t| \leq B$ and that $\sup_{f \in \mathcal{F}} \|f\|_\infty \leq B$. Then, Chaining EWA well-tuned (depending on γ , B , and n) satisfies:

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

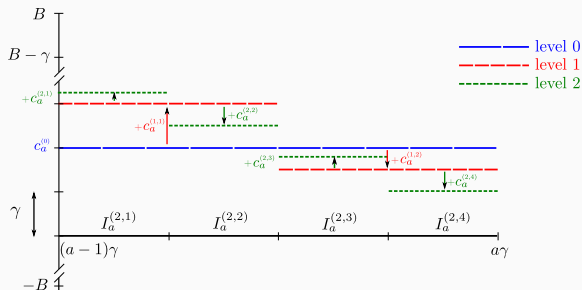
Remarks:

- the bound with $\varepsilon \neq 0$ and εn can also be obtained.
- calibrations in B and n should be possible by **doubling trick** and **clipping**

EFFICIENT IMPLEMENTATION FOR LIPSCHITZ FUNCTIONS

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
 \Rightarrow No need to update all coefficients \Rightarrow **manageable complexity**.

For Hölder functions: piecewise constant \rightarrow **piecewise polynomials**

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

Can be extended to Lipschitz on $[0, 1]^d$ functions at small cost.

We lose a log factor in the regret bound.

B. Guedj and his PhD student are implementing it for numerical experiments.

EXTENSIONS, CURRENT (AND FUTURE) WORK

Goal: minimize the regret $\text{Reg}_n = \sum_{t=1}^n \ell_t(\hat{f}_t) - \inf_{f \in \mathcal{F}} \sum_{t=1}^n \ell_t(f)$ for generic sequences of loss functions (ℓ_t) .

If the loss functions ℓ_t are convex and Lipschitz, we can achieve

$$\text{Reg}_n(\mathcal{F}) \lesssim \underbrace{\log \mathcal{N}_\infty(\mathcal{F}, \gamma)}_{\substack{\text{Large scale term not possible} \\ \text{(was thanks to exp-concavity)}}} + \sqrt{n} \int_\epsilon^1 \sqrt{\log \mathcal{N}_\infty(\mathcal{F}, \tau)} d\tau + \epsilon n$$

Lipschitz class on $[0, 1]^d$	Metric entropy	EWA Regret	Our Regret
$d = 1$	ϵ^{-1}	$n^{2/3}$	$n^{1/2}$
$d = 2$	ϵ^{-2}	$n^{3/4}$	$n^{1/2} \log n$
$d \geq 3$	ϵ^{-d}	$n^{(d+1)/(d+2)}$	$n^{(d-1)/d}$

First **constructive** algorithm to achieve the **optimal**⁶ rates.

The rate $n^{(d+1)/(d+2)}$ was achieved by G. and Baudin (2014) and Hazan and Megiddo (2007).

⁶ A. Rakhlin and K. Sridharan. "Online Nonparametric Regression with General Loss Functions". In: *arXiv* (2015).

Setting: at each step $t \geq 1$,

- simultaneously player plays $\hat{x}_t \in \mathcal{F}$ (possibly random) and environment chooses ℓ_t convex Lipschitz (hidden)
- then player suffers and **observes only** $\ell_t(\hat{x}_t)$

Goal: minimize the expected regret $\text{Reg}_n(\mathcal{F}) \stackrel{\text{def}}{=} \mathbb{E} \left[\sum_{t=1}^n \ell_t(\hat{x}_t) \right] - \inf_{x \in \mathcal{F}} \sum_{t=1}^n \ell_t(x)$

	Full information	Bandits Feedback
(Card $\mathcal{F} = K$) + EWA	$\sqrt{n \log K}$	\sqrt{nK}
ε -net + EWA	$\sqrt{n \log \mathcal{N}_\infty(\mathcal{F}, \varepsilon)} + \varepsilon n$	$\sqrt{n \mathcal{N}_\infty(\mathcal{F}, \varepsilon)} + \varepsilon n$
Chaining	$\sqrt{n} \int_\varepsilon^1 \sqrt{\log \mathcal{N}_\infty(\mathcal{F}, x)} dx + \varepsilon n$	$\sqrt{n} \int_\varepsilon^1 \sqrt{\mathcal{N}_\infty(\mathcal{F}, x)} dx + \varepsilon n?$

For $\mathcal{F} \subset [0, 1]$ this would lead to $\mathcal{O}(\sqrt{n})$ regret \rightarrow first constructive algorithm!

Difficulty: taking advantage of small loss ranges in bandits

Get the sequential entropy $\mathcal{N}^{\text{seq}}(\mathcal{F}, \varepsilon)$ instead of the metric entropy $\mathcal{N}_{\infty}(\mathcal{F}, \varepsilon)$

Efficient version for other function classes

- step-wise Lipschitz functions \rightarrow application to classification
- generalized additive models \rightarrow useful to predict electricity consumption

Similar results with other algorithms (Kernel regression)

THANK YOU !

REFERENCES



N. Cesa-Bianchi. “Analysis of two gradient-based algorithms for on-line regression”. In: *J. Comput. System Sci.* 59.3 (1999), pp. 392–411.



N. Cesa-Bianchi and G. Lugosi. *Prediction, learning, and games*. Cambridge University Press, 2006.



G. and P. Baudin. “A consistent deterministic regression tree for non-parametric prediction of time series”. <http://arxiv.org/abs/1405.1533>. 2014.



G. and S. Gerchinovitz. “A Chaining Algorithm for Online Nonparametric Regression”. In: *Proceedings of COLT’15*. Vol. 40. JMLR: Workshop and Conference Proceedings, 2015, pp. 764–796.



F. Gao, C-K. Ing, and Y. Yang. “Metric entropy and sparse linear approximation of ℓ_q -hulls for $0 < q \leq 1$ ”. In: *Journal of Approximation Theory* 166 (2013), pp. 42–55.



E. Hazan and N. Megiddo. “Online Learning with Prior Knowledge”. In: *Proceedings of the 20th Annual Conference on Learning Theory (COLT’07)*. Ed. by N. H. Bshouty and C. Gentile. Vol. 4539. Lecture Notes in Computer Science. Springer Berlin Heidelberg, 2007, pp. 499–513.

-  J. Kivinen and M.K. Warmuth. “Exponentiated Gradient Versus Gradient Descent for Linear Predictors”. In: *Information and Computation* 132.1 (1997), pp. 1–63.
-  N. Littlestone and M. K. Warmuth. “The Weighted Majority Algorithm”. In: *Information and Computation* 108.2 (1994), pp. 212–261.
-  G.G. Lorentz. “Metric Entropy, Widths, and Superpositions of Functions”. In: *Amer. Math. Monthly* 69.6 (1962), pp. 469–485.
-  A. Rakhlin and K. Sridharan. “Online Nonparametric Regression”. In: *COLT* 35 (2014), pp. 1232–1264.
-  A. Rakhlin and K. Sridharan. “Online Nonparametric Regression with General Loss Functions”. In: *arXiv* (2015).
-  V. Vovk. “Aggregating Strategies.” In: *Proceedings of the Third Workshop on Computational Learning Theory*. 1990, pp. 371–386.
-  V. Vovk. “Competitive on-line statistics”. In: *International Statistical Review* 69.2 (2001), pp. 213–248.