

Talk at CREST : October 5, 2015

Title: stochastic algorithm for online nonparametric regression

- Joint work with Pierre Gillberd
- paper: COLT 2015

I/ Setting: online (nonparametric) regression with individual sequences (7')

Online protocol: at each round  $t \in \mathbb{N}^*$ ,

- get  $x_t \in X$
- predict  $\hat{y}_t \in \mathbb{R}$
- observe  $y_t \in \mathbb{R} \mapsto$  loss:  $(y_t - \hat{y}_t)^2$

Goal: minimize the regret ( $\mathcal{F} \subseteq \mathbb{R}^X$ ):

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

Individual sequences: we want guarantees of the form

$$\sup_{\substack{(x_t)_{t \geq 1} \\ (y_t)_{t \geq 1} : |y_t| \leq B}} \text{Reg}_T(\mathcal{F}) \leq o(T)$$

$$(y_t)_{t \geq 1} : |y_t| \leq B$$

↑  
depends on the size of  $\mathcal{F}$

Ex: • finite  $\mathcal{F}$ :  $\sup_{(x_t), (y_t)} \text{Reg}_T(\mathcal{F}) \lesssim \log |\mathcal{F}|$

•  $\mathcal{F} := \left\{ \sum_{j=1}^d u_j \varphi_j : u \in \Delta_d \right\}$ :  $\sup \text{Reg}_T(\mathcal{F}) \lesssim d \log T$

•  $\mathcal{F} := \{f: [0,1] \rightarrow [-B,B] \text{ 1-Lip}\}$ :  $\sup \text{Reg}_T(\mathcal{F}) \lesssim T^{1/3}$

## II / Simplest case: finite $\mathcal{F}$ (5')

Algorithm: Exponentially Weighted Average forecaster (EWA)

$$\forall t \geq 1 \quad \hat{w}_{t,j} = \frac{\exp\left(-\eta \sum_{s=1}^{t-1} (y_s - f_j(x_s))^2\right)}{\sum_i \dots}, \quad 1 \leq j \leq |\mathcal{F}|$$

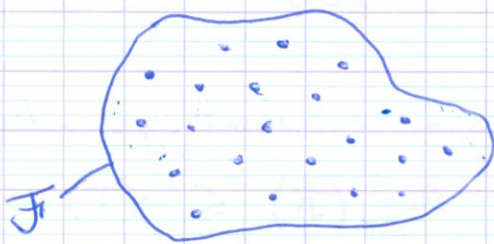
$$\hat{y}_t := \sum_{j=1}^{|\mathcal{F}|} \hat{w}_{t,j} f_j(x_t)$$

Thm (Kivinen and Warmuth 1995)

$$\left| \text{Reg}_{\mathcal{F}}(\mathcal{F}) \leq 8B^2 \log |\mathcal{F}| \right. \quad \text{for all } (x_t) \text{ and } (y_t) \\ \left. \text{s.t. } |y_t| \leq B \text{ and } |f_j(x_t)| \leq B. \right. \\ \left. \text{if } \eta := \frac{1}{8B^2}. \right.$$

## III / What if $\mathcal{F}$ is (uncountably) infinite? (20')

1) Natural idea: discretize



$\mathcal{F}^{(\epsilon)}$ : finite set s.t.

$$\forall f \in \mathcal{F}, \exists g \in \mathcal{F}^{(\epsilon)}:$$

$$\|f - g\|_{\infty} \leq \epsilon.$$

(we assume  $\|f\|_{\infty} \leq B \quad \forall f \in \mathcal{F}$ ,  
and actually:  $\mathcal{F}$  totally bounded).

Then, apply EWA to points in  $\mathcal{F}^{(\epsilon)}$ :

$$\sum_{t=1}^T (y_t - \hat{y}_t)^2 \leq \min_{f \in \mathcal{F}^{(\epsilon)}} \sum_{t=1}^T (y_t - f(x_t))^2 + 8B^2 \log |\mathcal{F}^{(\epsilon)}|$$

$$\leq \inf_{f \in \mathcal{F}} \sum_{t=1}^T (y_t - f(x_t))^2 + T\epsilon^2 + \underbrace{2T \cdot 2B\epsilon}$$

$$+ \underbrace{8B^2 \log |\mathcal{F}^{(\epsilon)}|}$$

$$\Rightarrow \text{Reg}_T(\mathcal{F}) \lesssim T\varepsilon + \log |\mathcal{F}^{(\varepsilon)}| \quad (*)$$

[introduce metric entropy here]

Ex:  $\mathcal{F} := \left\{ \sum_{j=1}^d \mu_j \varphi_j : \mu \in \Delta_d \right\} \Rightarrow \log |\mathcal{F}^{(\varepsilon)}| \lesssim d \log \frac{1}{\varepsilon}$

$$\begin{aligned} \text{Reg}_T(\mathcal{F}) &\lesssim T\varepsilon + d \log \frac{1}{\varepsilon} \\ &\lesssim d \log T \quad \text{for } \varepsilon \approx \frac{1}{T} \\ &\quad \text{Optimal!} \end{aligned}$$

### 2) Suboptimal bound if $\mathcal{F}$ is too large

Assume that  $\log \mathcal{N}_\infty^p(\mathcal{F}, \varepsilon) \approx \left(\frac{1}{\varepsilon}\right)^p$  (instead of  $d \log \frac{1}{\varepsilon}$ )

$\uparrow$  smallest possible = "metric entropy"       $\uparrow$  large! "nonparametric"

(\*) yields:  $\text{Reg}_T(\mathcal{F}) \lesssim T\varepsilon + \left(\frac{1}{\varepsilon}\right)^p$

$$\lesssim \underbrace{T^{\frac{p}{p+1}}}_{\text{suboptimal!}} \text{ for } \varepsilon \approx T^{-\frac{1}{p+1}}$$

Indeed:

Barkhtin and Lichner (COLT'14) proved that  $\text{Reg}_T(\mathcal{F})$  can be as small as:

Dudley-type regret bound

$$\boxed{\text{Reg}_T(\mathcal{F}) \lesssim \log \mathcal{N}_\infty^p(\mathcal{F}, \gamma) + \sqrt{T} \int_0^\gamma \sqrt{\log \mathcal{N}_\infty^p(\mathcal{F}, \varepsilon)} d\varepsilon}$$

$$\begin{aligned} &\lesssim \gamma^{-p} + \sqrt{T} \int_0^\gamma \varepsilon^{-p/2} d\varepsilon \\ \log \mathcal{N}_\infty^p(\mathcal{F}, \varepsilon) \approx \left(\frac{1}{\varepsilon}\right)^p &\lesssim \gamma^{-p} + \sqrt{T} \cdot \gamma^{1-p/2} \end{aligned}$$

(if  $p < 2$ )

$$\lesssim \underbrace{T^{\frac{p}{p+2}}}_{\text{optimal!}} \text{ for } \gamma \approx T^{-\frac{1}{p+2}}$$

(if  $p \in (0, 2)$ )

( $T^{1-1/p}$  if  $p > 2$ , with a slightly different Dudley bound)

Condition:  $\mathcal{F} := \left\{ f: [0,1] \rightarrow \mathbb{R} : \|f^{(k)}\|_{\infty} \leq B \ \forall k=0, \dots, p, \right.$   
 $\left. |f^{(p)}(x) - f^{(p)}(y)| \leq L|x-y|^{\alpha} \right\}$   
 $p \in \mathbb{N}, \alpha \in (0,1]$ .

Assume  $\beta := p + \alpha > 1/2$

$$\log \mathcal{N}_{\infty}(\mathcal{F}, \varepsilon) \asymp \left(\frac{1}{\varepsilon}\right)^{1/\beta} \implies p = \frac{1}{\beta}$$

$$\implies \text{Reg}_T(\mathcal{F}) \lesssim T^{\frac{p}{p+2}} = T^{\frac{1/\beta}{1/\beta+2}} = T^{\frac{1}{2\beta+1}}$$

$$\implies \frac{\text{Reg}_T(\mathcal{F})}{T} \lesssim T^{-\frac{2\beta}{2\beta+1}} \quad \boxed{\text{as for iid data!}}$$

(or  $T^{-\beta}$  if  $\beta \leq 2$ , as for iid data)

BUT: Bakker et al (2014): non-constructive regret bound  
 (no explicit algo).

Our contribution: explicit algo. (with metric entropy instead of sequential entropy)

#### IV / Linearizing the square loss on help belly (10')

EWA may be suboptimal; instead: EG:

Algo: Exponentiated Gradient (EG)

- \* loss functions:  $l_t(u)$ ,  $u \in \Delta_N$ , convex and diff.
- \* For all  $t \geq 1$ ,

$$\hat{u}_{t,i} = \frac{\exp\left(-\eta \sum_{s=1}^{t-1} \partial_{u_i} l_s(\hat{u}_s)\right)}{\sum_i \exp\left(-\eta \sum_{s=1}^{t-1} \partial_{u_i} l_s(\hat{u}_s)\right)}, \quad 1 \leq i \leq N$$

Regret bound (KW'93, CB'99)

Assume  $l_t$  convex, diff, and  $\|\nabla l_t\|_{\infty} \leq G$

Then:

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

for  $\eta = \frac{1}{G} \sqrt{\frac{2 \ln N}{T}}$

Application:  $l_t(u) := \left( y_{t_k} - f_0(x_t) - \sum_{j=1}^N u_j g_j(x_t) \right)^2, u \in \Delta_N$

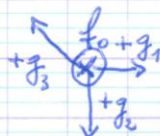
$$\sum_{t=1}^T \left( y_{t_k} - f_0(x_t) - \underbrace{\sum_{j=1}^N \hat{u}_j g_j(x_t)}_{=: \hat{f}_t} \right)^2 \leq \min_{1 \leq j \leq N} \sum_{t=1}^T \left( y_{t_k} - f_0(x_t) - g_j(x_t) \right)^2 + G \sqrt{2T \log N}$$

$$\leq B \max_j \|g_j\|_\infty \sqrt{2T \log N}$$

NB:  $G \sqrt{2T \log N} \ll B^2 \log N \iff G \ll B^2 \sqrt{\frac{\log N}{T}}$

EG "better" than EWA if  $\max_j \|g_j\|_\infty \ll B \sqrt{\frac{\log N}{T}}$

CL: linearizing the square loss can help if functions in  $\mathcal{F}$  are close.



(10')

V / Main algo: the Chaining Exponentially Weighted Average forecast

Chaining technique: approximate any  $f \in \mathcal{F}$  via  $\pi_0(f) \in \mathcal{F}^{(0)}, \pi_1(f) \in \mathcal{F}^{(1)}, \dots$

$$\sup_f \|\pi_k(f) - f\|_\infty \leq \frac{\delta}{2^k}, \quad |\mathcal{F}^{(k)}| = \mathcal{N}_\infty(\mathcal{F}, \delta/2^k)$$

Recursion:

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

Algo: two levels of aggregation in parallel:

- low-scale aggregation: new multi-scale version of EG to be competitive with all increments  $\pi_{k+1}(f) - \pi_k(f)$ ,  $k \geq 1$ , simultaneously w.r.t. one instance around each  $\pi_0(f) \in \mathcal{F}^{(0)}$ .



$\implies$  regret contribution  $\approx \sum_{k=1}^{+\infty} \frac{\delta B}{2^k} \sqrt{T \log \mathcal{N}_\infty(\mathcal{F}, \delta/2^k)}$

$\approx B \sqrt{T} \int_0^\delta \sqrt{\log \mathcal{N}_\infty(\mathcal{F}, \epsilon)} d\epsilon$

- high-scale aggregation: EWA to aggregate the resulting  $\mathcal{N}_\infty(\mathcal{F}, \delta)$  forecasts
- $\implies$  regret contribution  $\approx B^2 \log \mathcal{N}_\infty(\mathcal{F}, \delta)$
- $\implies$  Dudley type regret-bound

## References / previous work (10')

Cesa-Bianchi and Lugosi (2001): two-side aggregation  
 "Worst-case bounds for the  $\log$  loss of predictors"  $\log$  loss, online setting  
 chaining in the analysis

Balkekin et al. (to appear): square loss, batch setting  
 "Empirical entropy, minimax regret and  
 minimax risk" two-side aggregation  
 chaining in the analysis

Balkekin et al. (COLT'14): square loss, online  
 "Online nonparametric regression" chaining in the analysis  
 (but stronger notion of entropy)

Cesa-Bianchi (1999): absolute loss, online  
 "On prediction of individual seq" chaining in the algo.

## Second contribution: efficient algo for Hölder classes

function class	regret bound	time complexity	space complexity
Lipschitz on $[0,1]$	$T^{1/3} \log(T)$	$T^{4/3} \log T$	$T^{4/3} \log T$
$\beta$ -Hölder on $[0,1]$	$T^{\frac{1}{\beta+1}} \log^{3/2} T$	$\text{poly}(T)$	$\text{poly}(T)$

Main idea: dyadic discretization to only update few weights at each round.