

(joint with N. Gao-Bianchi, P. Gaillard, C. Gentile)

We will focus on how chaining (probabilistic tool) can be used to construct algorithms for nonparametric online learning of individual (non-stochastic) sequences.

1. "Classical chaining": a brief recap

Technique to control $\max_{t \in \mathcal{G}} X_t$ where (\mathcal{G}, d) : metric space

- $\mathbb{E}(X_t) = 0$

1.1. Small \mathcal{G}

- small increments $X_t - X_s$ when $d(t, s)$ small

Lemma: Let X_1, \dots, X_N be real r.v. such that: $\exists \nu \geq 0$,
 $\forall i, \forall \lambda \in \mathbb{R}, \mathbb{E}[e^{\lambda X_i}] \leq e^{\frac{\lambda^2 \nu}{2}}$ (subgaussianity)

Then, $\mathbb{E}\left[\max_{1 \leq i \leq N} X_i\right] \leq \sqrt{2\nu \log N}$

NB: for $X_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \nu)$,
 $\mathbb{E}\left[\max_{1 \leq i \leq N} X_i\right] \sim \sqrt{2\nu \log N}$
 $N \rightarrow +\infty$

Proof (Pisier-like argument):

$$\begin{aligned} \mathbb{E}\left[\max_{1 \leq i \leq N} X_i\right] &= \frac{1}{\lambda} \log \exp\left(\lambda \mathbb{E}\left[\max_{1 \leq i \leq N} X_i\right]\right) \quad \text{for any } \lambda > 0 \\ &\leq \frac{1}{\lambda} \log \mathbb{E}\left[\max_{1 \leq i \leq N} e^{\lambda X_i}\right] \quad \text{by Jensen's inequality} \\ &\leq \frac{1}{\lambda} \log\left(\sum_{i=1}^N \mathbb{E}\left[e^{\lambda X_i}\right]\right) \\ &\leq \frac{\log N}{\lambda} + \frac{\lambda \nu}{2} = \sqrt{2\nu \log N} \quad \text{for } \lambda := \sqrt{\frac{2 \log N}{\nu}} \end{aligned}$$

1.2. What about $\mathbb{E}[\max_{t \in \mathcal{G}} X_t]$ for finite but large \mathcal{G} ?

Assume that:

- (i) \mathcal{G} is finite (for simplicity)
- (ii) $\mathbb{E}(X_t) = 0 \quad \forall t \in \mathcal{G}$
- (iii) the increments are subgaussian:
 $\forall s, t \in \mathcal{G}, \forall \lambda \in \mathbb{R}, \mathbb{E}[e^{\lambda(X_t - X_s)}] \leq e^{\frac{\lambda^2}{2} d^2(s, t)}$

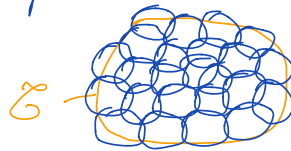
Letting $B := \max_{s, t \in \mathcal{G}} d(s, t)$, we can use the above lemma with $\sigma = B^2$ and get

$$\mathbb{E}\left[\max_{t \in \mathcal{G}} X_t\right] = \mathbb{E}\left[\max_{t \in \mathcal{G}} (X_t - X_{t_0})\right] \leq \sqrt{2B^2 \log |\mathcal{G}|}$$

Issue: This upper bound is usually very pessimistic!

Indeed, (iii) implies that $\text{Var}(X_t - X_s) \leq d^2(t, s) \rightarrow$ that X_t and X_s are highly correlated when $d(t, s)$ is small \rightarrow the cardinality $|\mathcal{G}|$ can be replaced with a much smaller number?

For $\varepsilon > 0$, define



$$\bigcup_{i=1}^N \overline{B}(t_i, \varepsilon) = \mathcal{G}$$

- ε -covering: any finite set $\{t_1, \dots, t_N\}$ s.t. the closed balls $\overline{B}(t_i, \varepsilon)$ cover \mathcal{G}
- $N(\mathcal{G}, d, \varepsilon)$: the smallest cardinality N of an ε -covering (∞ if none)

$\log N(\mathcal{G}, d, \varepsilon)$ is called "metric entropy" at scale ε . Quantifies richness of (\mathcal{G}, d) .

Ex $\log N([0, 1]^d, \|\cdot\|, \varepsilon) \approx \log \left(\frac{1}{\varepsilon}\right)^d \approx d \log \frac{1}{\varepsilon}$

Proposition (Dudley '67). *NB: there exist several refinements/extensions.*

Assume (i), (ii), and (iii) above. Then, setting $B := \max_{s, t \in \mathcal{G}} d(s, t)$,

$$\mathbb{E} \left[\max_{t \in \mathcal{G}} X_t \right] \leq 12 \int_0^{\frac{B}{2}} \sqrt{\log \mathcal{N}(\mathcal{G}, d, \varepsilon)} d\varepsilon$$

Proof: $\varepsilon_m := \frac{B}{2^m}$ \rightarrow discretization $\mathcal{G}_m \subseteq \mathcal{G}$, $|\mathcal{G}_m| = \mathcal{N}(\mathcal{G}, d, \varepsilon_m)$
 s.t. $\mathcal{G} = \bigcup_{t \in \mathcal{G}_m} \overline{B}(t, \varepsilon_m)$
 centers of balls of radius ε_m .

$$\mathcal{G}_0 = \{t_0\} \subseteq \mathcal{G}_1 \subseteq \dots \subseteq \mathcal{G}_M = \mathcal{G}.$$

We set $\pi_m(t) \in \operatorname{argmin}_{s \in \mathcal{G}_m} d(s, t)$. (approximation of t at scale ε_m)

$$\mathbb{E} \left[\max_{t \in \mathcal{G}} X_t \right] = \mathbb{E} \left[\max_{t \in \mathcal{G}} (X_t - X_{t_0}) \right]$$

$$= \mathbb{E} \left[\max_{t \in \mathcal{G}} \sum_{m=0}^{M-1} (X_{\pi_{m+1}(t)} - X_{\pi_m(t)}) \right]$$

$$\leq \sum_{m=0}^{M-1} \mathbb{E} \left[\max_{t \in \mathcal{G}} (X_{\pi_{m+1}(t)} - X_{\pi_m(t)}) \right]$$

$$\stackrel{\text{lemma}}{\leq} \sum_{m=0}^{M-1} \sqrt{2 B_m^2 \log(|\mathcal{F}_m| - |\mathcal{F}_{m+1}|)}$$

$$\text{with } B_m = \max_{t \in \mathcal{G}} d(\pi_{m+1}(t), \pi_m(t))$$

$$\leq \frac{B}{2^{m+1}} + \frac{B}{2^m} = \frac{3B}{2^{m+1}} \text{ by the triangle inequality}$$

$$\leq 6 \sum_{m=0}^{M-1} \frac{B}{2^{m+1}} \sqrt{\log \mathcal{N}(\mathcal{G}, d, \frac{B}{2^{m+1}})}$$

$$\leq 12 \sum_{m=0}^{M-1} \int_{\frac{B}{2^{m+2}}}^{\frac{B}{2^{m+1}}} \sqrt{\log \mathcal{N}(\hat{\sigma}_t, d, \varepsilon)} d\varepsilon$$

$$\leq 12 \int_0^{\frac{B}{2}} \sqrt{\log \mathcal{N}(\hat{\sigma}_t, d, \varepsilon)} d\varepsilon \quad \square$$

2. Chaining in (nonparametric) online learning

2.1. Online learning problem:

Sequence $(x_1, y_1), (x_2, y_2), \dots$ which is arbitrary.
 At any round $t \geq 1$:

- We observe $x_t \in \mathcal{X}$
- We predict $\hat{y}_t \in \mathcal{Y}$
- We observe $y_t \in \mathcal{Y}$ and suffer loss $l(\hat{y}_t, y_t)$
 for some loss function $l: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_+$.

Goal of the learner: minimize the regret

$$R_T(\mathcal{F}) := \sum_{t=1}^T l(\hat{y}_t, y_t) - \inf_{f \in \mathcal{F}} \sum_{t=1}^T l(f(x_t), y_t)$$

for some large (nonparametric) function set $\mathcal{F} \subseteq \mathcal{Y}^{\mathcal{X}}$,
 and for all sequences $(x_t, y_t)_{t \in [1, T]}$.

We have an analogue to Section 1, as shown below.

2.2. Simple case: \mathcal{F} is finite (and small)

Lemma: Assume \mathcal{F} is finite

$$\bullet |l(f(x), y) - l(g(x), y)| \leq B \quad \forall f, g \in \mathcal{F}, \forall x, y$$

Then there exists a randomized algorithm that chooses the \hat{f}_T in such a way that

$$\forall (x_t, y_t)_{t \leq T} \in (\mathcal{X} \times \mathcal{Y})^T, \quad \mathbb{E}[R_T(\mathcal{F})] \stackrel{(\text{or w.h.p.})}{\leq} B \sqrt{\frac{T}{2} \log |\mathcal{F}|}$$

- Also true if the elements of \mathcal{F} are themselves randomized algorithms.
- A famous algorithm doing so is the "Exponential Weights" algorithm (or "Hedge"): $\hat{f}_T \sim \sum_{f \in \mathcal{F}} p_t(f) f(x_t)$

$$\text{where } p_t(f) = \frac{\exp(-\gamma \sum_{s=1}^{t-1} l(f(x_s), y_s))}{\sum_{g \in \mathcal{F}} \exp(-\gamma \sum_{s=1}^{t-1} l(g(x_s), y_s))}$$

Proof also uses subgaussianity (of $l(\hat{f}_t, y)$ given the past).

- The lemma of Section 1 can also be used directly by reducing the learning problem to the estimation of $\mathbb{E}[\max_f X_f]$ for some well-chosen $(X_f)_{f \in \mathcal{F}}$.

[see Bakhlin, Tishbani, and Gensari, JMLR 2015]

2.3. Nonparametric case: \mathcal{F} is large

$(\mathcal{Y}, d) \mapsto (\mathcal{F}, d_\infty)$ where $d_\infty(f, g) := \sup_{x \in \mathcal{X}} d(f(x), g(x))$

$N_\infty(\mathcal{F}, \varepsilon) := N(\mathcal{F}, d_\infty, \varepsilon)$ = minimal number of d_∞ -balls to cover \mathcal{F} at scale ε .

Proposition: Assume that

- $\forall y \in \mathcal{Y}$, $l(\cdot, y)$ is L -Lipschitz
- $N_\infty(\mathcal{F}, \varepsilon) < +\infty \quad \forall \varepsilon > 0$

Let $M \geq 1$. There exists a (computationally horrible) randomized algorithm such that, setting $B := \sup_{f, g \in \mathcal{F}} d_\infty(f, g)$,

$$\forall (x_t, y_t)_{t \in \{1, \dots, T\}}, \quad \mathbb{E} \left[R_T(\mathcal{F}) \right] \leq 6L \sqrt{2T} \int_{B2^{-M-1}}^{\frac{B}{2}} \sqrt{\log N_\infty(\mathcal{F}, \varepsilon)} \, d\varepsilon + TLB2^{-M}$$

Alg: uses chaining in its construction

Partially inspired from Geo-Bianchi and Lugosi (1999) that was specific to $l(y, y') = |y - y'|$.

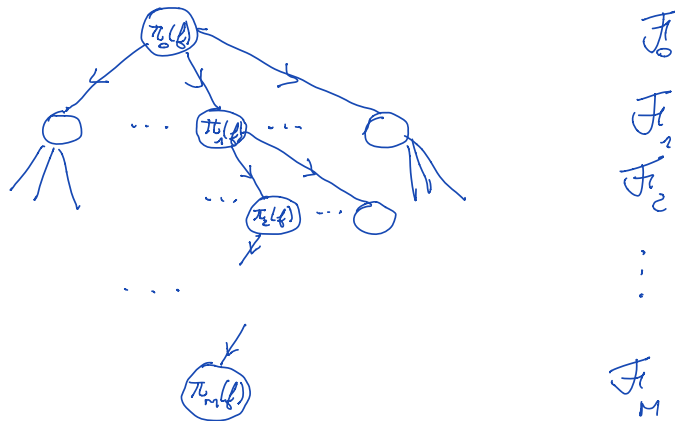
- Hierarchical discretization of (\mathcal{F}, d_∞) :

$$\mathcal{F}_0 = \{f_0\}$$

\mathcal{F}_m = smallest discretization of \mathcal{F} at scale $B/2^m$, $m \geq 1$ with cardinality $N_m := N(\mathcal{F}, d_\infty, B2^{-m})$

Level- m approximation: $\pi_m(f) \in \arg \min_{g \in \mathcal{F}_m} d_\infty(g, f)$

We build a Directed Acyclic Graph (NB: we could also build a tree), where level- m nodes are labelled with elements of \mathcal{F}_m , and where we connect $\pi_m(f)$ with $\pi_{m+1}(f)$ for all $f \in \mathcal{F}$ and $m = 0, \dots, M-1$.



→ At time t , each leaf $f \in \mathcal{F}_M$ recommends to play $f(x_t)$

→ On each internal node $f_m \in \mathcal{F}_m$, we place an instance $A_m(f_m)$ of the exponential weighting algorithm that learns the best of its children.

→ We play according to the algorithm sitting on the root.

Proof of regret bound:

Fix $f \in \mathcal{F}$. The regret against f can be rewritten as a sum of M regrets along the path $A_0(\pi_0(f)) \rightarrow A_1(\pi_1(f)) \rightarrow \dots \rightarrow A_M(\pi_M(f))$

$$\begin{aligned} & \mathbb{E} \left[\sum_{t=1}^T \ell(\hat{f}_t, y_t) - \sum_{t=1}^T \ell(f(x_t), y_t) \right] \\ & \leq \sum_{m=0}^{M-1} \mathbb{E} \left[\sum_{t=1}^T \ell(A_{m,t}(\pi_m(f)), y_t) - \sum_{t=1}^T \ell(A_{m,t}(\pi_{m+1}(f)), y_t) \right] + TLB2^{-M} \\ & \leq \sum_{m=0}^{M-1} B_m \sqrt{\frac{T}{2} \log N_{m+1}} + TLB2^{-M} \end{aligned}$$

where $B_m = \text{range of the bases of children of level-}m \text{ nodes}$
 $\leq L \times 2 \times 3B2^{-(m+1)}$

and $\log N_{m+1} = \log |F_{m+1}| = \log \mathcal{N}_\infty(F, B2^{-m-1})$

$$\begin{aligned} & \leq 6L \sqrt{2T} \sum_{m=0}^{M-1} B2^{-(m+2)} \sqrt{\log \mathcal{N}_\infty(F, B2^{-(m+1)})} + TLB2^{-M} \\ & \leq 6L \sqrt{2T} \int_{\frac{B}{2^{M+1}}}^{\frac{B}{2}} \sqrt{\log \mathcal{N}_\infty(F, \varepsilon)} d\varepsilon + TLB2^{-M} \quad \square \end{aligned}$$

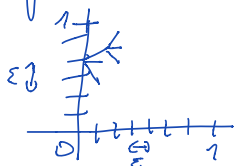
6.4. Example: Lipschitz functions

$$\mathcal{X} = [0, 1]^d, \quad \mathcal{Y} = \mathbb{R},$$

$$F = \left\{ f: [0, 1]^d \rightarrow \mathbb{R} : |f(x) - f(x')| \leq \|x - x'\| \quad \forall x, x' \in [0, 1]^d \right\}$$

↑
any norm

$$\log \mathcal{N}_\infty(F, \varepsilon) \asymp \left(\frac{1}{\varepsilon}\right)^d$$



$\frac{1}{\varepsilon} \times 3^{1/2}$ functions $\Rightarrow \log \mathcal{N}_\infty(F, \varepsilon) \asymp \frac{1}{\varepsilon}$ when $d=1$.

Plugging $\log N_{\infty}(F, \epsilon) \asymp \left(\frac{1}{\epsilon}\right)^d$ into the regret bound and taking $M = \left\lceil \frac{1}{d} \log_2 T \right\rceil$ (so that $\epsilon^{-M} \approx T^{-1/d}$), we get:

$$\mathbb{E}[R_T] \lesssim L \sqrt{T} \int_{BT^{-1/d}}^B \sqrt{\left(\frac{1}{\epsilon}\right)^d} d\epsilon + T^{1-1/d} \approx T^{1-1/d}$$

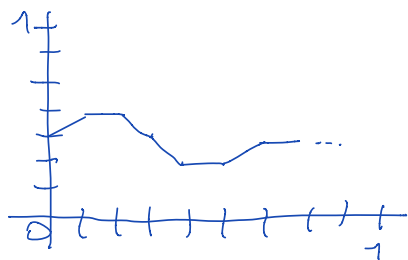
which matches the lower bound of Hazan and Megretski (2007).

NB: faster rates are possible under stronger assumptions on $\ell(\cdot; y)$ (e.g., faster rates with square loss or any exp-concave loss)

Efficient variants of the algorithm (when F = set of bounded and Lipschitz functions on $[0,1]^d$)

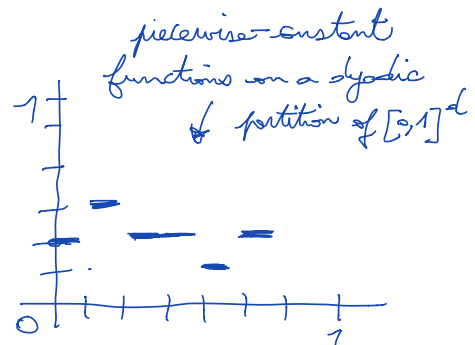
- The previous "exp" was super inefficient: requires $\exp(\text{poly}(T))$ updates at any round $t \geq 1$.
- Fortunately, there exist coverings F_m that are much more manageable from a computational viewpoint, while being almost optimal (up to log factors) from a statistical viewpoint.

Main idea in dimension $d=1$:



optimal covering
but computationally
inefficient

replaced
with



additional log factor in regret
but polynomial time algo.

independent of dimension d .