

Nonparametric Learning, Invariance, and Sketching Techniques for Online Algorithms

Nicolò Cesa-Bianchi

Università degli Studi di Milano

Joint work with:

- Rocco de Rosa (Roma) and Francesco Orabona (Yahoo NYC)
— nonparametric learning
- Alekh Agarwal (MSR-NYC), John Langford (MSR-NYC)
and Haipeng Luo (Princeton) — invariance and sketching



Online randomized classification in metric spaces

Setting

- Metric space (\mathcal{X}, ρ)
- Deterministic source $\sigma = ((x_1, y_1), (x_2, y_2), \dots)$
- Where $(x_t, y_t) \in \mathcal{X} \times \{0, 1\}$

For $t = 1, 2, \dots$

- 1 Receive next instance $x_t \in \mathcal{X}$
- 2 Output randomized prediction $\hat{Y}_t \in \{0, 1\}$ where $\mathbb{P}(\hat{Y}_t = 1) = p_t$
- 3 Observe true class label $y_t \in \{0, 1\}$
- 4 Pay loss $\mathbb{I}\{\hat{Y}_t \neq y_t\}$ with probability $\mathbb{P}(\hat{Y}_t \neq y_t) = |p_t - y_t|$



Nonparametric online learning

A nonparametric class

- $f : \mathcal{X} \rightarrow \mathbb{R}$ is **Lipschitz** if there exists $0 \leq L < \infty$ such that
$$|f(x) - f(x')| \leq L \rho(x, x') \quad x, x' \in \mathcal{X}$$
- $\text{Lip} = \text{Lip}(\mathcal{X})$ is the class of $[0, 1]$ -valued Lipschitz functions on \mathcal{X}
- How fast can we learn the best $f \in \text{Lip}$ for an arbitrary σ ?

Expected regret

$$R_T(\sigma) = \sum_{t=1}^T \ell_t(p_t) - \inf_{f \in \text{Lip}} \sum_{t=1}^T \ell_t(f(x_t))$$

where $\ell_t(p) = |p - y_t|$



Doubling dimension of metric spaces

Doubling dimension

A metric space has **doubling dimension** d if $\Theta(\varepsilon^{-d})$ nonoverlapping balls of radius ε can be packed in it

Nonparametric rates: k -NN

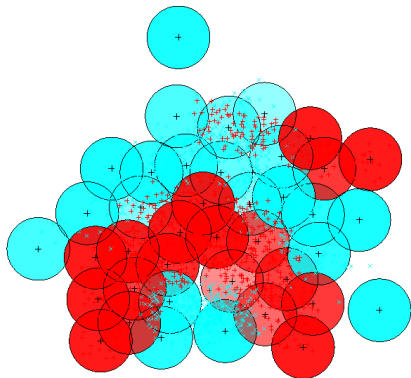
- **Stochastic source:** $X_1, X_2, \dots \in \mathcal{X}$ i.i.d.
- $\mathbb{P}(Y_t = 1 \mid X_t = x_t) = f^*(x_t)$ for some $f^* \in \text{Lip}$
- Then the excess expected mistake rate of k_T -NN w.r.t. $\mathbb{I}\{f^*(x) \geq \frac{1}{2}\}$ grows at optimal rate bounded by

$$\Theta\left(T^{\frac{d+1}{d+2}}\right)$$

where d is the doubling dimension of \mathcal{X}

- Can we match this rate for any deterministic source?

- The instance space is incrementally covered with “half-overlapping” balls
- Each ball has a local classifier sitting in it
- For prediction, local classifier is selected with a NN query
- For updating, ball radius determines relevant points
- Radius of each ball shrinks with time



Basic algorithm

Create a default cover \mathcal{C} containing a single ball with a default local classifier

For $t = 1, 2, \dots$

- 1 Receive next stream instance x_t and compute closest center $x_t^* \in \mathcal{C}$
- 2 Predict using local classifier sitting at x_t^*
- 3 Observe true class label y_t and update local classifier
- 4 If x_t is not contained in ball centered at x_t^* then add ball $B(x_t, \varepsilon_t)$ to \mathcal{C}
- 5 Shrink radii ε_t of all balls in \mathcal{C} according to law $\varepsilon_t = t^{-\frac{1}{d+2}}$

We use **majority voters** as local classifiers



Analysis (part I)

- Fix any Lipschitz $f : \mathcal{X} \rightarrow [0, 1]$ with Lipschitz constant L_f
- $p_s^{\text{opt}} \in \{0, 1\}$ is best prediction for all (x_t, y_t) such that $x_t \in \mathcal{B}(x_s, \varepsilon_t)$

$$\begin{aligned} R_T(\sigma) &= \sum_{x_s \in \mathcal{C}} \sum_{t: x_s = x_t^*} \left(\ell_t(p_t) - \ell_t(f(x_t)) \right) \\ &\leq \sum_{x_s \in \mathcal{C}} \sum_{t: x_s = x_t^*} \left(\ell_t(p_t) - \ell_t(p_s^{\text{opt}}) \right) + \sum_{x_s \in \mathcal{C}} \sum_{t: x_s = x_t^*} \left(\ell_t(f(x_s)) - \ell_t(f(x_t)) \right) \\ &\leq \sum_{x_s \in \mathcal{C}} \sqrt{T_s} + \sum_{x_s \in \mathcal{C}} \sum_{t: x_s = x_t^*} L_f \varepsilon_t \\ &\leq \sqrt{|\mathcal{C}|T} + L_f \sum_{t=1}^T t^{-\frac{1}{d+2}} \leq \sqrt{|\mathcal{C}|T} + L_f T^{\frac{d+1}{d+2}} \end{aligned}$$



Analysis (part II)

Since $\varepsilon_T = T^{-\frac{1}{d+2}}$, the packing number bound gives $|\mathcal{C}| \leq C_X T^{\frac{d}{d+2}}$

Therefore

$$\begin{aligned} R_T(\sigma) &\leq \sqrt{|\mathcal{C}|T} + L_f T^{\frac{d+1}{d+2}} \\ &\leq C_X \sqrt{T^{\frac{d}{d+2}}T} + L_f T^{\frac{d+1}{d+2}} \\ &\leq (C_X + L_f) T^{\frac{d+1}{d+2}} \end{aligned}$$



Regret bound

Theorem

For any deterministic source $\sigma = ((x_1, y_1), (x_2, y_2), \dots)$

$$\sum_{t=1}^T \ell_t(p_t) \leq \min_{f \in \text{Lip}} \left(\sum_{t=1}^T \ell_t(f(x_t)) + (C_X + L_f) T^{\frac{d+1}{d+2}} \right)$$

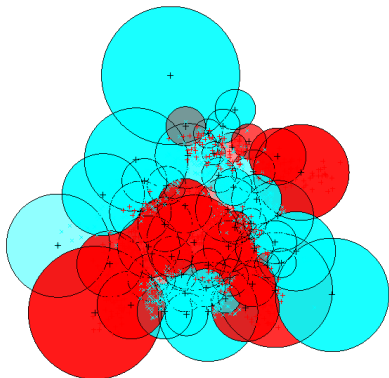
where d is the doubling dimension of \mathcal{X}

Remarks:

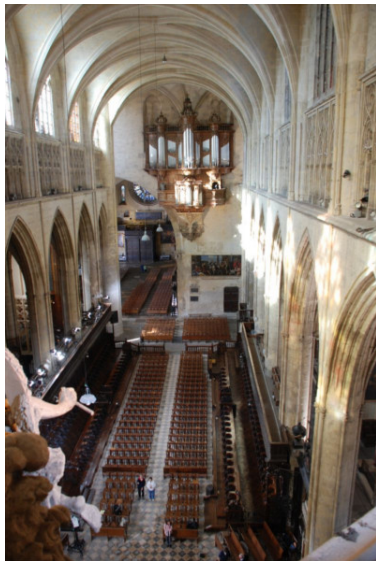
- **Oracle inequality:** best trade-off between fitness and complexity
- Doubling trick to learn the doubling dimension
- Easy to get **parametric rate** \sqrt{T} against any f such that $L_f = 0$
- Scaling $L_f T^{\frac{d+1}{d+2}}$ possible, but not uniformly over Lip

Variant: radius shrinks with mistake rate

- $\varepsilon_t = m_t(s)^{-1/(d+2)}$ where $m_t(s)$ is number of mistakes of classifier sitting at $x_s \in \mathcal{C}$
- Cover adapts to **local complexity** of data source
- We lose “half-overlapping property” and packing argument is gone



Entr'Acte



Online linear prediction

Setting

- Data source $\mathbf{x}_1, \mathbf{x}_2, \dots \in \mathbb{R}^d$
- Convex and differentiable loss sequence f_1, f_2, \dots on \mathbb{R}

For $t = 1, 2, \dots$

- 1 Receive next instance $\mathbf{x}_t \in \mathbb{R}^d$
- 2 Output linear prediction $\mathbf{p}_t = \mathbf{w}_t^\top \mathbf{x}_t$
- 3 Pay loss $\ell_t(\mathbf{w}_t) = f_t(\mathbf{w}_t^\top \mathbf{x}_t)$
- 4 Observe gradient $\mathbf{g}_t = \nabla \ell_t(\mathbf{w}_t) = f'_t(\mathbf{p}_t) \mathbf{x}_t$



Regret

$$R_T(\mathbf{w}) = \sum_{t=1}^T \ell_t(\mathbf{w}_t) - \sum_{t=1}^T \ell_t(\mathbf{w})$$

- We know that for arbitrary convex losses $R_T(\mathbf{w})$ scales with \sqrt{T}
- Dependence on \mathbf{w} , $\mathbf{x}_1, \dots, \mathbf{x}_T$, and f_1, \dots, f_T ?



First-order (Euclidean) methods

Online gradient descent

[Zinkevich, 2003]

$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t \mathbf{g}_t$ followed by projection to ball of radius U

$$R_T(\mathbf{w}) \leq (G_C U X) \sqrt{T}$$

for all

- \mathbf{w} with $\|\mathbf{w}\| \leq U$
- $\mathbf{x}_1, \dots, \mathbf{x}_T$ with $\max_t \|\mathbf{x}_t\| \leq X$
- f_1, \dots, f_T with $\max_t \max_{\mathbf{p}: \|\mathbf{p}\| \leq C} |f'_t(\mathbf{p})| \leq G_C$



Invariance to linear transformations

A natural comparison class

$$\mathcal{K}_C = \left\{ \mathbf{w} \in \mathbb{R}^d : |\mathbf{w}^\top \mathbf{x}_t| \leq C, t = 1, \dots, T \right\}$$

- Invariant to any (invertible) linear transformation of data sequence
- More general than $\|\mathbf{x}_t\| \leq X$ and $\|\mathbf{w}\| \leq U$
- Can we design an invariant algorithm whose regret scales with C ?



Second-order algorithms in a transductive setting

For any fixed $\mathbf{x}_1, \dots, \mathbf{x}_T$

$$\mathbf{w}'_t = \mathbf{w}_t - \mathbf{A}^{-1} \mathbf{g}_t \quad \mathbf{w}_{t+1} = \operatorname{argmin}_{\mathbf{w} \in \mathcal{K}_C} \|\mathbf{w} - \mathbf{w}'_t\|_A$$

For any $\mathbf{w} \in \mathcal{K}_C$

$$R_T(\mathbf{w}) \leq \|\mathbf{w}\|_A^2 + G_C^2 \sum_{t=1}^T \|\mathbf{x}_t\|_{A^{-1}}^2$$

Now pick the best A for the worst-case $\mathbf{w} \in \mathcal{K}_C$



Optimizing the bound

For any fixed $\mathbf{x}_1, \dots, \mathbf{x}_T$

The solution to $\inf_{A \succ 0} \sup_{\mathbf{w} \in \mathcal{K}'_C} \left(\|\mathbf{w}\|_A^2 + G_C^2 \sum_{t=1}^T \|\mathbf{x}_t\|_{A^{-1}}^2 \right)$

in the larger set $\mathcal{K}'_C = \left\{ \mathbf{w} \in \mathbb{R}^d : \sum_{t=1}^T (\mathbf{w}^\top \mathbf{x}_t)^2 \leq TC^2 \right\}$

is $A = \frac{G_C}{C} \left(\sqrt{\frac{d}{T}} \right) \sum_{t=1}^T \mathbf{x}_t \mathbf{x}_t^\top$

This gives the bound $R_T(\mathbf{w}) \leq CG_C \sqrt{dT}$

Can we get a similar bound in a **non-transductive** setting?



A nearly invariant online algorithm

Online Newton Step with projection

$$\mathbf{w}'_t = \mathbf{w}_t - \mathbf{A}_t^{-1} \mathbf{g}_t$$

$$\mathbf{w}_{t+1} = \operatorname{argmin}_{\mathbf{w} \in \mathcal{K}_{C,t+1}} \|\mathbf{w} - \mathbf{w}'_t\|_{\mathbf{A}_t}$$

$$\text{where } \mathbf{A}_t = \alpha \mathbf{I} + \sum_{s=1}^t \eta_s \mathbf{g}_s \mathbf{g}_s^\top$$

$$\text{and } \mathcal{K}_{C,t+1} = \left\{ \mathbf{w} \in \mathbb{R}^d : |\mathbf{w}^\top \mathbf{x}_{t+1}| \leq C \right\}$$

Regret bound

$$R_T(\mathbf{w}) \leq \alpha \|\mathbf{w}\|^2 + \left(\ln \frac{T}{\alpha} \right) C G_C \sqrt{dT} \quad \text{for all } \mathbf{w} \in \mathcal{K}_{C,T}$$

A truly invariant online algorithm

Online Newton Step with projection and pseudoinverse

$$\mathbf{w}'_t = \mathbf{w}_t - \mathbf{A}_t^+ \mathbf{g}_t$$

$$\mathbf{w}_{t+1} = \operatorname{argmin}_{\mathbf{w} \in \mathcal{K}_{C,t+1}} \|\mathbf{w} - \mathbf{w}'_t\|_{\mathbf{A}_t}$$

$$\text{where } \mathbf{A}_t = \sum_{s=1}^t \eta_s \mathbf{g}_s \mathbf{g}_s^\top$$

$$\text{and } \mathcal{K}_{C,t+1} = \left\{ \mathbf{w} \in \mathbb{R}^d : |\mathbf{w}^\top \mathbf{x}_{t+1}| \leq C \right\}$$

Regret bound

$$R_T(\mathbf{w}) \leq C G_C \left(\sqrt{dT} + (R^2 \ln T) \sqrt{\frac{T}{d}} \right)$$

for all $\mathbf{w} \in \mathcal{K}_{C,T}$ where R is the rank of the data matrix

- Second-order algorithms require updating $A_t^{-1} = \alpha I + G_t G_t^T$ where G_t is $d \times t$
- This is done in time $\Theta(d^2)$
- Want faster algorithms using $B_t = \alpha I + S_t S_t^T$ where S_t is $d \times m$ for $m \ll d$
- Since $B_t^{-1} = \frac{1}{\alpha} \left(I + S_t \underbrace{(\alpha I + S_t S_t^T)^{-1}}_H S_t^T \right)$ where H is $m \times m$
- We can update B_t^{-1} in time $\mathcal{O}(d^2 + md) = \mathcal{O}(md)$



Sketched Online Newton Step with projection

$$\mathbf{w}'_t = \mathbf{w}_t - A_t^{-1} \mathbf{g}_t$$

$$\mathbf{w}_{t+1} = \operatorname{argmin}_{\mathbf{w} \in \mathcal{K}_{C,t+1}} \|\mathbf{w} - \mathbf{w}'_t\|_{A_t}$$

Replace $A_t = \alpha I + \sum_{s=1}^t \mathbf{g}_s \mathbf{g}_s^\top$

with $B_t = \alpha I + S_t S_t^\top$

where S_t is $d \times m$

Regret bound for sketched ONS

$$R_T(\mathbf{w}) \leq \alpha \|\mathbf{w}\|^2 + \sum_{t=1}^T \|\mathbf{g}_t\|_{B_t^{-1}}^2 + \sum_{t=1}^T \|\mathbf{w}_t - \mathbf{w}\|_{B_t - B_{t-1}}^2$$

Gaussian Random Projection

$S_t = S_{t-1} + \mathbf{g}_t \mathbf{Z}_t^\top$ where \mathbf{Z}_t are m i.i.d. draws from $\mathcal{N}(0, m^{-1/2})$

- 1 $(1 - \varepsilon)B_t \preceq A_t \preceq (1 + \varepsilon)B_t$ w.h.p. for $m \geq \frac{R + \ln T}{\varepsilon^2}$
- 2 $\mathbb{E}_t[B_t - B_{t-1}] = \mathbf{g}_t \mathbf{g}_t^\top$

This leads to the same regret bound as before (but only in expectation)

$$R_T(\mathbf{w}) \leq \alpha \|\mathbf{w}\|^2 + \sum_{t=1}^T \underbrace{\|\mathbf{g}_t\|_{B_t^{-1}}^2}_{B_t^{-1} \preceq \frac{1}{1-\varepsilon} A_t^{-1}} + \sum_{t=1}^T \underbrace{\|\mathbf{w}_t - \mathbf{w}\|_{B_t - B_{t-1}}^2}_{= \mathbf{g}_t \mathbf{g}_t^\top}$$

$\mathcal{O}(md)$ space, $\mathcal{O}(md)$ time per round

Item frequency approximation for streams

- 1 Keep m counters for d types of objects, $m \ll d$
- 2 If type of next incoming object has a counter, increase it, otherwise pick a counter at zero and set it to 1
- 3 If there are no counters at zero, then decrease all counters until the smallest hits zero

Matrix approximations

- 1 Keep m orthogonal vectors in \mathbb{R}^d , $m \ll d$
- 2 Add the next incoming vector and make the set orthogonal again via SVD
- 3 Shrink all vectors until the shortest hits zero

$\mathcal{O}(md)$ space, $\mathcal{O}(md)$ time per round

Regret bound

$$R_T(\mathbf{w}) \leq \alpha \|\mathbf{w}\|^2 + \left(\left(\ln \frac{T}{\alpha} \right) CG_C + \frac{\Gamma}{\alpha} \right) m \sqrt{\frac{T}{d}}$$

where $\Gamma_T = \sum_{i=m+1}^d \lambda_i(S_T S_T^T)$ is the eigenvalue tail of the sketch

- When $m = d$ we recover the standard bound
- When data sequence has rank at most m we win



Additional results

- Faster rates for exp-concave losses
- Sketching using Oja's algorithm
- Sparsity-preserving sketching
- Experiments running

