
Supplementary Material for “Nonparametric Online Regression while Learning the Metric”

Ilja Kuzborskij
EPFL
Switzerland
ilja.kuzborskij@gmail.com

Nicolò Cesa-Bianchi
Dipartimento di Informatica
Università degli Studi di Milano
Milano 20135, Italy
nicolo.cesa-bianchi@unimi.it

1 Nonparametric gradient learning

In this section we describe a nonparametric gradient learning algorithm introduced in [2]. Throughout this section, we assume instances \mathbf{x}_t are realizations of i.i.d. random variables \mathbf{X}_t drawn according to some fixed and unknown distribution μ which has a continuous density on its support \mathcal{X} . Labels y_t are generated according to the noise model $y_t = f(\mathbf{x}_t) + \nu(\mathbf{x}_t)$, where $\nu(\mathbf{x})$ is a subgaussian zero-mean random variable for all $\mathbf{x} \in \mathcal{X}$. The algorithm computes a sequence of estimates $\hat{f}_1, \hat{f}_2, \dots$ of the regression function f through kernel regression. Let $\mathcal{X}_n \equiv \{\mathbf{x}_1, \dots, \mathbf{x}_n\} \subset \mathcal{X}$ be the data observed so far and let y_1, \dots, y_n their corresponding labels. Let $K : \mathbb{R}_+ \rightarrow \mathbb{R}_+$ be a nonincreasing kernel, strictly positive on $[0, 1)$, and such that $K(1) = 0$. Then the estimate at time n is defined by

$$\hat{f}_n(\mathbf{x}) = \sum_{t=1}^n y_t \omega_t(\mathbf{x}) \quad \text{where} \quad \omega_t(\mathbf{x}) = \begin{cases} \frac{K(\|\mathbf{x} - \mathbf{x}_t\|/\varepsilon_n)}{\sum_{s=1}^n K(\|\mathbf{x} - \mathbf{x}_s\|/\varepsilon_n)} & \text{if } \mathcal{B}(\mathbf{x}, \varepsilon_n) \cap \mathcal{X}_n \neq \emptyset, \\ 1/n & \text{otherwise} \end{cases}$$

where $\varepsilon_n > 0$ is the kernel scaling parameter. We then approximate the gradient of \hat{f} at any given point through the finite difference method

$$\Delta_i(\mathbf{x}) = \frac{1}{2\tau_n} \left(\hat{f}(\mathbf{x} + \tau_n \mathbf{e}_i) - \hat{f}(\mathbf{x} - \tau_n \mathbf{e}_i) \right) \quad \text{for } i = 1, \dots, d$$

where $\tau_n > 0$ is a parameter. Let further

$$A_i(\mathbf{x}) = \mathbb{I} \left\{ \min_{b \in \{-\tau_n, \tau_n\}} \mu_n(\mathcal{B}(\mathbf{x} + b\mathbf{e}_i, \varepsilon/2)) \geq \frac{2d}{n} (\ln 2n) \right\} \quad \text{for } i = 1, \dots, d$$

where μ_n is the empirical distribution of μ after observing \mathcal{X}_n , and define the gradient estimate

$$\hat{\nabla} f(\mathbf{x}_t) = \left(\Delta_1(\mathbf{x}_t) A_1(\mathbf{x}_t), \dots, \Delta_d(\mathbf{x}_t) A_d(\mathbf{x}_t) \right).$$

The algorithm outputs at time n the gradient outer product estimate

$$\hat{\mathbf{G}}_n = \frac{1}{n} \sum_{t=1}^n \hat{\nabla} f(\mathbf{x}_t) \hat{\nabla} f(\mathbf{x}_t)^\top$$

Let $\mathbf{G} = \mathbb{E} [\nabla f(\mathbf{X}) \nabla f(\mathbf{X})^\top]$ be the expected gradient outer product, where \mathbf{X} has law μ . The next lemma states that, under Assumption 1, $\hat{\mathbf{G}}_n$ is a consistent estimate of \mathbf{G} .

Lemma (Consistency of the Expected Gradient Outerproduct Estimator [2, Theorem 1]). *If Assumption 1 holds, then there exists a nonnegative and nonincreasing sequence $\{\gamma_n\}_{n \geq 1}$ such that for all n , the estimated gradient outerproduct (7) computed with parameters $\varepsilon_n > 0$, and $0 < \tau_n < \tau_0$ satisfies $\|\hat{\mathbf{G}}_n - \mathbf{G}\|_2 \leq \gamma_n$ with high probability with respect to the random draw of $\mathbf{X}_1, \dots, \mathbf{X}_n$. Moreover, if $\tau_n = \Theta(\varepsilon_n^{1/4})$, $\varepsilon_n = \Omega\left((\ln n)^{\frac{2}{d}} n^{-\frac{1}{d}}\right)$, and $\varepsilon_n = \mathcal{O}\left(n^{-\frac{1}{2(d+1)}}\right)$ then $\gamma_n \rightarrow 0$ as $n \rightarrow \infty$.*

The actual rate of convergence depends, in a complicated way, on parameters related to the distribution μ and the regression function f . In our application of Lemma 4 we assume $\gamma_n \leq n^{-\alpha}$ for all n large enough and for some $\alpha > 0$. Note also that the convergence of \hat{G}_n to G holds in probability with respect to the random draw of $\mathbf{X}_1, \dots, \mathbf{X}_n$. Hence there is a confidence parameter δ which is not shown here. However, the dependence of the convergence rate on $\frac{1}{\delta}$ is only polylogarithmic and therefore not problematic for our applications.

2 Proofs from Section 3

Lemma (Volumetric packing bound). *Consider a pair of norms $\|\cdot\|, \|\cdot\|'$ and let $B, B' \subset \mathbb{R}^d$ be the corresponding unit balls. Then*

$$\mathcal{M}(B, \varepsilon, \|\cdot\|') \leq \frac{\text{vol}\left(B + \frac{\varepsilon}{2}B'\right)}{\text{vol}\left(\frac{\varepsilon}{2}B'\right)}.$$

Proof. Let $\{\mathbf{x}_1, \dots, \mathbf{x}_M\}$ be a maximal ε -packing of B according to $\|\cdot\|'$. Since we have a packing, the $\|\cdot\|'$ -balls of radius $\varepsilon/2$ and centers $\mathbf{x}_1, \dots, \mathbf{x}_M$ are disjoint, and their union is contained in $B + \frac{\varepsilon}{2}B'$. Thus,

$$M \text{vol}\left(\frac{\varepsilon}{2}B'\right) \leq \text{vol}\left(B + \frac{\varepsilon}{2}B'\right)$$

which concludes the proof. \square

Lemma (Ellipsoid packing bound). *If B is the unit Euclidean ball then*

$$\mathcal{M}(B, \varepsilon, \|\cdot\|_M) \leq \left(\frac{8\sqrt{2}}{\varepsilon}\right)^s \prod_{i=1}^s \sqrt{\lambda_i} \quad \text{where} \quad s = \max \left\{ i : \sqrt{\lambda_i} \geq \varepsilon, i = 1, \dots, d \right\}.$$

Proof. The change of variable $\mathbf{x}' = \mathbf{M}^{1/2}\mathbf{x}$ implies $\|\mathbf{x}\|_2 = \|\mathbf{x}'\|_{M^{-1}}$ and $\|\mathbf{x}\|_M = \|\mathbf{x}'\|_2$. Therefore $\mathcal{M}(B, \varepsilon, \|\cdot\|_M) = \mathcal{M}(E, \varepsilon, \|\cdot\|_2)$ where $E \equiv \{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x}\|_{M^{-1}} \leq 1\}$ is the unit ball in the norm $\|\cdot\|_{M^{-1}}$. Next, we write the coordinates (x_1, \dots, x_d) of any point $\mathbf{x} \in \mathbb{R}^d$ using the orthonormal basis $\mathbf{u}_1, \dots, \mathbf{u}_d$. Consider the truncated ellipsoid $\tilde{E} \equiv \{\mathbf{x} \in E : x_i = 0, i = s+1, \dots, d\}$. By adapting an argument from [3], we prove that any ε -cover of \tilde{E} according to $\|\cdot\|_2$ is also a $(\varepsilon\sqrt{2})$ -cover of E according to the same norm. Indeed, let $\tilde{S} \subset \tilde{E}$ be a ε -cover of \tilde{E} . Fix any $\mathbf{x} \in E$ and let

$$\begin{aligned} \min_{\tilde{\mathbf{x}} \in \tilde{S}} \|\mathbf{x} - \tilde{\mathbf{x}}\|_2^2 &= \min_{\tilde{\mathbf{x}} \in \tilde{S}} \sum_{j=1}^s (x_j - \tilde{x}_j)^2 + \sum_{j=s+1}^d x_j^2 \\ &\leq \varepsilon^2 + \sum_{j=s+1}^d x_j^2 && \text{(since } \tilde{S} \text{ is a } \varepsilon\text{-covering of } \tilde{E}) \\ &\leq \varepsilon^2 + \lambda_{s+1} \sum_{j=s+1}^d \frac{x_j^2}{\lambda_j} && \text{(since } \lambda_{s+1}/\lambda_j \geq 1 \text{ for } j = s+1, \dots, d) \\ &\leq 2\varepsilon^2 \end{aligned}$$

where the last inequality holds since $\lambda_{s+1} \leq \varepsilon^2$ and since $\|\mathbf{x}\|_{M^{-1}}^2 = \sum_{i=1}^d x_i^2/\lambda_i \leq 1$ for any $\mathbf{x} \in E$, where $x_i = \mathbf{u}_i^\top \mathbf{x}$ for all $i = 1, \dots, d$. Let $B' \subset \mathbb{R}^d$ be the unit Euclidean ball, and let $\tilde{B}' \equiv \{\mathbf{x} \in B' : x_i = 0, i = s+1, \dots, d\}$ be its truncated version. Since $\lambda_i \geq \varepsilon^2$ for $i = 1, \dots, s$ we have that for all $\mathbf{x} \in \varepsilon\tilde{B}'$, $x_1^2 + \dots + x_s^2 \leq \varepsilon^2$ and so

$$\|\mathbf{x}\|_{M^{-1}}^2 = \sum_{i=1}^s \frac{x_i^2}{\lambda_i} \leq \sum_{i=1}^s \frac{\varepsilon^2}{\lambda_i} \leq 1.$$

Therefore $\varepsilon \tilde{B}' \subseteq \tilde{E}$ which implies $\text{vol}(\tilde{E} + \frac{\varepsilon}{2} \tilde{B}') \leq \text{vol}(2\tilde{E})$.

$$\begin{aligned}
\mathcal{M}(E, 2\varepsilon\sqrt{2}, \|\cdot\|_2) &\leq \mathcal{N}(E, \varepsilon\sqrt{2}, \|\cdot\|_2) \\
&\leq \mathcal{N}(\tilde{E}, \varepsilon, \|\cdot\|_2) \\
&\leq \mathcal{M}(\tilde{E}, \varepsilon, \|\cdot\|_2) \\
&\leq \frac{\text{vol}(\tilde{E} + \frac{\varepsilon}{2} \tilde{B}')}{\text{vol}(\frac{\varepsilon}{2} \tilde{B}')} \quad (\text{by Lemma 1}) \\
&\leq \frac{\text{vol}(2\tilde{E})}{\text{vol}(\frac{\varepsilon}{2} \tilde{B}')} = \left(\frac{4}{\varepsilon}\right)^s \frac{\text{vol}(\tilde{E})}{\text{vol}(\tilde{B}')}
\end{aligned}$$

Now, using the standard formula for the volume of an ellipsoid,

$$\text{vol}(\tilde{E}) = \text{vol}(\tilde{B}') \prod_{i=1}^s \sqrt{\lambda_i}.$$

This concludes the proof. \square

The following lemma states that whenever f has bounded partial derivatives with respect to the eigenbase of M , then f is Lipschitz with respect to $\|\cdot\|_M$.

Lemma (Bounded derivatives imply Lipschitzness in M -metric). *Let $f : \mathcal{X} \rightarrow \mathbb{R}$ be everywhere differentiable. Then for any $\mathbf{x}, \mathbf{x}' \in \mathcal{X}$,*

$$|f(\mathbf{x}) - f(\mathbf{x}')| \leq \|\mathbf{x} - \mathbf{x}'\|_M \sqrt{\sum_{i=1}^d \frac{\|\nabla_{\mathbf{u}_i} f\|_\infty^2}{\lambda_i}}.$$

Proof. By the mean value theorem, there exists a \mathbf{z} on the segment joining \mathbf{x} and \mathbf{y} such that $f(\mathbf{x}) - f(\mathbf{y}) = \nabla f(\mathbf{z})^\top (\mathbf{x} - \mathbf{y})$. Hence

$$\begin{aligned}
f(\mathbf{x}) - f(\mathbf{y}) &= \nabla f(\mathbf{z})^\top (\mathbf{x} - \mathbf{y}) \\
&= \sum_{i=1}^d \nabla f(\mathbf{z})^\top \mathbf{u}_i \mathbf{u}_i^\top (\mathbf{x} - \mathbf{y}) \\
&\leq \sum_{i=1}^d \left(\sup_{\mathbf{z}' \in \mathcal{X}} \nabla f(\mathbf{z}')^\top \mathbf{u}_i \right) \mathbf{u}_i^\top (\mathbf{x} - \mathbf{y}) \\
&= \sum_{i=1}^d \frac{\|\nabla_{\mathbf{u}_i} f\|_\infty}{\sqrt{\lambda_i}} \left(\sqrt{\lambda_i} \mathbf{u}_i^\top (\mathbf{x} - \mathbf{y}) \right) \\
&\leq \sqrt{\sum_{i=1}^d \frac{\|\nabla_{\mathbf{u}_i} f\|_\infty^2}{\lambda_i}} \sqrt{\sum_{i=1}^d \lambda_i (\mathbf{u}_i^\top (\mathbf{x} - \mathbf{y}))^2} \quad (\text{by the Cauchy-Schwarz inequality}) \\
&= \|\mathbf{x} - \mathbf{y}\|_M \sqrt{\sum_{i=1}^d \frac{\|\nabla_{\mathbf{u}_i} f\|_\infty^2}{\lambda_i}}.
\end{aligned}$$

By symmetry, we can upper bound $f(\mathbf{y}) - f(\mathbf{x})$ with the same quantity. \square

Now we are ready to prove the regret bound.

Theorem (Regret with Fixed Metric). *Suppose Algorithm 1 is run with a positive definite matrix M with eigenbasis $\mathbf{u}_1, \dots, \mathbf{u}_d$ and eigenvalues $1 = \lambda_1 \geq \dots \geq \lambda_d > 0$. Then, for any differentiable $f : \mathcal{X} \rightarrow \mathbb{R}$ we have that*

$$R_T(f) \stackrel{\tilde{O}}{=} \left(\sqrt{\det_\kappa(M)} + \sqrt{\sum_{i=1}^d \frac{\|\nabla_{\mathbf{u}_i} f\|_\infty^2}{\lambda_i}} \right) T^{\frac{\rho_T}{1+\rho_T}}$$

where $\kappa = \kappa(\rho_T, T) \leq \rho_T \leq d$.

Proof. Let S_t be the value of the variable S at the end of time t . Hence $S_0 = \emptyset$. The functions $\pi_t : \mathcal{X} \rightarrow \{1, \dots, t\}$ for $t = 1, 2, \dots$ map each data point \mathbf{x} to its closest (in norm $\|\cdot\|_M$) center in S_{t-1} ,

$$\pi_t(\mathbf{x}) = \begin{cases} \arg \min_{s \in S_{t-1}} \|\mathbf{x} - \mathbf{x}_s\|_M & \text{if } S_{t-1} \neq \emptyset \\ t & \text{otherwise.} \end{cases}$$

The set T_s contain all data points \mathbf{x}_t that at time t belonged to the ball with center \mathbf{x}_s and radius ε_t ,

$$T_s \equiv \{t : \|\mathbf{x}_t - \mathbf{x}_s\|_M \leq \varepsilon_t, t = s, \dots, T\}.$$

Finally, y_s^* is the best fixed prediction for all examples (\mathbf{x}_t, y_t) such that $t \in T_s$,

$$y_s^* = \arg \min_{y \in \mathcal{Y}} \sum_{t \in T_s} \ell_t(y) = \frac{1}{|T_s|} \sum_{t \in T_s} y_t. \quad (1)$$

We proceed by decomposing the regret into a local (estimation) and a global (approximation) term,

$$R_T(f) = \sum_{t=1}^T (\ell_t(\hat{y}_t) - \ell_t(f(\mathbf{x}_t))) = \sum_{t=1}^T (\ell_t(\hat{y}_t) - \ell_t(y_{\pi_t(\mathbf{x}_t)}^*)) + \sum_{t=1}^T (\ell_t(y_{\pi_t(\mathbf{x}_t)}^*) - \ell_t(f(\mathbf{x}_t))).$$

The estimation term is bounded as

$$\sum_{t=1}^T (\ell_t(\hat{y}_t) - \ell_t(y_{\pi_t(\mathbf{x}_t)}^*)) = \sum_{s \in S_T} \sum_{t \in T_s} (\ell_t(\hat{y}_t) - \ell_t(y_s^*)) \leq 8 \sum_{s \in S_T} \ln(e|N_s|) \leq 8 \ln(eT) |S_T|.$$

The first inequality is a known bound on the regret under square loss [1, page 43]. We upper bound the size of the final packing S_T using Lemma 2,

$$|S_T| \leq \mathcal{M}(B, \varepsilon_T, \|\cdot\|_M) \leq \left(\frac{8\sqrt{2}}{\varepsilon_T} \right)^\kappa \prod_{i=1}^\kappa \sqrt{\lambda_i} \leq (8\sqrt{2})^\kappa \sqrt{\det_\kappa(\mathbf{M})} T^{\frac{\kappa}{1+\rho_T}}$$

where $\kappa = \kappa(\rho_T, T)$. Therefore, since $\rho_T \geq \kappa(\rho_T, T)$,

$$\sum_{t=1}^T (\ell_t(\hat{y}_t) - \ell_t(y_{\pi_t(\mathbf{x}_t)}^*)) \leq 8 \ln(eT) (8\sqrt{2})^{\rho_T} \sqrt{\det_\kappa(\mathbf{M})} T^{\frac{\rho_T}{1+\rho_T}}. \quad (2)$$

Next, we bound the approximation term. Using (1) we have

$$\sum_{t=1}^T (\ell_t(y_{\pi_t(\mathbf{x}_t)}^*) - \ell_t(f(\mathbf{x}_t))) \leq \sum_{t=1}^T (\ell_t(f(\mathbf{x}_{\pi_t(\mathbf{x}_t)})) - \ell_t(f(\mathbf{x}_t))).$$

Note that ℓ_t is 2-Lipschitz because $y_t, \hat{y}_t \in [0, 1]$. Hence, using Lemma 3,

$$\begin{aligned} \ell_t(f(\mathbf{x}_{\pi_t(\mathbf{x}_t)})) - \ell_t(f(\mathbf{x}_t)) &\leq 2 \|f(\mathbf{x}_{\pi_t(\mathbf{x}_t)}) - f(\mathbf{x}_t)\| \\ &\leq 2 \|\mathbf{x}_t - \mathbf{x}_{\pi_t(\mathbf{x}_t)}\|_M \sqrt{\sum_{i=1}^d \frac{\|\nabla_{\mathbf{u}_i} f\|_\infty^2}{\lambda_i}} \\ &\leq 2\varepsilon_t \sqrt{\sum_{i=1}^d \frac{\|\nabla_{\mathbf{u}_i} f\|_\infty^2}{\lambda_i}}. \end{aligned}$$

Recalling $\varepsilon_t = t^{-\frac{1}{1+\rho_t}}$ where $\rho_t \leq \rho_{t+1}$, we write

$$\sum_{t=1}^T t^{-\frac{1}{1+\rho_t}} \leq \sum_{t=1}^T t^{-\frac{1}{1+\rho_T}} \leq \int_0^T \tau^{-\frac{1}{1+\rho_T}} d\tau = \left(1 + \frac{1}{\rho_T}\right) T^{\frac{\rho_T}{1+\rho_T}} \leq 2T^{\frac{\rho_T}{1+\rho_T}}.$$

Thus we may write

$$\sum_{t=1}^T (\ell_t(y_{\pi_t(\mathbf{x}_t)}^*) - \ell_t(f(\mathbf{x}_t))) \leq 4 \left(\sqrt{\sum_{i=1}^d \frac{\|\nabla_{\mathbf{u}_i} f\|_\infty^2}{\lambda_i}} \right) T^{\frac{\rho_T}{1+\rho_T}}.$$

The proof is concluded by combining the above with (2). \square

References

- [1] N. Cesa-Bianchi and G. Lugosi. *Prediction, Learning, and Games*. Cambridge University Press, 2006.
- [2] S. Trivedi, J. Wang, S. Kpotufe, and G. Shakhnarovich. A consistent Estimator of the Expected Gradient Outerproduct. In *Conference on Uncertainty in Artificial Intelligence (UAI)*, 2014.
- [3] M. Wainwright. *High-Dimensional Statistics: A Non-Asymptotic Viewpoint (In Preparation)*. 2017.