

# Data sparse nonparametric regression with $\epsilon$ -insensitive losses: supplementary material

**Maxime Sangnier**

MAXIME.SANGNIER@UPMC.FR

*Sorbonne Universités, UPMC Univ Paris 06, CNRS, Paris, France*

**Olivier Fercoq**

OLIVIER.FERCOQ@TELECOM-PARISTECH.FR

*Université Paris-Saclay, Télécom ParisTech, LTCI, Paris, France*

**Florence d’Alché-Buc**

FLORENCE.DALCHE@TELECOM-PARISTECH.FR

*Université Paris-Saclay, Télécom ParisTech, LTCI, Paris, France*

**Editors:** Yung-Kyun Noh and Min-Ling Zhang

## Appendix A. Remarks

**Remark 1** *It is easy to see that, for a unidimensional loss  $\ell$ ,  $\ell_\epsilon$  is obtained in the following manner:*

$$\begin{aligned} \forall \xi \in \mathbb{R}: \quad \ell_\epsilon(\xi) &= \ell(\xi - \min(|\xi|, \epsilon) \operatorname{sign}(\xi)) \\ &= \begin{cases} 0 & \text{if } |\xi| \leq \epsilon \\ \ell\left(\xi\left(1 - \frac{\epsilon}{|\xi|}\right)\right) & \text{otherwise.} \end{cases} \end{aligned}$$

*Consequently, when a multivariate loss  $\ell$  is separable, that is  $\ell(\boldsymbol{\xi}) = \sum_{j=1}^p \ell^{(j)}(\xi_j)$  (for some unidimensional losses  $\ell^{(j)}$ ), it is tempting to consider each component separately and to define  $\ell_\epsilon = \sum_{j=1}^p \ell_\epsilon^{(j)}$ . Basically, this boils down to replacing  $\|\cdot\|_{\ell_2}$  by  $\|\cdot\|_\infty$  in the general  $\epsilon$ -loss introduced in this paper.*

*However, this is not a good idea since this definition would result in adding an  $\ell_1$ -norm  $\sum_{i=1}^n \|\boldsymbol{\alpha}_i\|_{\ell_1}$  instead of an  $\ell_1/\ell_2$ -norm in the dual. As a consequence, we would obtain sparse vectors  $\boldsymbol{\alpha}_i$ , which is not the data sparsity we pursue since  $\boldsymbol{\alpha}_i$  could have null components but could be different from 0, forcing us to keep the points  $\mathbf{x}_i$  for prediction.*

**Remark 2** *In the body of the text, omitting the intercept  $\mathbf{b}$  in Problem (P2) comes down to removing the linear constraint in Problem (P3). This practice is common for [support vector regression \(SVR\)](#) with a Gaussian kernel, but is excluded for [quantile regression \(QR\)](#) (Takeuchi et al., 2006; Sangnier et al., 2016).*

**Example 1** *Examples of scalar and matrix kernels are:*

$$k(\mathbf{x}, \mathbf{x}') = (1 + \langle \mathbf{x}, \mathbf{x}' \rangle_{\ell_2})^d \quad (\text{polynomial}),$$

*where  $d > 0$  is the degree (Mohri et al., 2012), and*

$$K(\mathbf{x}, \mathbf{x}') = \left[ \left( 1 + \langle T_i(\mathbf{x}), T_j(\mathbf{x}') \rangle_{\ell_2} \right)^d \right]_{1 \leq i, j \leq p} \quad (\text{transformable}),$$

where  $T_i: \mathbb{R}^p \rightarrow \mathbb{R}^p$  are any transformations (Alvarez et al., 2012).

**Remark 3** As it is standard for coordinate descent methods, our implementation uses efficient updates for the computation of both  $\sum_{j=1}^n K(\mathbf{x}_i, \mathbf{x}_j) \boldsymbol{\alpha}_j$  and  $\bar{\boldsymbol{\theta}}^l$ . In addition, convergence of Algorithm 1 can be assessed by duality gap (objective of (P2) minus objective of (P3) in the body of the text). Yet, even though we do not have closed-form expressions for the primal loss  $\ell_\epsilon$ , the duality gap can be over-estimated by upper-bounding  $\ell_\epsilon$  in the following manner:

$$\forall \boldsymbol{\xi} \in \mathbb{R}^p: \quad \ell_\epsilon(\boldsymbol{\xi}) \leq \ell \left( \boldsymbol{\xi} \left( 1 - \frac{\min(\epsilon, \|\boldsymbol{\xi}\|_{\ell_2})}{\|\boldsymbol{\xi}\|_{\ell_2}} \right) \right).$$

This is true since  $\left\| \frac{\min(\epsilon, \|\boldsymbol{\xi}\|_{\ell_2})}{\|\boldsymbol{\xi}\|_{\ell_2}} \boldsymbol{\xi} \right\|_{\ell_2} \leq \epsilon$ .

**Remark 4** Contrarily to QR, expectile regression involves a differentiable mapping  $\ell^*$ . Consequently, it can be easily incorporated to the quadratic contribution of (P7) (see body of the text). Nevertheless, it can also be considered jointly with  $\|\cdot\|_{\ell_2}$ , in the same manner as for QR. In this case, the differentiable part remains the same for expectile and quantile regression, only the non-differentiable part changes. The proximal operator needed is given in the following proposition.

**Proposition 5** Let  $\psi: \mathbf{y} \in \mathbb{R}^p \mapsto \frac{1}{2} \sum_{i=1}^p |\tau_j - \mathbf{I}_{y_j < 0}|^{-1} y_j^2$ . Then

$$\forall \mathbf{y} \in \mathbb{R}^p, \forall j \in [p] \\ \left[ \text{prox}_{\lambda(\|\cdot\|_{\ell_2} + \psi)}(\mathbf{y}) \right]_j = \left( 1 + \frac{\lambda}{\mu} + \lambda |\tau_j - \mathbf{I}_{y_j < 0}|^{-1} \right)^{-1} y_j,$$

if  $\|\mathbf{y}\|_{\ell_2} > \lambda$ , where  $\mu > 0$  is solution to:

$$\sum_{j=1}^p \frac{y_j^2}{\left( \mu \left( 1 + \lambda |\tau_j - \mathbf{I}_{y_j < 0}|^{-1} \right) + \lambda \right)^2} = 1, \quad (1)$$

(such a solution exists) and  $\text{prox}_{\lambda(\|\cdot\|_{\ell_2} + \psi)}(\mathbf{y}) = 0$  if  $\|\mathbf{y}\|_{\ell_2} \leq \lambda$ .

Similarly to Equation 1 in the body of the text, the scaling factor  $\mu$  in Proposition 5 can be easily obtained by a bisection of a Newton-Raphson method.

## Appendix B. Technical details

### B.1. Convexity and redefinition of $\ell_\epsilon$

For the sake of simplicity, let us first define:

$$\forall \boldsymbol{\xi} \in \mathbb{R}^p: \quad \tilde{\ell}_\epsilon(\boldsymbol{\xi}) = \inf_{\mathbf{u} \in \mathbb{R}^p: \|\mathbf{u}\|_{\ell_2} \leq \epsilon} \ell(\boldsymbol{\xi} - \mathbf{u}). \quad (2)$$

Since  $\ell$  is convex,  $(\xi, \mathbf{u}) \mapsto \ell(\xi - \mathbf{u}) + \chi_{\|\mathbf{u}\|_{\ell_2} \leq \epsilon}$  is jointly convex with respect to  $\xi$  and  $\mathbf{u}$ . Therefore,  $\tilde{\ell}_\epsilon$  is convex as the coordinate infimum of a jointly convex function (Boyd and Vandenberghe, 2004).

Let us now show that  $\tilde{\ell}_\epsilon = \ell_\epsilon$ . First, since for any  $\xi$ , Slater's constraint qualification are satisfied for (2), strong duality holds, that is,

$$\forall \xi \in \mathbb{R}^p, \exists \lambda \geq 0 : \tilde{\ell}_\epsilon(\xi) = \inf_{\mathbf{u} \in \mathbb{R}^p} \ell(\xi - \mathbf{u}) + \lambda \|\mathbf{u}\|_{\ell_2} - \lambda\epsilon,$$

and thanks to the lower semi-continuity of the objective, the infimum is attained at, let us say,  $\hat{\mathbf{u}}$ . Then, when  $\|\xi\|_{\ell_2} \leq \epsilon$ , we can chose  $\hat{\mathbf{u}} = \xi$  and we get  $\tilde{\ell}_\epsilon(\xi) = 0$ , which is the infimum of  $\ell$ . On the other hand, when  $\|\xi\|_{\ell_2} > \epsilon$ , let us consider the Karush-Kuhn-Tucker (KKT) conditions. By complementary slackness, either  $\lambda = 0$  and  $\|\hat{\mathbf{u}}\|_{\ell_2} \leq \epsilon$ , or  $\|\hat{\mathbf{u}}\|_{\ell_2} = \epsilon$ . In the first situation ( $\lambda = 0$  and  $\|\hat{\mathbf{u}}\|_{\ell_2} \leq \epsilon$ ),  $\tilde{\ell}_\epsilon(\xi) = \inf_{\mathbf{u} \in \mathbb{R}^p} \ell(\xi - \mathbf{u}) = \ell(\xi - \hat{\mathbf{u}}) = \ell(0) = 0$  and  $\hat{\mathbf{u}} = \xi$  (by uniqueness of the minimizer of  $\ell$ ). Thus,  $\|\xi\|_{\ell_2} \leq \epsilon$ , which is contradictory. Consequently, we have necessarily,  $\|\hat{\mathbf{u}}\|_{\ell_2} = \epsilon$ . To summarize:

$$\forall \xi \in \mathbb{R}^p : \tilde{\ell}_\epsilon(\xi) = \begin{cases} 0 & \text{if } \|\xi\|_{\ell_2} \leq \epsilon \\ \inf_{\mathbf{u} \in \mathbb{R}^p : \|\mathbf{u}\|_{\ell_2} = \epsilon} \ell(\xi - \mathbf{u}) & \text{otherwise,} \end{cases}$$

which is exactly the definition of  $\ell_\epsilon$ .

## B.2. Dual and representer theorem

Since  $\ell_\epsilon$  is convex and can be replaced by (2), Problem (P2) from the body of the text can be reformulated in (Lagrange multipliers are indicated on the right):

$$\begin{aligned} & \underset{\substack{h \in \mathcal{H}, b \in \mathbb{R}^p, \\ \forall i \in [n], \xi_i \in \mathbb{R}^p, \mathbf{r}_i \in \mathbb{R}^p}}{\text{minimize}} \quad \frac{\lambda}{2} \|h\|_{\mathcal{H}}^2 + \frac{1}{n} \sum_{i=1}^n \ell(\xi_i) \\ & \text{s. t.} \quad \begin{cases} \forall i \in [n], \\ \frac{\mathbf{y}_i - (h(\mathbf{x}_i) + b)}{n} = \frac{\mathbf{r}_i + \xi_i}{n} & : \alpha_i \in \mathbb{R}^p \\ \frac{\|\mathbf{r}_i\|_{\ell_2}^2}{2\epsilon n} \leq \frac{\epsilon}{2n} & : \mu_i \in \mathbb{R}_+ \end{cases} \end{aligned} \tag{P1}$$

Let us compute a dual to Problem (P1). The Lagrangian reads:

$$\begin{aligned}
 \mathfrak{L}(h, \mathbf{b}, (\boldsymbol{\xi}_i)_{1 \leq i \leq n}, (\mathbf{r}_i)_{1 \leq i \leq n}, (\boldsymbol{\alpha}_i)_{1 \leq i \leq n}, (\mu_i)_{1 \leq i \leq n}) \\
 = \frac{\lambda}{2} \|h\|_{\mathcal{H}}^2 + \frac{1}{n} \sum_{i=1}^n \ell(\boldsymbol{\xi}_i) + \frac{1}{n} \sum_{i=1}^n \langle \boldsymbol{\alpha}_i, \mathbf{y}_i - (h(\mathbf{x}_i) + \mathbf{b}) - \mathbf{r}_i - \boldsymbol{\xi}_i \rangle_{\ell_2} \\
 + \frac{1}{2\epsilon n} \sum_{i=1}^n \mu_i \|\mathbf{r}_i\|_{\ell_2}^2 - \frac{\epsilon}{2n} \sum_{i=1}^n \mu_i \\
 = \frac{1}{n} \sum_{i=1}^n (\ell(\boldsymbol{\xi}_i) - \langle \boldsymbol{\alpha}_i, \boldsymbol{\xi}_i \rangle_{\ell_2}) + \frac{\lambda}{2} \|h\|_{\mathcal{H}}^2 - \left\langle \frac{1}{n} \sum_{i=1}^n E_{\mathbf{x}_i}^* \boldsymbol{\alpha}_i, h \right\rangle_{\mathcal{H}} \\
 - \left\langle \frac{1}{n} \sum_{i=1}^n \boldsymbol{\alpha}_i, \mathbf{b} \right\rangle_{\ell_2} + \frac{1}{n} \sum_{i=1}^n \left( \frac{\mu_i}{2\epsilon} \|\mathbf{r}_i\|_{\ell_2}^2 - \langle \boldsymbol{\alpha}_i, \mathbf{r}_i \rangle_{\ell_2} \right) + \frac{1}{n} \sum_{i=1}^n \langle \boldsymbol{\alpha}_i, \mathbf{y}_i \rangle_{\ell_2} \\
 - \frac{\epsilon}{2n} \sum_{i=1}^n \mu_i.
 \end{aligned}$$

The objective function of the dual problem to (P1) is obtained by minimizing the Lagrangian with respect to the primal variables  $h$ ,  $\mathbf{b}$ ,  $(\boldsymbol{\xi}_i)_{1 \leq i \leq n}$  and  $(\mathbf{r}_i)_{1 \leq i \leq n}$ . For this purpose, let us remark that minimizing on  $\boldsymbol{\xi}_i$  boils down to introducing the Fenchel-Legendre transform of  $\ell$ :  $\ell^* : \boldsymbol{\alpha} \in \mathbb{R}^p \mapsto \sup_{\boldsymbol{\xi} \in \mathbb{R}^p} \langle \boldsymbol{\alpha}, \boldsymbol{\xi} \rangle_{\ell_2} - \ell(\boldsymbol{\xi})$ . Thus, it remains to compute:

$$\begin{aligned}
 \mathfrak{L}_D((\boldsymbol{\alpha}_i)_{1 \leq i \leq n}, (\mu_i)_{1 \leq i \leq n}) \\
 = \inf_{\substack{h \in \mathcal{H}, \mathbf{b} \in \mathbb{R}^p, \\ \forall i \in [n], \mathbf{r}_i \in \mathbb{R}^p}} \left\{ -\frac{1}{n} \sum_{i=1}^n \ell^*(\boldsymbol{\alpha}_i) + \frac{\lambda}{2} \|h\|_{\mathcal{H}}^2 - \left\langle \frac{1}{n} \sum_{i=1}^n E_{\mathbf{x}_i}^* \boldsymbol{\alpha}_i, h \right\rangle_{\mathcal{H}} \right. \\
 \left. - \left\langle \frac{1}{n} \sum_{i=1}^n \boldsymbol{\alpha}_i, \mathbf{b} \right\rangle_{\ell_2} + \frac{1}{n} \sum_{i=1}^n \left( \frac{\mu_i}{2\epsilon} \|\mathbf{r}_i\|_{\ell_2}^2 - \langle \boldsymbol{\alpha}_i, \mathbf{r}_i \rangle_{\ell_2} \right) + \frac{1}{n} \sum_{i=1}^n \langle \boldsymbol{\alpha}_i, \mathbf{y}_i \rangle_{\ell_2} \right. \\
 \left. - \frac{\epsilon}{2n} \sum_{i=1}^n \mu_i \right\}.
 \end{aligned}$$

Since  $\mathcal{H}$  is unbounded in all directions, the minimum of  $\mathfrak{L}$  with respect to  $h$ ,  $\mathbf{b}$  and  $(\mathbf{r}_i)_{i=1}^n$  is obtained by setting the gradients to 0, which leads to  $h = \frac{1}{\lambda n} \sum_{i=1}^n E_{\mathbf{x}_i}^* \boldsymbol{\alpha}_i$ ,  $\sum_{i=1}^n \boldsymbol{\alpha}_i = 0$  and  $\mathbf{r}_i = \frac{\epsilon}{\mu_i} \boldsymbol{\alpha}_i$ ,  $\forall i \in [n]$ . Thus, the dual objective reads:

$$\begin{aligned}
 \mathfrak{L}_D((\boldsymbol{\alpha}_i)_{1 \leq i \leq n}, (\mu_i)_{1 \leq i \leq n}) \\
 = -\frac{1}{n} \sum_{i=1}^n \ell^*(\boldsymbol{\alpha}_i) - \frac{1}{2\lambda n^2} \sum_{i,j=1}^n \left\langle \boldsymbol{\alpha}_i, E_{\mathbf{x}_i} E_{\mathbf{x}_j}^* \boldsymbol{\alpha}_j \right\rangle_{\ell_2} + \frac{1}{n} \sum_{i=1}^n \langle \boldsymbol{\alpha}_i, \mathbf{y}_i \rangle_{\ell_2} \\
 - \frac{\epsilon}{2n} \sum_{i=1}^n \left( \frac{1}{\mu_i} \|\boldsymbol{\alpha}_i\|_{\ell_2}^2 + \mu_i \right).
 \end{aligned}$$

Then, the dual optimization problem consists in maximizing  $\mathfrak{L}_D$  subject to the constraints  $\sum_{i=1}^n \boldsymbol{\alpha}_i = 0$  and  $\mu_i \geq 0$ ,  $\forall i \in [n]$ . Remarking that  $\inf_{\forall i \in [n], \mu_i \in \mathbb{R}_+} \frac{1}{2} \sum_{i=1}^n \left( \frac{1}{\mu_i} \|\boldsymbol{\alpha}_i\|_{\ell_2}^2 + \mu_i \right) =$

$\sum_{i=1}^n \|\boldsymbol{\alpha}_i\|_{\ell_2}$  (Bach et al., 2012), a dual to Problem (P1) is:

$$\begin{aligned} & \underset{\forall i \in [n], \boldsymbol{\alpha}_i \in \mathbb{R}^p}{\text{minimize}} \quad \frac{1}{n} \sum_{i=1}^n \ell^*(\boldsymbol{\alpha}_i) + \frac{1}{2\lambda n^2} \sum_{i,j=1}^n \left\langle \boldsymbol{\alpha}_i, E_{\mathbf{x}_i} E_{\mathbf{x}_j}^* \boldsymbol{\alpha}_j \right\rangle_{\ell_2} \\ & \quad - \frac{1}{n} \sum_{i=1}^n \langle \boldsymbol{\alpha}_i, \mathbf{y}_i \rangle_{\ell_2} + \frac{\epsilon}{n} \sum_{i=1}^n \|\boldsymbol{\alpha}_i\|_{\ell_2} \\ & \text{s. t.} \quad \sum_{i=1}^n \boldsymbol{\alpha}_i = 0. \end{aligned} \quad (\text{P2})$$

### B.3. Generalization

Let  $P: f \in \mathcal{F} \mapsto \mathbb{E}[\ell(Y - f(X))]$  and  $P_n: f \in \mathcal{F} \mapsto \frac{1}{n} \sum_{i=1}^n \ell_\epsilon(Y - f(X))$ , as well as respective twins  $P_\epsilon$  and  $P_{n,\epsilon}$  obtained by substituting  $\ell_\epsilon$  to  $\ell$ . Let us decompose  $P\hat{f}_\epsilon - Pf^\dagger$ :

$$P\hat{f}_\epsilon - Pf^\dagger = (P\hat{f}_\epsilon - P_n\hat{f}_\epsilon) + (P_n\hat{f}_\epsilon - P_nf^\dagger) + (P_nf^\dagger - Pf^\dagger).$$

First, by concentration inequalities (Bartlett and Mendelson, 2002; Maurer, 2016; Sangnier et al., 2016), we have, with probability greater than  $1 - \delta$ :

$$P\hat{f}_\epsilon - P_n\hat{f}_\epsilon \leq \sup_{f \in \mathcal{F}} (Pf - P_nf) \leq 2\sqrt{2}L\mathcal{R}_n(\mathcal{F}) + LM\sqrt{\frac{\log(1/\delta)}{2n}}.$$

Second, let us decompose  $P_n\hat{f}_\epsilon - P_nf^\dagger$ :

$$P_n\hat{f}_\epsilon - P_nf^\dagger = (P_n\hat{f}_\epsilon - P_{n,\epsilon}\hat{f}_\epsilon) + (P_{n,\epsilon}\hat{f}_\epsilon - P_{n,\epsilon}f^\dagger) + (P_{n,\epsilon}f^\dagger - P_nf^\dagger).$$

By Lipschitz continuity, we have:

$$\forall \boldsymbol{\xi}, \mathbf{u} \in \mathbb{R}^p, \|\mathbf{u}\|_{\ell_2} \leq \epsilon: \ell(\boldsymbol{\xi}) - \ell(\boldsymbol{\xi} - \mathbf{u}) \leq L \|\boldsymbol{\xi} - (\boldsymbol{\xi} - \mathbf{u})\|_{\ell_2} \leq L\epsilon.$$

Consequently,  $\ell(\boldsymbol{\xi}) - \ell_\epsilon(\boldsymbol{\xi}) \leq L\epsilon$  and  $P_n\hat{f}_\epsilon - P_{n,\epsilon}\hat{f}_\epsilon \leq L\epsilon$ . In addition  $P_{n,\epsilon}\hat{f}_\epsilon - P_{n,\epsilon}f^\dagger \leq 0$  since  $\hat{f}_\epsilon$  is a minimizer of  $P_{n,\epsilon}$  over  $\mathcal{F}$ , and  $f^\dagger \in \mathcal{F}$ . Finally,  $P_{n,\epsilon}f^\dagger - P_nf^\dagger \leq 0$  since  $\ell$  upper bounds  $\ell_\epsilon$ . To summarize the second point,  $P_n\hat{f}_\epsilon - P_nf^\dagger \leq L\epsilon$ .

Third and last, by Hoeffding's inequality (Boucheron et al., 2013), with probability at least  $1 - \delta$ :

$$P_nf^\dagger - Pf^\dagger \leq LM\sqrt{\frac{\log(1/\delta)}{2n}}.$$

Gathering these three points with a union bound concludes the proof.

#### B.4. Algorithms

**Proof** (Lemma 3, body of the text) Let  $\phi: \mu \in [0, 1] \mapsto \left(1 + \frac{\lambda}{\|\mu\mathbf{y}\|_{-\mathbf{a}}\|_{\ell_2}}\right) \mu$ . First,  $\phi(1) = 1 + \frac{\lambda}{\|\mathbf{y}\|_{-\mathbf{a}}\|_{\ell_2}} \geq 1$ . Second, for  $\mu \geq 0$  sufficiently close to 0,  $[\mu\mathbf{y}]_{-\mathbf{a}}^{\mathbf{b}} = \mu\mathbf{y}$  (since entries of  $\mathbf{a}$  and  $\mathbf{b}$  are positive). Therefore  $\phi(0) = \lim_{\mu \downarrow 0} \left(\mu + \frac{\lambda\mu}{\mu\|\mathbf{y}\|_{\ell_2}}\right) = \frac{\lambda}{\|\mathbf{y}\|_{\ell_2}} \leq 1$ . Finally, since  $\phi$  is a continuous mapping on  $[0, 1]$  and  $1 \in [\phi(0), \phi(1)]$ , then the equation  $\phi(\mu) = 1$  has a solution in  $[0, 1]$ .  $\blacksquare$

**Proof** (Proposition 4, body of the text) The proof is in two part. First, we write optimality conditions for the proximal operator of interest, then we show that  $[\mu\mathbf{y}]_{-\mathbf{a}}^{\mathbf{b}}$  satisfies these optimality conditions when  $\mu$  is appropriately defined. From now on, let  $\mathbf{y} \in \mathbb{R}^p$ .

**Optimality conditions** Let  $\mathbf{x}^* = \text{prox}_{\lambda\|\cdot\|_{\ell_2} + \chi_{-\mathbf{a} \preccurlyeq \cdot \preccurlyeq \mathbf{b}}}(\mathbf{y}) = \arg \min_{-\mathbf{a} \preccurlyeq \mathbf{x} \preccurlyeq \mathbf{b}} \lambda \|\mathbf{x}\|_{\ell_2} + \frac{1}{2} \|\mathbf{y} - \mathbf{x}\|_{\ell_2}^2$ .

1. Assume that  $\mathbf{x}^* \neq 0$ . Then,  $\lambda \|\cdot\|_{\ell_2} + \frac{1}{2} \|\mathbf{y} - \cdot\|_{\ell_2}^2$  is differentiable at  $\mathbf{x}^*$  and for each coordinate  $j \in [p]$ , either:

- (a)  $-a_j < x_j^* < b_j$  and  $\left(1 + \frac{\lambda}{\|\mathbf{x}^*\|_{\ell_2}}\right) x_j^* = y_j$ ;
- (b) or  $x_j^* = b_j$  and  $\left(1 + \frac{\lambda}{\|\mathbf{x}^*\|_{\ell_2}}\right) x_j^* \leq y_j$ ;
- (c) or  $x_j^* = -a_j$  and  $\left(1 + \frac{\lambda}{\|\mathbf{x}^*\|_{\ell_2}}\right) x_j^* \geq y_j$ .

Gathering Conditions 1a-1c gives  $\|\mathbf{x}^*\|_{\ell_2} + \lambda \leq \|\mathbf{y}\|_{\ell_2}$ . Since  $\mathbf{x}^* \neq 0$ , we get  $\lambda < \|\mathbf{y}\|_{\ell_2}$ . Conversely, if  $\|\mathbf{y}\|_{\ell_2} \leq \lambda$ , then  $\mathbf{x}^* = 0$ .

2. If  $\mathbf{x}^* = 0$ , then  $\forall \delta > 0$  such that  $-\mathbf{a} \preccurlyeq \delta\mathbf{y} \preccurlyeq \mathbf{b}$ ,  $\lambda \|\delta\mathbf{y}\|_{\ell_2} + \frac{1}{2} \|\mathbf{y} - \delta\mathbf{y}\|_{\ell_2}^2 \geq \frac{1}{2} \|\mathbf{y}\|_{\ell_2}^2$ , that is  $\lambda \|\mathbf{y}\|_{\ell_2} \geq (1 - \frac{\delta}{2}) \|\mathbf{y}\|_{\ell_2}^2$ . Thus, by continuity when  $\delta \downarrow 0$ , we have  $\lambda \geq \|\mathbf{y}\|_{\ell_2}$ . To sum up,  $\mathbf{x}^* = 0$  if and only if  $\|\mathbf{y}\|_{\ell_2} \leq \lambda$ .

**Proximal solution** Let  $\mathbf{x} = [\mu\mathbf{y}]_{-\mathbf{a}}^{\mathbf{b}}$ , where  $\mu$  is defined in Proposition 4 from the body of the text. Assume that  $\|\mathbf{y}\|_{\ell_2} \leq \lambda$ , then  $\mu = 0$  and  $\mathbf{x} = 0$  satisfies the optimality conditions.

On the other hand, if  $\|\mathbf{y}\|_{\ell_2} > \lambda$ , then  $\left(1 + \frac{\lambda}{\|\mathbf{x}\|_{\ell_2}}\right) \mu = 1$ . As a result, either:

1.  $-a_j < x_j < b_j$ , so necessarily  $x_j = \mu y_j$  (otherwise it would be clipped to  $b_j$  or  $-a_j$ ). Therefore  $\left(1 + \frac{\lambda}{\|\mathbf{x}\|_{\ell_2}}\right) x_j = \left(1 + \frac{\lambda}{\|\mathbf{x}\|_{\ell_2}}\right) (\mu y_j) = y_j$ ;
2. or  $x_j = b_j$ , meaning that  $\mu y_j \geq b_j$ . So  $\left(1 + \frac{\lambda}{\|\mathbf{x}\|_{\ell_2}}\right) (\mu y_j) \geq \left(1 + \frac{\lambda}{\|\mathbf{x}\|_{\ell_2}}\right) b_j$ , that is  $\left(1 + \frac{\lambda}{\|\mathbf{x}\|_{\ell_2}}\right) x_j \leq y_j$ ;

3. or  $x_j = -a_j$ , meaning that  $\mu y_j \leq -a_j$ . So  $\left(1 + \frac{\lambda}{\|\mathbf{x}\|_{\ell_2}}\right)(\mu y_j) \leq \left(1 + \frac{\lambda}{\|\mathbf{x}\|_{\ell_2}}\right)(-a_j)$ , that is  $\left(1 + \frac{\lambda}{\|\mathbf{x}\|_{\ell_2}}\right)x_j \geq y_j$ .

Thus, when  $\|\mathbf{y}\|_{\ell_2} > \lambda$ ,  $\mathbf{x}$  satisfies the optimality conditions. This concludes the proof.  $\blacksquare$

**Corollary 6** Let two  $n$ -tuples  $\mathbf{A} = (\mathbf{a}_1, \dots, \mathbf{a}_n)$  and  $\mathbf{B} = (\mathbf{b}_1, \dots, \mathbf{b}_n)$  of vectors from  $\mathbb{R}^p$  with positive entries. For any  $n$ -tuple  $\mathbf{Y} = (\mathbf{y}_1, \dots, \mathbf{y}_n)$  of vectors from  $\mathbb{R}^p$ , let:

$$(\mathbf{x}_1, \dots, \mathbf{x}_n) = \text{prox}_{\lambda \|\cdot\|_{\ell_1/\ell_2} + \chi_{-\mathbf{A} \preccurlyeq \cdot \preccurlyeq \mathbf{B}}}(\mathbf{Y}),$$

where  $\|\mathbf{Y}\|_{\ell_1/\ell_2} = \sum_{i=1}^n \|\mathbf{y}_i\|_{\ell_2}$ . Then,  $\forall i \in [n]$ :

$$\mathbf{x}_i = \text{prox}_{\lambda \|\cdot\|_{\ell_2} + \chi_{-\mathbf{a}_i \preccurlyeq \cdot \preccurlyeq \mathbf{b}_i}}(\mathbf{y}_i).$$

**Proof** This is a direct consequence of the separability of  $\lambda \|\cdot\|_{\ell_1/\ell_2} + \chi_{-\mathbf{A} \preccurlyeq \cdot \preccurlyeq \mathbf{B}}$ .  $\blacksquare$

**Proof** (Proposition 5) The proof is similar to the one for Proposition 4 (see body of the text). Let  $\mathbf{y} \in \mathbb{R}^p$ .

**Optimality conditions** Let  $\mathbf{x}^* = \text{prox}_{\lambda(\|\cdot\|_{\ell_2} + \psi)}(\mathbf{y}) = \arg \min_{\mathbf{x} \in \mathbb{R}^p} \lambda \|\mathbf{x}\|_{\ell_2} + \lambda \psi(\mathbf{x}) + \frac{1}{2} \|\mathbf{y} - \mathbf{x}\|_{\ell_2}^2$ .

1. Assume that  $\mathbf{x}^* \neq 0$ . Then,  $\lambda \|\cdot\|_{\ell_2} + \lambda \psi + \frac{1}{2} \|\mathbf{y} - \cdot\|_{\ell_2}^2$  is differentiable at  $\mathbf{x}^*$  and for each coordinate  $j \in [p]$ :

$$y_j = \left( \frac{\lambda}{\|\mathbf{x}^*\|_{\ell_2}} + \lambda \left| \tau_j - \mathbf{I}_{x_j^* < 0} \right|^{-1} + 1 \right) x_j^*.$$

It appears that  $x_j^*$  and  $y_j$  have same sign. Therefore,  $\mathbf{I}_{x_j^* < 0} = \mathbf{I}_{y_j < 0}$  and

$$x_j^* = \left( \frac{\lambda}{\|\mathbf{x}^*\|_{\ell_2}} + \lambda \left| \tau_j - \mathbf{I}_{y_j < 0} \right|^{-1} + 1 \right)^{-1} y_j.$$

Now, the previous relation implies:

$$\|\mathbf{x}^*\|_{\ell_2}^2 = \sum_{j=1}^p \frac{y_j^2}{\left( \frac{\lambda}{\|\mathbf{x}^*\|_{\ell_2}} + \lambda \left| \tau_j - \mathbf{I}_{y_j < 0} \right|^{-1} + 1 \right)^2}.$$

Since  $\mathbf{x}^* \neq 0$ , we get:

$$1 = \sum_{j=1}^p \frac{y_j^2}{\left( \lambda + \|\mathbf{x}^*\|_{\ell_2} \left( \lambda \left| \tau_j - \mathbf{I}_{y_j < 0} \right|^{-1} + 1 \right) \right)^2}.$$

But  $\|\mathbf{x}^*\|_{\ell_2} > 0$ , so:

$$1 = \sum_{j=1}^p \frac{y_j^2}{\left(\lambda + \|\mathbf{x}^*\|_{\ell_2} \left(\lambda |\tau_j - \mathbf{I}_{y_j < 0}|^{-1} + 1\right)\right)^2} < \sum_{j=1}^p \frac{y_j^2}{\lambda^2},$$

that is  $\lambda < \|\mathbf{y}\|_{\ell_2}$ . Conversely, if  $\|\mathbf{y}\|_{\ell_2} \leq \lambda$ , then  $\mathbf{x}^* = 0$ .

2. If  $\mathbf{x}^* = 0$ , then  $\forall \delta > 0$ ,  $\lambda \|\delta \mathbf{y}\|_{\ell_2} + \lambda \psi(\delta \mathbf{y}) + \frac{1}{2} \|\mathbf{y} - \delta \mathbf{y}\|_{\ell_2}^2 \geq \frac{1}{2} \|\mathbf{y}\|_{\ell_2}^2$ , that is  $\lambda(\|\mathbf{y}\|_{\ell_2} + \delta \psi(\mathbf{y})) \geq (1 - \frac{\delta}{2}) \|\mathbf{y}\|_{\ell_2}^2$ . Thus, by continuity when  $\delta \downarrow 0$ , we have  $\lambda \geq \|\mathbf{y}\|_{\ell_2}$ . To sum up,  $\mathbf{x}^* = 0$  if and only if  $\|\mathbf{y}\|_{\ell_2} \leq \lambda$ .

**Proximal solution** If  $\|\mathbf{y}\|_{\ell_2} \leq \lambda$ , then  $\mathbf{x} = 0$  is satisfies trivially the optimality conditions.

On the other hand, if  $\|\mathbf{y}\|_{\ell_2} > \lambda$ , then  $\sum_{j=1}^p \frac{y_j^2}{\lambda^2} > 1$  and  $\lim_{\mu \rightarrow +\infty} \sum_{j=1}^p \frac{y_j^2}{\left(\mu \left(1 + \lambda |\tau_j - \mathbf{I}_{y_j < 0}|^{-1}\right) + \lambda\right)^2} = 0$ .

Thus, by continuity, Equation 1 has a solution  $\mu > 0$ . Let  $\mu$  be such a solution and let  $\mathbf{x} \in \mathbb{R}^p$  such that for each coordinate  $j \in [p]$ ,

$$x_j = \left(1 + \frac{\lambda}{\mu} + \lambda |\tau_j - \mathbf{I}_{y_j < 0}|^{-1}\right)^{-1} y_j.$$

Then:

$$\frac{\|\mathbf{x}\|_{\ell_2}^2}{\mu^2} = \sum_{j=1}^p \frac{y_j^2}{\mu^2 \left(1 + \frac{\lambda}{\mu} + \lambda |\tau_j - \mathbf{I}_{y_j < 0}|^{-1}\right)^2} = 1.$$

Consequently

$$x_j = \left(1 + \frac{\lambda}{\|\mathbf{x}\|_{\ell_2}} + \lambda |\tau_j - \mathbf{I}_{y_j < 0}|^{-1}\right)^{-1} y_j.$$

and  $\mathbf{x}$  satisfies the optimality conditions. This concludes the proof. ■

## Appendix C. Numerical experiments

Table 1 reports the average empirical loss (scaled by 100) along with the standard deviations. It completes Table 2 from the body of the text. For each dataset, the bold-face numbers are the two lowest losses. These values should be compared to the loss for  $\epsilon = 0$ .

## References

M.A. Alvarez, L. Rosasco, and N.D. Lawrence. Kernels for Vector-Valued Functions: a Review. *Foundations and Trends in Machine Learning*, 4(3):195–266, 2012.

F. Bach, R. Jenatton, J. Mairal, and G. Obozinski. Optimization with Sparsity-Inducing Penalties. *Foundations and Trends in Machine Learning*, 4(1):1–106, January 2012.

Table 1: Empirical pinball loss  $\times 100$  along with percentage of support vectors (the less, the better).

Data set	$\epsilon = 0$	$\epsilon = 0.05$	$\epsilon = 0.1$	$\epsilon = 0.2$	$\epsilon = 0.5$	$\epsilon = 0.75$	$\epsilon = 1$	$\epsilon = 1.5$	$\epsilon = 2$	$\epsilon = 3$
caution	67.50 $\pm$ 12.96	<b>67.40</b> $\pm$ 13.12	<b>67.17</b> $\pm$ 12.68	67.54 $\pm$ 13.00	69.93 $\pm$ 12.51	73.24 $\pm$ 13.10	76.80 $\pm$ 12.51	83.42 $\pm$ 14.97	100.22 $\pm$ 16.63	142.19 $\pm$ 16.11
ftcollinsnow	<b>109.07</b> $\pm$ 5.88	109.12 $\pm$ 5.95	109.14 $\pm$ 6.00	109.15 $\pm$ 6.00	109.11 $\pm$ 6.30	110.39 $\pm$ 7.10	<b>109.05</b> $\pm$ 6.72	110.90 $\pm$ 6.87	109.81 $\pm$ 6.22	113.50 $\pm$ 9.29
highway	79.29 $\pm$ 17.06	78.10 $\pm$ 16.15	76.75 $\pm$ 15.42	76.66 $\pm$ 19.11	<b>75.09</b> $\pm$ 18.82	<b>70.94</b> $\pm$ 21.58	75.10 $\pm$ 21.02	97.67 $\pm$ 22.91	112.30 $\pm$ 20.10	112.09 $\pm$ 20.12
heights	91.05 $\pm$ 1.12	91.00 $\pm$ 1.20	90.98 $\pm$ 1.19	<b>90.98</b> $\pm$ 1.21	91.18 $\pm$ 1.09	91.21 $\pm$ 1.27	<b>90.98</b> $\pm$ 1.13	91.09 $\pm$ 1.02	91.51 $\pm$ 1.20	93.34 $\pm$ 1.84
sniffer	32.34 $\pm$ 5.14	<b>31.40</b> $\pm$ 4.81	32.31 $\pm$ 5.19	<b>31.40</b> $\pm$ 2.98	34.64 $\pm$ 3.85	39.84 $\pm$ 5.03	41.82 $\pm$ 4.14	52.06 $\pm$ 5.88	62.21 $\pm$ 11.77	103.76 $\pm$ 16.42
snowgeese	<b>49.62</b> $\pm$ 19.38	<b>50.51</b> $\pm$ 18.23	51.25 $\pm$ 18.64	51.08 $\pm$ 17.69	52.88 $\pm$ 15.11	53.81 $\pm$ 13.97	62.81 $\pm$ 19.66	90.15 $\pm$ 23.97	107.53 $\pm$ 23.13	94.25 $\pm$ 24.65
ufc	57.87 $\pm$ 3.09	57.90 $\pm$ 3.07	<b>57.78</b> $\pm$ 2.99	57.84 $\pm$ 3.01	<b>57.67</b> $\pm$ 2.84	57.84 $\pm$ 2.76	58.19 $\pm$ 2.92	61.04 $\pm$ 3.68	66.81 $\pm$ 4.00	86.23 $\pm$ 4.79
birthwt	99.93 $\pm$ 8.51	99.95 $\pm$ 8.53	99.93 $\pm$ 8.63	99.70 $\pm$ 8.64	<b>99.25</b> $\pm$ 8.66	100.50 $\pm$ 10.31	99.80 $\pm$ 11.29	<b>98.71</b> $\pm$ 10.04	99.56 $\pm$ 9.43	103.39 $\pm$ 8.83
crabs	8.59 $\pm$ 0.66	<b>8.52</b> $\pm$ 0.68	<b>8.49</b> $\pm$ 0.73	9.44 $\pm$ 0.57	19.94 $\pm$ 1.38	23.08 $\pm$ 1.59	31.44 $\pm$ 3.75	44.08 $\pm$ 4.64	53.45 $\pm$ 5.65	86.91 $\pm$ 9.48
GAGurine	44.30 $\pm$ 5.85	<b>44.26</b> $\pm$ 5.79	<b>44.25</b> $\pm$ 5.76	44.86 $\pm$ 6.04	46.20 $\pm$ 5.35	49.87 $\pm$ 4.88	52.88 $\pm$ 3.94	57.06 $\pm$ 3.47	65.89 $\pm$ 3.88	103.32 $\pm$ 24.62
geyser	<b>77.81</b> $\pm$ 5.36	78.15 $\pm$ 5.39	<b>78.12</b> $\pm$ 5.38	78.45 $\pm$ 5.35	78.40 $\pm$ 5.77	78.28 $\pm$ 5.88	78.54 $\pm$ 5.82	80.55 $\pm$ 6.34	85.15 $\pm$ 6.18	99.92 $\pm$ 8.65
gilgais	<b>32.96</b> $\pm$ 4.09	<b>33.12</b> $\pm$ 3.99	33.27 $\pm$ 4.11	33.42 $\pm$ 3.88	35.08 $\pm$ 3.35	36.62 $\pm$ 3.59	37.94 $\pm$ 3.68	48.17 $\pm$ 9.44	94.65 $\pm$ 4.98	104.12 $\pm$ 5.92
topo	47.49 $\pm$ 7.93	48.93 $\pm$ 7.43	48.74 $\pm$ 7.10	48.17 $\pm$ 7.01	<b>41.65</b> $\pm$ 5.60	<b>45.24</b> $\pm$ 3.53	51.19 $\pm$ 7.92	53.68 $\pm$ 8.39	58.21 $\pm$ 13.35	80.57 $\pm$ 15.18
BostonHousing	<b>34.54</b> $\pm$ 3.34	34.68 $\pm$ 3.46	34.70 $\pm$ 3.39	<b>34.09</b> $\pm$ 3.37	35.27 $\pm$ 3.02	37.65 $\pm$ 3.18	41.31 $\pm$ 3.41	55.04 $\pm$ 5.61	73.39 $\pm$ 12.35	112.22 $\pm$ 12.91
CobarOre	<b>0.50</b> $\pm$ 0.38	<b>5.05</b> $\pm$ 1.90	8.75 $\pm$ 3.44	12.47 $\pm$ 4.27	23.84 $\pm$ 6.03	35.82 $\pm$ 8.20	47.35 $\pm$ 10.94	66.15 $\pm$ 14.56	84.51 $\pm$ 17.70	106.89 $\pm$ 15.52
engel	43.57 $\pm$ 6.05	43.50 $\pm$ 6.02	<b>43.47</b> $\pm$ 6.08	<b>43.44</b> $\pm$ 5.99	57.36 $\pm$ 46.14	43.98 $\pm$ 5.37	46.31 $\pm$ 6.29	53.15 $\pm$ 5.45	69.43 $\pm$ 9.22	100.48 $\pm$ 11.63
mcycle	<b>63.95</b> $\pm$ 5.25	<b>63.88</b> $\pm$ 5.20	64.26 $\pm$ 5.99	64.90 $\pm$ 6.68	65.89 $\pm$ 5.89	67.29 $\pm$ 6.13	70.11 $\pm$ 7.65	74.78 $\pm$ 6.43	86.49 $\pm$ 6.77	109.79 $\pm$ 12.67
BigMac2003	<b>49.94</b> $\pm$ 12.85	<b>49.97</b> $\pm$ 12.84	50.00 $\pm$ 12.83	50.27 $\pm$ 13.19	51.16 $\pm$ 13.37	51.44 $\pm$ 10.57	53.63 $\pm$ 14.29	77.40 $\pm$ 24.48	106.38 $\pm$ 13.97	136.76 $\pm$ 61.70
UN3	71.27 $\pm$ 4.69	<b>70.94</b> $\pm$ 4.57	<b>71.03</b> $\pm$ 4.68	71.49 $\pm$ 5.06	71.37 $\pm$ 5.01	71.53 $\pm$ 5.90	72.68 $\pm$ 6.17	76.72 $\pm$ 6.13	84.50 $\pm$ 7.00	109.59 $\pm$ 4.71
cpus	<b>11.31</b> $\pm$ 9.32	<b>13.32</b> $\pm$ 8.95	15.57 $\pm$ 9.16	20.16 $\pm$ 8.06	25.88 $\pm$ 8.93	35.66 $\pm$ 11.61	55.27 $\pm$ 14.69	65.05 $\pm$ 9.70	65.02 $\pm$ 9.65	

P.L. Bartlett and S. Mendelson. Rademacher and Gaussian Complexities: Risk Bounds and Structural Results. *Journal of Machine Learning Research*, 3(Nov):463–482, 2002.

S. Boucheron, G. Lugosi, and P. Massart. *Concentration Inequalities: A Nonasymptotic Theory of Independence*. Oxford University Press, Oxford, New York, 2013.

S.P. Boyd and L. Vandenberghe. *Convex Optimization*. Cambridge University Press, 2004.

A. Maurer. A vector-contraction inequality for Rademacher complexities. In *Proceedings of The 27th International Conference on Algorithmic Learning Theory*, 2016.

M. Mohri, A. Rostamizadeh, and A. Talwalkar. *Foundations of Machine Learning*. The MIT Press, 2012.

M. Sangnier, O. Fercoq, and F. d’Alché Buc. Joint quantile regression in vector-valued RKHSs. In *Advances in Neural Information Processing Systems 29*, 2016.

I. Takeuchi, Q.V. Le, T.D. Sears, and A.J. Smola. Nonparametric Quantile Estimation. *Journal of Machine Learning Research*, 7:1231–1264, 2006.