

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

**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

## Abstract

Leveraging the celebrated [support vector regression \(SVR\)](#) method, we propose a unifying framework in order to deliver regression machines in [reproducing kernel Hilbert spaces \(RKHSs\)](#) with data sparsity. The central point is a new definition of  $\epsilon$ -insensitivity, valid for many regression losses (including quantile and expectile regression) and their multivariate extensions. We show that the dual optimization problem to empirical risk minimization with  $\epsilon$ -insensitivity involves a data sparse regularization. We also provide an analysis of the excess of risk as well as a randomized coordinate descent algorithm for solving the dual. Numerical experiments validate our approach.

**Keywords:** Quantile regression, Expectile regression, Operator-valued kernel.

## 1. Introduction

In supervised learning, sparsity is a common feature ([Hastie et al., 2015](#)). When it is put on the table, the first reaction is certainly to invoke feature selection ([Tibshirani, 1996](#)). In many domains such as finance and biology, feature selection aims at discovering the covariates that really explain the outcomes (hopefully in a causal manner). However, when dealing with nonparametric modeling, sparsity may also occur in the orthogonal direction: data, which support ultimately the desired estimator. Data sparsity is mainly valuable for two reasons: first, controlling the number of training data used in the model can prevent from overfitting. Second, data sparsity implies less data to store and less time needed for prediction as well as potentially less time needed for training. The main evidence of this phenomenon lies in the celebrated [support vector machines \(SVMs\)](#) ([Boser et al., 1992](#); [Vapnik, 2010](#); [Mohri et al., 2012](#)). The principle of margin maximization results in selecting only a few data points that support the whole prediction function, leaving the others sink into oblivion.

Data sparsity is particularly known for classification thanks to [SVMs](#) but is less predominant for regression. Indeed, while classifying relies on the few points close to the frontier (they are likely to be misclassified), regression is supported by points for which the outcome is far from the conditional expectation. Those points, that are likely to be poorly characterized by the prediction, gather most of the training points. Thus, by nature, regressors tend

to be less data sparse than classifiers. In order to obtain the valued data sparsity for a wide class of regression problems, we consider regression relying on empirical risk minimization (Vapnik, 2010) with non-linear methods based on RKHSs for uni and multidimensional outcomes (Micchelli and Pontil, 2005).

The companion of SVM for regression, SVR (Drucker et al., 1997), is probably the most relevant representative of data sparse regression methods in the frequentist domain (in the Bayesian literature, relevance vector machines can deliver sparse regressors (Tipping, 2001)). SVR makes use of the celebrated  $\epsilon$ -insensitive loss  $\xi \in \mathbb{R} \mapsto \max(0, |\xi| - \epsilon)$ , where  $\epsilon > 0$ , in order to produce data sparsity. In practice, residues close to 0 are not penalized, and consequently corresponding points are not considered in the prediction function. Later, Park and Kim (2011) extended this  $\epsilon$ -insensitive loss to quantile regression (QR) (Koenker, 2005) in an *ad hoc* manner.

These approaches are in deep contrast compared to kernel ridge regression (using the loss  $\xi \in \mathbb{R} \mapsto \frac{1}{2}\xi^2$ ), that does not benefit from data sparsity. To remedy that lack of sparsity, Lim et al. (2014) leveraged the representer theorem (Schölkopf et al., 2001) for ridge regression and penalized the coefficients associated to each training point. The same procedure was used by Zhang et al. (2016) for QR but with a sparse constraint instead of a sparse regularization. Contrarily to  $\epsilon$ -insensitive losses, the latter approaches constrain the hypotheses while keeping the vanilla regression losses.

The previous works are either restricted to unidimensional losses (Drucker et al., 1997; Park and Kim, 2011; Zhang et al., 2016) or compelled to assume a finite representation of the optimal regressor (Lim et al., 2014; Zhang et al., 2016). Our contribution in this context is to provide a novel and unifying view of  $\epsilon$ -insensitive losses for data sparse kernel regression, that is theoretically well founded (Section 3). In the spirit of SVR, we show that, for a large number of regression problems (uni or multivariate, quantile or traditional regression), data sparsity can be obtained by properly defining a  $\epsilon$ -insensitive variant of the loss in the data-fitting term which turns to be a data sparse regularization in the dual. This framework contains SVR as a special case but also a new loss for QR (that is different from the one introduced by Park and Kim (2011)) and expectile regression (Newey and Powell, 1987), as well as their multivariate extensions (Section 4) in the context of vector-valued RKHSs. We also provide an analysis of the excess of risk (Section 5). Special attention is paid to the dual form of the optimization problem and a randomized primal-dual block coordinate descent algorithm for estimating our data sparse regressor (Section 6) is derived. Numerical experiments conclude the discussion (Section 7).

## 2. Framework

**Notations** In the whole paper, we denote vectors with bold-face letters. Moreover,  $\tau \in (0, 1)$  is a quantile level,  $\boldsymbol{\tau} \in (0, 1)^p$  is a vector of quantile levels,  $n$  and  $p$  are two positive integers (samples size and dimension),  $[n]$  is the range of integers between 1 and  $n$ ,  $\boldsymbol{\alpha}^l \in \mathbb{R}^n$  denotes the  $l^{\text{th}}$  row vector of any  $\boldsymbol{\alpha} = (\boldsymbol{\alpha}_i)_{1 \leq i \leq n} \in (\mathbb{R}^p)^n$ ,  $\mathbb{I}$  is the indicator function (set to 1 when the condition is fulfilled and 0 otherwise),  $\mathbf{1}$  denotes the all-ones vector,  $\chi$  is the characteristic function (set to 0 when the condition is met and to  $\infty$  otherwise),  $\preceq$  (respectively  $\prec$ ) denotes the pointwise (respectively strict) inequality,  $\text{diag}$  is the operator mapping a vector to a diagonal matrix,  $\partial\psi$  is the subdifferential of any function  $\psi$  (when it

exists),  $\text{prox}_\psi$  denotes its proximal operator (Bauschke and Combettes, 2011),  $\text{proj}_{\mathbf{1}}$  is the projector onto the span of  $\mathbf{1}$ .

**General framework** Let  $\mathcal{X}$  be a nonempty input space and  $\mathbb{R}^p$  be the output Hilbert space (for a given integer  $p$ ). Let also  $\mathcal{H}$  be a Hilbert space of hypotheses  $h: \mathcal{X} \rightarrow \mathbb{R}^p$  and  $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{1 \leq i \leq n}$  be a training sample of  $n$  couples  $(\mathbf{x}_i, \mathbf{y}_i) \in \mathcal{X} \times \mathbb{R}^p$ . We consider the setting of regularized empirical risk minimization for regression, based on a real-valued convex loss function  $\ell: \mathbb{R}^p \rightarrow \mathbb{R}$ . This is to provide a minimizer to the optimization problem:

$$\underset{h \in \mathcal{H}, \mathbf{b} \in \mathbb{R}^p}{\text{minimize}} \quad \frac{\lambda}{2} \|h\|_{\mathcal{H}}^2 + \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{y}_i - (h(\mathbf{x}_i) + \mathbf{b})), \quad (\text{P1})$$

where  $\lambda > 0$  is a trade-off parameter and  $\|\cdot\|_{\mathcal{H}}$  is the norm associated to  $\mathcal{H}$ . Problem (P1) consists in minimizing a data-fitting term  $\frac{1}{n} \sum_{i=1}^n \ell(\boldsymbol{\xi}_i)$ , where the residues  $\boldsymbol{\xi}_i = \mathbf{y}_i - (h(\mathbf{x}_i) + \mathbf{b})$  are driven close to zero. Here, the prediction function (or regressor) is  $f: \mathbf{x} \in \mathcal{X} \mapsto h(\mathbf{x}) + \mathbf{b}$ . It is made of a functional part  $h \in \mathcal{H}$  and of an intercept  $\mathbf{b} \in \mathbb{R}^p$  (such as what we can find in SVMs). The last component of Problem (P1),  $\frac{\lambda}{2} \|\cdot\|_{\mathcal{H}}^2$ , is a regularizer, which penalizes functions with high complexities, allowing for generalizing on unknown observations (Mohri et al., 2012).

Let, for any  $\mathbf{x} \in \mathcal{X}$ ,  $E_{\mathbf{x}}: h \in \mathcal{H} \mapsto h(\mathbf{x})$  be the evaluation map and  $E_{\mathbf{x}}^*$  be its adjoint operator. A common theorem in nonparametric regression, known as the representer theorem (Schölkopf et al., 2001), states that any estimator based on (P1) is a finite expansion of local functions supported by the training examples. More formally, the functional part of any solution  $(\hat{h}, \hat{\mathbf{b}})$  of Problem (P1) admits a representation of the form  $\hat{h} = \sum_{i=1}^n E_{\mathbf{x}_i}^* \boldsymbol{\alpha}_i$ , where  $\boldsymbol{\alpha}_i \in \mathbb{R}^p$  for all  $i \in [n]$ .

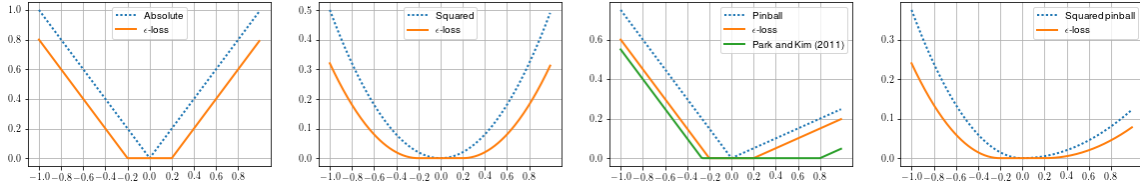
A proof of this statement based on Karush-Kuhn-Tucker (KKT) conditions can be found in Supplementary B, but many proofs already exist for kernel methods in the scalar ( $p = 1$ ) and the vectorial ( $p > 1$ ) cases (Schölkopf et al., 2001; Brouard et al., 2016). The representer theorem states that  $\hat{h}$  can be parametrized by a finite number of  $\boldsymbol{\alpha}_i \in \mathbb{R}^p$  (which boil down to be Lagrange multipliers), such that each  $\boldsymbol{\alpha}_i$  is associated to a training point  $\mathbf{x}_i$ . Therefore, as soon as  $\boldsymbol{\alpha}_i = 0$ ,  $\mathbf{x}_i$  does not appear in  $\hat{h}$  (that is,  $\mathbf{x}_i$  is not a support vector). When many  $\boldsymbol{\alpha}_i = 0$ , we say that  $\hat{h}$  is data sparse. This feature is very attractive since it makes possible to get rid of useless training points when computing  $\hat{h}(\mathbf{x})$ , and speeds up the prediction task. The next section introduces a way to get such a sparsity.

### 3. $\epsilon$ -insensitive losses and the associated minimization problem

In the whole discussion, we will assume that the loss function  $\ell$  is convex, lower semi-continuous and has a unique minimum at 0 with null value. Then, we develop the idea that data sparsity can be obtained by substituting the original loss function  $\ell$  by a slightly different version  $\ell_\epsilon$ , where  $\ell_\epsilon$  is a kind of *soft-thresholding* of  $\ell$ . More formally, let  $\epsilon$  be a positive parameter and define the  $\epsilon$ -insensitive the loss  $\ell_\epsilon$  by:

$$\forall \boldsymbol{\xi} \in \mathbb{R}^p: \quad \ell_\epsilon(\boldsymbol{\xi}) = \begin{cases} 0 & \text{if } \|\boldsymbol{\xi}\|_{\ell_2} \leq \epsilon \\ \inf_{\mathbf{d} \in \mathbb{R}^p: \|\mathbf{d}\|_{\ell_2} = 1} \ell(\boldsymbol{\xi} - \epsilon \mathbf{d}) & \text{otherwise.} \end{cases}$$

	NAME	LOSS $\ell(\boldsymbol{\xi})$	$\epsilon$ -INSENSITIVE LOSS $\ell_\epsilon(\boldsymbol{\xi})$	CONJUGATE LOSS $\ell^*(\boldsymbol{\alpha})$
UNIDIMENSIONAL ( $\boldsymbol{\xi} \in \mathbb{R}, \boldsymbol{\alpha} \in \mathbb{R}$ )	Absolute	$ \boldsymbol{\xi} $	$\max(0,  \boldsymbol{\xi}  - \epsilon)$	$\chi_{-1 \leq \alpha \leq 1}$
	Squared	$\frac{1}{2} \boldsymbol{\xi}^2$	$\frac{1}{2} (\max(0,  \boldsymbol{\xi}  - \epsilon))^2$	$\frac{1}{2} \alpha^2$
	Pinball	$(\tau - \mathbb{I}_{\boldsymbol{\xi} < 0}) \boldsymbol{\xi}$	$\begin{cases} 0 & \text{if }  \boldsymbol{\xi}  \leq \epsilon \\ \tau(\boldsymbol{\xi} - \epsilon) & \text{if } \boldsymbol{\xi} \geq \epsilon \\ (\tau - 1)(\boldsymbol{\xi} + \epsilon) & \text{if } \boldsymbol{\xi} \leq -\epsilon \end{cases}$	$\chi_{\tau - 1 \leq \alpha \leq \tau}$
	Squared pinball	$\frac{1}{2}  \tau - \mathbb{I}_{\boldsymbol{\xi} < 0}  \boldsymbol{\xi}^2$	$\begin{cases} 0 & \text{if }  \boldsymbol{\xi}  \leq \epsilon \\ \frac{\tau}{2} (\boldsymbol{\xi} - \epsilon)^2 & \text{if } \boldsymbol{\xi} \geq \epsilon \\ \frac{1-\tau}{2} (\boldsymbol{\xi} + \epsilon)^2 & \text{if } \boldsymbol{\xi} \leq -\epsilon \end{cases}$	$\frac{1}{2}  \tau - \mathbb{I}_{\alpha < 0} ^{-1} \alpha^2$
MULTIDIMENSIONAL ( $\boldsymbol{\xi} \in \mathbb{R}^p, \boldsymbol{\alpha} \in \mathbb{R}^p$ )	$\ell_1$ -norm	$\ \boldsymbol{\xi}\ _{\ell_1}$	No closed-form found	$\chi_{-1 \leq \alpha \leq 1}$
	$\ell_2$ -norm	$\frac{1}{2} \ \boldsymbol{\xi}\ _{\ell_2}^2$	$\frac{1}{2} \left\  \boldsymbol{\xi} - \min(\ \boldsymbol{\xi}\ _{\ell_2}, \epsilon) \frac{\boldsymbol{\xi}}{\ \boldsymbol{\xi}\ _{\ell_2}} \right\ _{\ell_2}^2$	$\frac{1}{2} \ \boldsymbol{\alpha}\ _{\ell_2}^2$
	Multiple pinball	$\sum_{j=1}^p (\tau_j - \mathbb{I}_{\boldsymbol{\xi}_j < 0}) \boldsymbol{\xi}_j$	No closed-form found	$\chi_{\tau - 1 \leq \alpha \leq \tau}$
	Multiple squared pinball	$\frac{1}{2} \sum_{j=1}^p  \tau_j - \mathbb{I}_{\boldsymbol{\xi}_j < 0}  \boldsymbol{\xi}_j^2$	No closed-form found	$\frac{1}{2} \sum_{j=1}^p  \tau_j - \mathbb{I}_{\alpha_j < 0} ^{-1} \alpha_j^2$

 Table 1: Examples of  $\epsilon$ -insensitive losses.

 Figure 1: Examples of unidimensional  $\epsilon$ -losses ( $\epsilon = 0.2, \tau = 0.25$ )

Put another way, the new loss value associated to a residue  $\boldsymbol{\xi}$  is set to 0 if the magnitude of the residue is sufficiently small, or to the smallest loss value in a neighborhood centered at the original residue. Let us remark that, by the previous assumptions,  $\ell_\epsilon$  is convex (see Supplementary B). To illustrate this definition, Table 1 as well as Figure 1 provide several examples of  $\epsilon$ -insensitive losses, particularly for regression, QR and expectile regression. Before developing further these examples in the next section, we explain why using this  $\epsilon$ -insensitive loss helps in getting data sparsity for the regressor.

**Proposition 1 (Dual optimization problem)** *A dual to the learning problem*

$$\underset{h \in \mathcal{H}, \mathbf{b} \in \mathbb{R}^p}{\text{minimize}} \quad \frac{\lambda}{2} \|h\|_{\mathcal{H}}^2 + \frac{1}{n} \sum_{i=1}^n \ell_\epsilon(\mathbf{y}_i - (h(\mathbf{x}_i) + \mathbf{b})) \quad (\text{P2})$$

is

$$\begin{aligned} \underset{\forall i \in [n], \boldsymbol{\alpha}_i \in \mathbb{R}^p}{\text{minimize}} \quad & \frac{1}{n} \sum_{i=1}^n \ell^*(\boldsymbol{\alpha}_i) + \frac{1}{2\lambda n^2} \sum_{i,j=1}^n \langle \boldsymbol{\alpha}_i, E_{\mathbf{x}_i} E_{\mathbf{x}_j}^* \boldsymbol{\alpha}_j \rangle_{\ell_2} - \frac{1}{n} \sum_{i=1}^n \langle \boldsymbol{\alpha}_i, \mathbf{y}_i \rangle_{\ell_2} + \frac{\epsilon}{n} \sum_{i=1}^n \|\boldsymbol{\alpha}_i\|_{\ell_2} \\ \text{s. t.} \quad & \sum_{i=1}^n \boldsymbol{\alpha}_i = \mathbf{0}, \end{aligned} \quad (\text{P3})$$

where  $\ell^*: \boldsymbol{\alpha} \in \mathbb{R}^p \mapsto \sup_{\boldsymbol{\xi} \in \mathbb{R}^p} \langle \boldsymbol{\alpha}, \boldsymbol{\xi} \rangle_{\ell_2} - \ell(\boldsymbol{\xi})$  is the Fenchel-Legendre transform of  $\ell$ . In addition, let  $(\hat{\boldsymbol{\alpha}}_i)_{i \in [n]}$  be solutions to Problem (P3). Then, the function  $\hat{h} = \frac{1}{\lambda n} \sum_{i=1}^n E_{\mathbf{x}}^* \hat{\boldsymbol{\alpha}}_i$  ( $\hat{h} \in \mathcal{H}$ ) is solution to Problem (P2) for a given intercept  $\hat{\mathbf{b}} \in \mathbb{R}^p$ .

We see that turning  $\ell$  into an  $\epsilon$ -insensitive loss leads to blending a sparsity regularization in the learning process: the  $\ell_1/\ell_2$ -norm  $\sum_{i=1}^n \|\alpha_i\|_{\ell_2}$ , which is known to induce sparsity in the vectors  $\alpha_i$  (Bach et al., 2012). Even though a proof of this proposition is proposed in Supplementary B, we can grasp this phenomenon by remarking that  $\ell_\epsilon$  turns out to be an infimal convolution (Bauschke and Combettes, 2011), noted  $\square$ , of two functions. Thus, its Fenchel-Legendre transform is the sum of Fenchel-Legendre transforms of each contribution (Bauschke and Combettes, 2011):

$$\ell_\epsilon^* = \left( \ell \square \chi_{\|\cdot\|_{\ell_2} \leq \epsilon} \right)^* = \ell^* + \left( \chi_{\|\cdot\|_{\ell_2} \leq \epsilon} \right)^* = \ell^* + \epsilon \|\cdot\|_{\ell_2}.$$

The next section discusses several applications of  $\epsilon$ -insensitive losses.

#### 4. Examples of $\epsilon$ -insensitive losses

To illustrate the link between  $\epsilon$ -insensitive losses and data sparsity, we focus on RKHSs  $\mathcal{H}$ , based either on a scalar-valued kernel  $k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  (Steinwart and Christmann, 2008) or on a matrix-valued kernel (MVK)  $K: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}^{p \times p}$  (Micchelli and Pontil, 2005). Therefore, the base functions  $E_{\mathbf{x}}^* \alpha_i$ , appearing in Proposition 1 become  $\mathbf{x}' \in \mathcal{X} \mapsto \alpha_i k(\mathbf{x}', \mathbf{x}) \in \mathbb{R}$  ( $\alpha_i \in \mathbb{R}$ ) and  $\mathbf{x}' \in \mathcal{X} \mapsto K(\mathbf{x}', \mathbf{x}) \alpha_i \in \mathbb{R}^p$  ( $\alpha_i \in \mathbb{R}^p$ ), respectively for the uni and multidimensional case.

##### 4.1. Least absolute deviation

The notion of  $\epsilon$ -insensitive loss is very well illustrated in SVR: let us study the problem of least absolute deviation, that is, solving (P1), where there is a single output ( $h: \mathcal{X} \rightarrow \mathbb{R}$ ,  $y_i \in \mathbb{R}, b \in \mathbb{R}, \alpha_i \in \mathbb{R}$ ) and the loss to be minimized is  $\ell(\xi) = |\xi|$ . In this situation, it is easy to check that  $\ell_\epsilon$  boils down to:  $\ell_\epsilon(\xi) = \max(0, |\xi| - \epsilon)$ , which is the well known  $\epsilon$ -insensitive loss considered in SVR (Hastie et al., 2009).

Now, since  $\ell^*(\alpha) = 0$  when  $-1 \leq \alpha \leq 1$  and  $+\infty$  otherwise, Problem (P3) becomes (up to normalizing the objective function by  $n$ ):

$$\begin{aligned} & \underset{\forall i \in [n], \alpha_i \in \mathbb{R}}{\text{minimize}} && \frac{1}{2\lambda n} \sum_{i,j=1}^n \alpha_i \alpha_j k(\mathbf{x}_i, \mathbf{x}_j) - \sum_{i=1}^n \alpha_i y_i + \epsilon \sum_{i=1}^n |\alpha_i| \\ & \text{s. t.} && \forall i \in [n], -1 \leq \alpha_i \leq 1, \quad \sum_{i=1}^n \alpha_i = 0. \end{aligned}$$

Using that  $|\alpha| = \inf_{\alpha^+, \alpha^- \geq 0} \{\alpha^+ + \alpha^-\}$ , we can introduce auxiliary variables  $\alpha_i^+ \geq 0$  and  $\alpha_i^- \geq 0$ , decouple  $\alpha_i$  in  $\alpha_i^+ - \alpha_i^-$  and replace  $|\alpha_i|$  by  $\alpha_i^+ + \alpha_i^-$ . This turns the previous optimization problem into:

$$\begin{aligned} & \underset{\forall i \in [n], \alpha_i^+, \alpha_i^- \in \mathbb{R}}{\text{minimize}} && \frac{1}{2\lambda n} \sum_{i,j=1}^n (\alpha_i^+ - \alpha_i^-)(\alpha_j^+ - \alpha_j^-) k(\mathbf{x}_i, \mathbf{x}_j) - \sum_{i=1}^n (\alpha_i^+ - \alpha_i^-) y_i + \epsilon \sum_{i=1}^n (\alpha_i^+ + \alpha_i^-) \\ & \text{s. t.} && \forall i \in [n], 0 \leq \alpha_i^+ \leq 1, \quad 0 \leq \alpha_i^- \leq 1, \quad \sum_{i=1}^n \alpha_i^+ - \alpha_i^- = 0, \end{aligned}$$

which is the well known dual of **SVR**. In this perspective, **SVR** is a data sparse version of least absolute deviation, for which the sparsity inducing norm appears in the dual.

## 4.2. Least mean squares

Besides **SVR**, Kernel Ridge Regression is also well known among kernel methods. In particular, for the general case of multivariate regression ( $\mathbf{y}_i \in \mathbb{R}^p$ ), we aim at minimizing the squared loss  $\ell(\cdot) = \frac{1}{2} \|\cdot\|_{\ell_2}^2$ . Referring to Table 1, Problem (P3) then becomes:

$$\begin{aligned} & \underset{\forall i \in [n], \boldsymbol{\alpha}_i \in \mathbb{R}^p}{\text{minimize}} && \frac{1}{2} \sum_{i,j=1}^n \left\langle \boldsymbol{\alpha}_i, \left( \mathbf{I}_p + \frac{1}{\lambda n} K(\mathbf{x}_i, \mathbf{x}_j) \right) \boldsymbol{\alpha}_j \right\rangle_{\ell_2} - \sum_{i=1}^n \langle \boldsymbol{\alpha}_i, \mathbf{y}_i \rangle_{\ell_2} + \epsilon \sum_{i=1}^n \|\boldsymbol{\alpha}_i\|_{\ell_2} \\ & \text{s. t.} && \sum_{i=1}^n \boldsymbol{\alpha}_i = 0, \end{aligned} \quad (\text{P4})$$

where  $\mathbf{I}_p$  is the  $p$ -dimensional identity matrix. This is essentially similar to the dual of kernel ridge regression but with an extra  $\ell_1/\ell_2$ -norm on Lagrange multipliers.

As far as we know, data sparsity for multivariate kernel ridge regression has been first introduced by [Lim et al. \(2014\)](#). In this work, the authors assume first that the optimal predictor  $\hat{h}$  has a finite representation  $\hat{h} = \frac{1}{\lambda n} \sum_{i=1}^n K(\cdot, \mathbf{x}_i) \mathbf{c}_i$  without relying on a representer theorem and regularize the vectors of weights  $\mathbf{c}_i$  associated to each point thanks to an  $\ell_1/\ell_2$ -norm. They also dismiss the intercept  $\mathbf{b}$  for simplicity. The learning problem, as presented in ([Lim et al., 2014](#)) (with normalization chosen at purpose), is:

$$\begin{aligned} & \underset{\substack{h \in \mathcal{H}, \\ \forall i \in [n], \mathbf{c}_i \in \mathbb{R}^p}}{\text{minimize}} && \frac{\lambda}{2} \|h\|_{\mathcal{H}}^2 + \frac{1}{2n} \sum_{i=1}^n \|\mathbf{y}_i - h(\mathbf{x}_i)\|_{\ell_2}^2 + \frac{\epsilon}{\lambda n^2} \sum_{i=1}^n \|\mathbf{c}_i\|_{\ell_2} \\ & \text{s. t.} && h = \frac{1}{\lambda n} \sum_{i=1}^n K(\cdot, \mathbf{x}_i) \mathbf{c}_i. \end{aligned}$$

For the sake of readability, we write  $\underline{\mathbf{c}} = (\mathbf{c}_1^\top, \dots, \mathbf{c}_n^\top)^\top$  the vector of all weights and  $\|\underline{\mathbf{c}}\|_{\ell_1/\ell_2} = \sum_{i=1}^n \|\mathbf{c}_i\|_{\ell_2}$ . Let also  $\underline{K} = (K(\mathbf{x}_i, \mathbf{x}_j))_{1 \leq i, j \leq n}$  be the **symmetric positive semi-definite (PSD)** kernel matrix. Fermat's rule indicates that a solution  $\hat{\underline{\mathbf{c}}}$  of the previous learning problem must verify  $\underline{K}(\mathbf{y} - \frac{1}{\lambda n} \underline{K} \hat{\underline{\mathbf{c}}} - \hat{\underline{\mathbf{c}}}) \in \partial(\epsilon \|\cdot\|_{\ell_1/\ell_2})(\hat{\underline{\mathbf{c}}})$ .

To compare [Lim et al.](#)'s approach to the framework proposed in this paper, let us dismiss the intercept  $\mathbf{b}$ , which means omitting the dual constraint  $\sum_{i=1}^n \boldsymbol{\alpha}_i = 0$  in (P4), and let  $\hat{\underline{\boldsymbol{\alpha}}} = (\hat{\boldsymbol{\alpha}}_1^\top, \dots, \hat{\boldsymbol{\alpha}}_n^\top)^\top$  be a solution to (P4). Then Fermat's rule gives:  $(\mathbf{y} - \frac{1}{\lambda n} \underline{K} \hat{\underline{\boldsymbol{\alpha}}} - \hat{\underline{\boldsymbol{\alpha}}}) \in \partial(\epsilon \|\cdot\|_{\ell_1/\ell_2})(\hat{\underline{\boldsymbol{\alpha}}})$ . It appears that both approaches induce data sparsity thanks to a structured norm and are equivalent up to a normalization by  $\underline{K}$ .

## 4.3. Quantile regression

A slight generalization of least absolute deviation consists in changing the slope of the absolute loss  $\ell(\cdot) = |\cdot|$ , asymmetrically around 0. For a quantile level  $\tau \in (0, 1)$ , the so called *pinball loss* is defined by:

$$\forall \xi \in \mathbb{R}, \quad \rho_\tau(\xi) = (\tau - \mathbf{I}_{\xi < 0}) \xi = \max(\tau \xi, (\tau - 1)\xi).$$

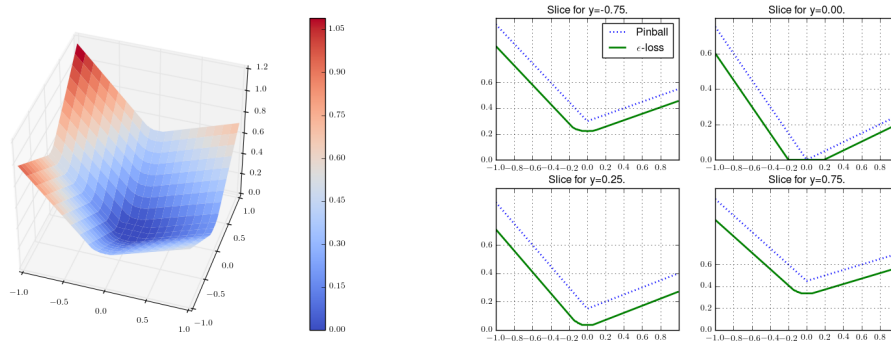


Figure 2: Multiple  $\epsilon$ -pinball loss and slices of it ( $\epsilon = 0.2$ ,  $\tau = (0.25, 0.6)$ ).

Such a loss is used to estimate a conditional quantile of the output random variable, instead of a conditional median such as in standard regression (Koenker, 2005).

In order to improve prediction, some works propose to estimate and to predict simultaneously several quantiles, which is referred to as *joint QR* (Takeuchi et al., 2013; Sangnier et al., 2016). In this context, close to multivariate regression, the output space is  $\mathbb{R}^p$  and the loss is  $\ell(\boldsymbol{\xi}) = \sum_{j=1}^p \rho_{\tau_j}(\xi_j)$ , where  $\boldsymbol{\tau}$  is a vector of  $p$  quantile levels  $\tau_j \in (0, 1)$ . Even though we do not have a closed-form expression for the corresponding  $\epsilon$ -loss, Figure 2 provides graphical representations of  $\ell_\epsilon$ .

Besides the loss described above, the multiple QR framework comes with  $\mathbf{y}_i \in \mathbb{R}^p$ ,  $\boldsymbol{\alpha}_i \in \mathbb{R}^p$ ,  $\mathbf{b} \in \mathbb{R}^p$  and  $h: \mathcal{X} \rightarrow \mathbb{R}^p$ , such that the  $j^{\text{th}}$  component of the prediction value,  $f_j(\mathbf{x}) = h_j(\mathbf{x}) + b_j$ , estimates the  $\tau_j$ -conditional quantile of the output random variable.

In a manner very similar to least absolute deviation, the Fenchel-Legendre transformation of the loss of interest is  $\ell^*(\boldsymbol{\alpha}) = 0$  when  $\boldsymbol{\tau} - \mathbf{1} \preceq \boldsymbol{\alpha} \preceq \boldsymbol{\tau}$  and  $+\infty$  otherwise. Therefore, Problem (P3) becomes:

$$\begin{aligned} & \underset{\forall i \in [n], \boldsymbol{\alpha}_i \in \mathbb{R}^p}{\text{minimize}} && \frac{1}{2\lambda n} \sum_{i,j=1}^n \langle \boldsymbol{\alpha}_i, K(\mathbf{x}_i, \mathbf{x}_j) \boldsymbol{\alpha}_j \rangle_{\ell_2} - \sum_{i=1}^n \langle \boldsymbol{\alpha}_i, \mathbf{y}_i \rangle_{\ell_2} + \epsilon \sum_{i=1}^n \|\boldsymbol{\alpha}_i\|_{\ell_2} \\ & \text{s. t.} && \forall i \in [n], \boldsymbol{\tau} - \mathbf{1} \preceq \boldsymbol{\alpha}_i \preceq \boldsymbol{\tau}, \quad \sum_{i=1}^n \boldsymbol{\alpha}_i = 0, \end{aligned} \tag{P5}$$

where the sparsity inducing term is the  $\ell_1/\ell_2$ -norm  $\sum_{i=1}^n \|\boldsymbol{\alpha}_i\|_{\ell_2}$  (Bach et al., 2012).

Very recently, Park and Kim (2011) also proposed an  $\epsilon$ -insensitive loss for *single QR* (that is, in the unidimensional case, when estimating only a single  $\tau$ -conditional quantile), based on a generalization of the one used for SVR. The proposition made by Park and Kim (2011) is:

$$\ell_\epsilon^{\text{PK}}(\xi) = \max(0, \rho_\tau(\xi) - \epsilon) = \begin{cases} 0 & \text{if } \rho_\tau(\xi) \leq \epsilon \\ \rho_\tau(\xi) - \epsilon & \text{otherwise.} \end{cases}$$

That is different from the  $\epsilon$ -loss proposed in this paper, as illustrated in Figure 1. As a comparison, the single QR dual in (Park and Kim, 2011) is similar to (P5) but with  $\max\left(\frac{\alpha_i}{\tau}, \frac{\alpha_i}{\tau-1}\right)$  instead of  $|\alpha_i|$ . Let us remark that, because of the asymmetry, it is difficult to get a sparsity intuition for this modified  $\ell_1$ -penalization.

The advantage of our approach over (Park and Kim, 2011) is to provide a unifying view of  $\epsilon$ -losses in the multidimensional setting, which matches SVR as a special case. Moreover, the possibility to extend Park and Kim’s loss to joint QR is unclear and could not guarantee data sparsity: as explained in Remark 1 (Supplementary A), this goal is achieved thanks to the  $\ell_1/\ell_2$ -norm in the dual, which is equivalent to our proposition of  $\epsilon$ -loss.

#### 4.4. Expectile regression

As for QR, it is worth considering weighting the squared loss asymmetrically around 0. The resulting framework is called *expectile regression* (Newey and Powell, 1987), by analogy with quantiles and conditional expectation estimation with the squared loss. Despite the fact that the expectile regression estimator has no statistical interpretation, it can be seen as a smooth (and less robust) version of a conditional quantile. This explains that this topic, at least in the unidimensional context, attracts a growing interest in the learning community (Farooq and Steinwart, 2015; Yang et al., 2015; Farooq and Steinwart, 2017).

Following the guiding principle of QR (Sangnier et al., 2016), we can imagine learning simultaneously several expectiles in a *joint expectile regression* framework. As far as we know, this multivariate setting did not appear in the literature yet. Therefore, we consider the multiple squared pinball loss:

$$\forall \xi \in \mathbb{R}, \quad \ell(\boldsymbol{\xi}) = \sum_{j=1}^p \psi_{\tau_j}(\xi_j), \quad \text{with } \psi_{\tau}(\xi) = \frac{1}{2} |\tau - \mathbf{I}_{\xi < 0}| \xi^2.$$

Similarly to the squared loss, the Fenchel-Legendre transformation of  $\ell$  is  $\ell^*: \boldsymbol{\alpha} \in \mathbb{R}^p \mapsto \frac{1}{2} \boldsymbol{\alpha}^\top \Delta(\boldsymbol{\alpha}) \boldsymbol{\alpha}$ , where  $\Delta(\boldsymbol{\alpha})$  is the  $p \times p$  diagonal matrix whose entries are  $\left( \left| \tau_j - \mathbf{I}_{(\boldsymbol{\alpha}_i)_j < 0} \right|^{-1} \right)_{1 \leq j \leq p}$ . Therefore, a dual optimization problem of sparse joint expectile regression is:

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

We see that the only difference compared to multivariate regression is the anisotropic ridge  $\Delta(\boldsymbol{\alpha}_i)$ . Moreover, compared to other kernelized approaches for expectiles estimation such as (Farooq and Steinwart, 2015; Yang et al., 2015), this setting goes beyond both in its multivariate feature as well as its data sparsity.

## 5. Statistical guarantees

In this section, we focus on the excess of risk for learning with the proposed  $\epsilon$ -insensitive loss. As it is usual in statistical learning, we consider the constrained version of Problem (P2), which is legitimated by convexity of  $\ell_\epsilon$  (Tikhonov and Arsenin, 1977):

$$\underset{\substack{h \in \mathcal{H} : \|h\|_{\mathcal{H}} \leq c \\ \mathbf{b} \in \mathbb{R}^p}}{\text{minimize}} \quad \frac{1}{n} \sum_{i=1}^n \ell_\epsilon(\mathbf{y}_i - (h(\mathbf{x}_i) + \mathbf{b})), \quad (\text{P6})$$



where  $c$  is a positive parameter linked to  $\lambda$ . This change of optimization problem is motivated by the fact that any solution to (P2) is also solution to (P6) for a well chosen parameter  $c$ .

Let  $(X, Y)$  be the couple of random variables of interest, following an unknown joint distribution. We assume being provided with  $n$  independent and identically distributed (*iid*) copies of  $(X, Y)$ , denoted  $(X_i, Y_i)_{1 \leq i \leq n}$ . Let us now fix an intercept  $\mathbf{b} \in \mathbb{R}^p$ .  $\mathcal{F} = \{h(\cdot) + \mathbf{b} : h \in \mathcal{H}, \|h\|_{\mathcal{H}} \leq c\}$  is a class of hypotheses. We denote  $\mathcal{R}_n(\mathcal{F})$  the Rademacher average of the class  $\mathcal{F}$  (Bartlett and Mendelson, 2002). Let  $f^\dagger \in \arg \min_{f \in \mathcal{F}} \mathbb{E}[\ell(Y - f(X))]$  be the target function, where the expectation is computed jointly on  $X$  and  $Y$ . For instance, in the context of single QR, if we assume that the conditional quantile lives in  $\mathcal{F}$ , then the target  $f^\dagger$  can be this conditional quantile (up to uniqueness). Moreover, let  $\hat{f}_\epsilon \in \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \ell_\epsilon(Y_i - f(X_i))$  be the empirical estimator of  $f^\dagger$ .

**Theorem 2 (Generalization)** *Let us assume that the loss  $\ell$  is  $L$ -Lipschitz ( $L > 0$ ) and that residues are bounded:  $\forall f \in \mathcal{F}, \|Y - f(X)\|_{\ell_2} \leq M$  almost surely (where  $M > 0$ ). Then, for any  $\delta \in (0, 1]$ , with probability at least  $1 - \delta$  over the random sample  $(X_i, Y_i)_{1 \leq i \leq n}$ :*

$$\mathbb{E} \left[ \ell \left( Y - \hat{f}_\epsilon(X) \right) \right] - \mathbb{E} \left[ \ell \left( Y - f^\dagger(X) \right) \right] \leq 2\sqrt{2}L\mathcal{R}_n(\mathcal{F}) + 2LM\sqrt{\frac{\log(2/\delta)}{2n}} + L\epsilon.$$

A proof of this theorem is given in Supplementary B. The bound is similar to what we usually observe (Boucheron et al., 2005; Mohri et al., 2012), but suffers from an extra term, which is linear in  $\epsilon$ . The latter embodies the bias induced by soft-thresholding the original loss  $\ell$  to induce data sparsity.

## 6. A primal-dual training algorithm

### 6.1. Description

For traditional kernel methods (real or vector-valued ones), pairwise coordinate descent as well as deterministic and randomized coordinate descents are popular and efficient training algorithms (Platt, 1999; Shalev-Shwartz and Zhang, 2013; Minh et al., 2016). However, few of current algorithms are able to handle multiple non-differentiable contributions in the objective value (such as the ones introduced by  $\ell_\epsilon^*$ ) and multiple linear constraints (coming from considering an intercept in our regressor, which is mandatory for QR (Takeuchi et al., 2006; Sangnier et al., 2016)). For these reasons, we propose to use a randomized primal-dual coordinate descent (PDCD) technique, introduced by Fercoq and Bianchi (2015) and utterly workable for the problem at hand. Moreover, PDCD has been proved favorably competitive with several state-of-the-art approaches (Fercoq and Bianchi, 2015).

The learning problem we are interested in is:

$$\underset{\forall i \in [n], \alpha_i \in \mathbb{R}^p}{\text{minimize}} \underbrace{s(\alpha_1, \dots, \alpha_n) + \sum_{i=1}^n \ell^*(\alpha_i)}_{\substack{\text{differentiable} \\ \text{not differentiable}}} + \underbrace{\epsilon \sum_{i=1}^n \|\alpha_i\|_{\ell_2}}_{\text{not differentiable}} + \underbrace{\chi_{\sum_{i=1}^n \alpha_i = 0}}_{\text{not differentiable}}, \quad (\text{P7})$$

where  $s(\boldsymbol{\alpha}_1, \dots, \boldsymbol{\alpha}_n) = \frac{1}{2\lambda n} \sum_{i,j=1}^n \langle \boldsymbol{\alpha}_i, K(\mathbf{x}_i, \mathbf{x}_j) \boldsymbol{\alpha}_j \rangle_{\ell_2} - \frac{1}{n} \sum_{i=1}^n \langle \boldsymbol{\alpha}_i, \mathbf{y}_i \rangle_{\ell_2}$  is a quadratic (differentiable) function. Depending on the kind of regression (see in Table 1), the mapping  $\ell^*$  may be differentiable or not. In any case, the objective function in (P7) can be written as the summation of three components: one differentiable and two non-differentiable. Given this three-term decomposition, PDCD (see Algorithm 1) dualizes the second non-differentiable component and deploys a randomized block coordinate descent, in which each iteration involves the proximal operators of the first non-differential function and of the Fenchel-Legendre transformation of the second one. We can see in Algorithm 1 that PDCD uses dual variables  $\boldsymbol{\theta} \in (\mathbb{R}^p)^n$  (which are updated during the descent) and has two sets of parameters  $\boldsymbol{\nu} \in \mathbb{R}^n$  and  $\boldsymbol{\mu} \in \mathbb{R}^n$ , that verify  $\forall i \in [n]: \mu_i < \frac{1}{\lambda_{\max,i} + \nu_i}$ , where  $\lambda_{\max,i}$  is the largest eigenvalue of  $K(\mathbf{x}_i, \mathbf{x}_i)$ . In practice, we keep the same parameters as in (Fercoq and Bianchi, 2015):  $\nu_i = 10\lambda_{\max,i}$  and  $\mu_i$  equal to 0.95 times the bound.

---

**Algorithm 1** Primal-Dual Coordinate Descent.
 

---

Initialize  $\boldsymbol{\alpha}_i, \boldsymbol{\theta}_i \in \mathbb{R}^p$  ( $\forall i \in [n]$ ).

**repeat**

    Choose  $i \in [n]$  uniformly at random.

    Set  $\bar{\boldsymbol{\theta}}^l \leftarrow \text{proj}_{\mathbb{R}^p}(\boldsymbol{\theta}^l + \text{diag}(\boldsymbol{\nu})\boldsymbol{\alpha}^l)$  for all  $l \in [p]$ .

    Set  $\mathbf{d}_i \leftarrow \nabla_{\boldsymbol{\alpha}_i} s(\boldsymbol{\alpha}_1, \dots, \boldsymbol{\alpha}_n) + 2\bar{\boldsymbol{\theta}}_i - \boldsymbol{\theta}_i$ .

    Set  $\bar{\boldsymbol{\alpha}}_i \leftarrow \text{prox}_{\mu_i(\epsilon\|\cdot\|_{\ell_2} + \ell^*)}(\boldsymbol{\alpha}_i - \mu_i \mathbf{d}_i)$ .

    Update coordinate  $i$ :  $\boldsymbol{\alpha}_i \leftarrow \bar{\boldsymbol{\alpha}}_i, \boldsymbol{\theta}_i \leftarrow \bar{\boldsymbol{\theta}}_i$ ,  
 and keep other coordinates unchanged.

**until** convergence

---

PDCD is a workable algorithm for non-differentiable  $\ell^*$  as soon as the proximal operator  $\text{prox}_{\epsilon\|\cdot\|_{\ell_2} + \ell^*}$  can be computed (when  $\ell^*$  is differentiable, it can be moved to  $s$  and there is only  $\text{prox}_{\epsilon\|\cdot\|_{\ell_2}}$ , which is easy to compute (Bach et al., 2012)). For the situations of interest in this paper, non-differentiable  $\ell^*$  appear for multiple QR and take the form of the characteristic function of box constraints (see Table 1). Nevertheless, as far as we know, computing the proximal operator of the sum of the  $\ell_2$ -norm and box constraints has not been done yet. Therefore, the following section provides a way to compute it for box constraints intersecting both the negative and the positive orthants (proofs are in Supplementary B.4).

## 6.2. Proximal operator for multiple quantile regression

Let  $\mathbf{a}$  and  $\mathbf{b}$  be two vectors from  $\mathbb{R}^p$  with positive entries, and  $\lambda > 0$ . From now on, let us denote  $[\cdot]_{-\mathbf{a}}^{\mathbf{b}} : \mathbb{R}^p \rightarrow \mathbb{R}^p$  the *clip operator*, defined by:

$$\forall \mathbf{y} \in \mathbb{R}^p, \forall j \in [p]: \quad \left([\mathbf{y}]_{-\mathbf{a}}^{\mathbf{b}}\right)_j = \begin{cases} y_j & \text{if } -a_j < y_j < b_j \\ b_j & \text{if } y_j \geq b_j \\ -a_j & \text{if } y_j \leq -a_j. \end{cases}$$

**Lemma 3** Let  $\mathbf{y} \in \mathbb{R}^p$  such that  $\|\mathbf{y}\|_{\ell_2} \geq \lambda$  be a fixed vector and  $\mu \in [0, 1]$ . The equation

$$\left(1 + \frac{\lambda}{\|[\mu\mathbf{y}]_{-\mathbf{a}}^{\mathbf{b}}\|_{\ell_2}}\right) \mu = 1 \quad (1)$$

is well defined and has a solution  $\mu \in [0, 1]$ .

**Proposition 4**

$$\forall \mathbf{y} \in \mathbb{R}^p, \quad \text{prox}_{\lambda\|\cdot\|_{\ell_2} + \chi_{-\mathbf{a} \preceq \cdot \preceq \mathbf{b}}}(\mathbf{y}) = [\mu\mathbf{y}]_{-\mathbf{a}}^{\mathbf{b}},$$

where  $\mu = 0$  if  $\|\mathbf{y}\|_{\ell_2} \leq \lambda$  and solution to Equation (1) otherwise.

Proposition 4 states that  $\text{prox}_{\lambda\|\cdot\|_{\ell_2} + \chi_{-\mathbf{a} \preceq \cdot \preceq \mathbf{b}}}$  can be computed as the composition of a clip and a scaling operator, for which the scaling factor can be easily obtained by a bisection of a Newton-Raphson method. As a straightforward consequence of Proposition 4, we can express the proximal operator of the sum of an  $\ell_1/\ell_2$ -norm and box constraints as the concatenation of  $\text{prox}_{\lambda\|\cdot\|_{\ell_2} + \chi_{-\mathbf{a} \preceq \cdot \preceq \mathbf{b}}}$  (see Supplementary B for a formal statement), which involves as many scaling factors as groups ( $n$  in this case).

### 6.3. Active set for multiple quantile regression

Learning sparse models gains a lot in identifying early the *active points* and optimizing only over them. This is a way to speed up learning, that is particularly topical (Ndiaye et al., 2015; Shibagaki et al., 2016; Ndiaye et al., 2016). *Active points* may have different meanings depending on the context (sparse regression, SVR, SVM) but the unifying concept is to detect optimizing variables for which the optimal value cannot be figured out beforehand.

In the context of QR (but this also holds true for SVR and SVM (Shibagaki et al., 2016)), *active points* correspond to every data point  $\mathbf{x}_i$  for which the dual vector  $\boldsymbol{\alpha}_i$  is neither null nor on the border of the box constraint. These active points can be identified thanks to optimality conditions.

Let  $\hat{f} = \hat{h} + \hat{\mathbf{b}}$  be an optimal solution to (P2) and  $(\hat{\boldsymbol{\alpha}}_i)_{1 \leq i \leq n}$  be optimal dual vectors to (P3). In the general setting, primal feasibility and stationarity from KKT conditions yield:

$$\forall i \in [n]: \quad \mathbf{y}_i - \hat{f}(\mathbf{x}_i) \in \partial(\ell_\epsilon^*)(\hat{\boldsymbol{\alpha}}_i), \quad (2)$$

where, in the case of QR,  $\ell_\epsilon^* = \epsilon \|\cdot\|_{\ell_2} + \chi_{\tau-1 \preceq \cdot \preceq \tau}$ . Therefore,  $\forall \boldsymbol{\alpha} \in \mathbb{R}^p$ :

$$\partial(\ell_\epsilon^*)(\boldsymbol{\alpha}) = \begin{cases} \{\mathbf{d} : \mathbf{d} \in \mathbb{R}^p, \|\mathbf{d}\|_{\ell_2} \leq \epsilon\} & \text{if } \boldsymbol{\alpha} = 0 \\ \frac{\epsilon}{\|\boldsymbol{\alpha}\|_{\ell_2}} \boldsymbol{\alpha} & \text{if } \boldsymbol{\alpha} \neq 0, \tau - 1 \prec \boldsymbol{\alpha} \prec \tau \\ \left\{ \mathbf{d} \in \partial(\chi_{\tau-1 \preceq \cdot \preceq \tau})(\boldsymbol{\alpha}) : \left\langle \mathbf{d}, \frac{\boldsymbol{\alpha}}{\|\boldsymbol{\alpha}\|_{\ell_2}} \right\rangle_{\ell_2} \geq \epsilon \right\} & \text{otherwise.} \end{cases}$$

Consequently, optimality condition in Equation (2) indicates that:

$$1. \quad \left\| \mathbf{y}_i - \hat{f}(\mathbf{x}_i) \right\|_{\ell_2} < \epsilon \implies \hat{\boldsymbol{\alpha}}_i = 0;$$

$$2. \left\| \mathbf{y}_i - \hat{f}(\mathbf{x}_i) \right\|_{\ell_2} > \epsilon \implies \forall j \in [p], (\hat{\boldsymbol{\alpha}}_i)_j = \tau_j \text{ or } (\hat{\boldsymbol{\alpha}}_i)_j = \tau_j - 1.$$

Therefore, for each situation (1. or 2.), if both conditions are fulfilled for the current estimates  $f$  and  $\boldsymbol{\alpha}_i$ , the corresponding dual vector  $\boldsymbol{\alpha}_i$  is put aside (at least temporally) and is not updated until optimality conditions are violated. In section 7.1, we will show that this strategy dramatically speed up the learning process when  $\epsilon$  is large enough.

## 7. Numerical experiments

This section presents two numerical experiments in the context of multiple QR with data sparsity. The first experiment deals with training time. We compare an implementation of Algorithm 1 with an off-the-shelf solver and study the impact of the active set strategy. The second experiment analyses the effect of  $\epsilon$  on quantile prediction and on the number of support vectors (data points  $\mathbf{x}_i$  for which  $\boldsymbol{\alpha}_i \neq 0$ ).

Following Sangnier et al. (2016), we use a matrix-valued kernel of the form  $K(\mathbf{x}, \mathbf{x}') = k(\mathbf{x}, \mathbf{x}')\mathbf{B}$ , where  $\mathbf{B} = (\exp(-\gamma(\tau_i - \tau_j)^2))_{1 \leq i, j \leq p}$  and  $k(\mathbf{x}, \mathbf{x}') = \exp(-\|\mathbf{x} - \mathbf{x}'\|_{\ell_2}^2 / 2\sigma^2)$ . In the first experiment (simulated data),  $\gamma = 1$  and  $\sigma^2 = 1$ , while in the second one (real data)  $\gamma$  is set to 0.1 and  $\sigma$  is the 0.7-quantile of the pairwise distances of the training data  $\{\mathbf{x}_i\}_{1 \leq i \leq n}$ . Quantile levels of interest are  $\boldsymbol{\tau} = (0.25, 0.5, 0.75)$ .

### 7.1. Training time

We aim at comparing the execution time of four approaches for solving (P5): i) Algorithm 1 with active set (*PDCD with AS*); ii) Algorithm 1 without active set (*PDCD without AS*); iii) an off-the-shelf solver based on an interior-point method for quadratic cone programming (*CVXOPT (CONEQP)*) and iv) quadratic programming (*CVXOPT (QP)*) (Anderson et al., 2012). For the last approach, we leverage a variational formulation of the sparsity constraint (Bach et al., 2012):  $\sum_{i=1}^n \|\boldsymbol{\alpha}_i\|_{\ell_2} = \inf_{\forall i \in [n], \mu_i \in \mathbb{R}_+} \frac{1}{2} \sum_{i=1}^n \left( \frac{1}{\mu_i} \|\boldsymbol{\alpha}_i\|_{\ell_2}^2 + \mu_i \right)$ , and alternate minimization with respect to  $(\boldsymbol{\alpha}_i)_{1 \leq i \leq n}$  and  $(\mu_i)_{1 \leq i \leq n}$ . By convexity and differentiability of the objective, alternate minimization converges to an optimal solution (Rockafellar, 1970).

We fix  $1/(\lambda n) = 10^2$  and we use a synthetic dataset for which  $X \in [0, 1.5]$ . The target  $Y$  is computed as a sine curve at 1 Hz modulated by a sine envelope at 1/3 Hz and mean 1. Moreover, this pattern is distorted with a random Gaussian noise with mean 0 and a linearly decreasing standard deviation from 1.2 at  $X = 0$  to 0.2 at  $X = 1.5$ .

Figure 3 depicts the dual gap achieved by the competitors with respect to the CPU time in different configurations ( $\epsilon \in \{0.1, 1, 5\}$  and sample size  $n \in \{50, 500\}$ ). These curves are obtained by increasing the maximal number of iterations of each solver and recording primal and dual objective values attained along with the CPU time. Given a configuration ( $\epsilon$  and  $n$  fixed), all dual gaps are computed with respect to the lowest primal obtained. This comparison procedure is motivated by the fact that for multiple QR, we cannot compute the exact duality gap (we cannot compute the primal loss  $\ell_\epsilon$ ), but only over-estimate it (see Remark 3 in Supplementary A). This explains why plots in Figure 3 do not necessarily converge to 0 (we conjecture that when a method faces no progress, the exact dual gap is almost null). Overall, the curves portray the behavior of the dual objective function rather than the exact duality gap.

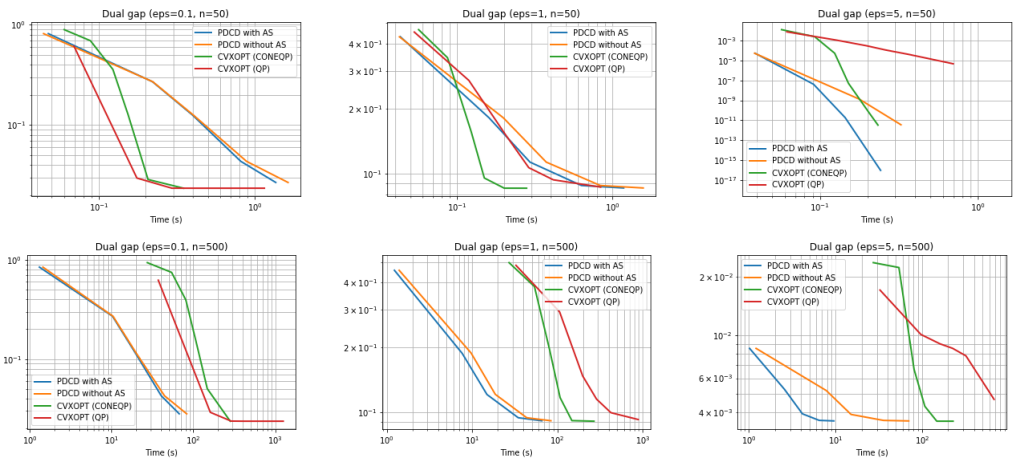


Figure 3: Training time (synthetic dataset).

For small-scale data and small values of  $\epsilon$ , CVXOPT is more efficient than our approach, but its supremacy caves in when increasing  $\epsilon$  and/or the sample size  $n$ . The impact of the problem size  $n$  on the comparison between interior-point methods (such as CVXOPT) and first-order ones (such as PDCD) is known for a long time, so this first finding is comforting. Also, increasing the sparsity regularizer  $\epsilon$  requires more and more calls to the alternate procedure CVXOPT (QP), so the behavior with respect to  $\epsilon$  was expected.

Analyzing the active-set strategy, we remark that it is never a drawback but becomes a real asset for hard sparsity constraints. This comes from many  $\hat{\alpha}_i$  to be null and thus put aside. Yet, it appears that, contrarily to the unidimensional case, only a few  $\alpha_i$  are on the border of the box constraint. Therefore, the active-set strategy provides only a limited (but real) gain compared to the vanilla version of PDCD for small values of  $\epsilon$ .

## 7.2. Data sparsity

This section investigates the actual data sparsity introduced by the proposed  $\epsilon$ -insensitive loss, along with its impact on the generalization error. For this purpose, we consider the datasets used in (Takeuchi et al., 2006; Zhang et al., 2016), coming from the UCI repository and three R packages: quantreg, alr3 and MASS. The sample sizes  $n$  vary from 38 (CobarOre) to 1375 (heights) and the numbers of explanatory variables vary from 1 (6 sets) to 12 (BostonHousing). The datasets are standardized coordinate-wise to have zero mean and unit variance. In addition, the generalization error is estimated by the mean over 10 trials of the empirical loss  $\frac{1}{m} \sum_{i=1}^m \ell(\mathbf{y}_i^{\text{test}} - f(\mathbf{x}_i^{\text{test}}))$  computed on a test set. For each trial, the whole dataset is randomly split in a train and a test set with ratio 0.7-0.3. The parameter  $1/(\lambda n)$  is chosen by a 3-fold cross-validation (minimizing the pinball loss) on a logarithmic scale between  $10^{-3}$  and  $10^3$  (10 values). At last, a training point  $\mathbf{x}_i$  is considered a support vector if  $\|\alpha_i\|_{\ell_2} / p > 10^{-3}$ .

Table 2 reports the average empirical test loss (scaled by 100) along with the percentage of support vectors (standard deviations appear in Supplementary C). For each dataset, the bold-face numbers are the two lowest losses, to be compared to the loss for  $\epsilon = 0$ . We

Table 2: Empirical (test) pinball loss  $\times 100$  (percentage of support vectors in parentheses).

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 (100)	<b>67.40</b> (99)	<b>67.17</b> (95)	67.54 (87)	69.93 (65)	73.24 (49)	76.80 (39)	83.42 (24)	100.22 (17)	142.19 (6)
ftcollinssnow	<b>109.07</b> (100)	109.12 (100)	109.14 (100)	109.15 (100)	109.11 (100)	110.39 (98)	<b>109.05</b> (82)	110.90 (53)	109.81 (34)	113.50 (8)
highway	79.29 (100)	78.10 (97)	76.75 (97)	76.66 (91)	<b>75.09</b> (64)	<b>70.94</b> (51)	75.10 (37)	97.67 (31)	112.30 (0)	112.09 (1)
heights	91.05 (100)	91.00 (100)	90.98 (100)	<b>90.98</b> (100)	91.18 (99)	91.21 (86)	<b>90.98</b> (67)	91.09 (40)	91.51 (23)	93.34 (5)
sniffer	32.34 (100)	<b>31.40</b> (99)	32.31 (97)	<b>31.40</b> (84)	34.64 (43)	39.84 (24)	41.82 (12)	52.06 (6)	62.21 (4)	103.76 (3)
snowgeese	<b>49.62</b> (100)	<b>50.51</b> (98)	51.25 (87)	51.08 (70)	52.88 (37)	53.81 (23)	62.81 (17)	90.15 (15)	107.53 (14)	94.25 (1)
ufc	57.87 (100)	57.90 (100)	<b>57.78</b> (100)	57.84 (100)	<b>57.67</b> (78)	57.84 (52)	58.19 (35)	61.04 (15)	66.81 (5)	86.23 (2)
birthwt	99.93 (100)	99.95 (100)	99.93 (100)	99.70 (100)	<b>99.25</b> (100)	100.50 (87)	99.80 (67)	<b>98.71</b> (43)	99.56 (24)	103.39 (6)
crabs	8.59 (100)	<b>8.52</b> (91)	<b>8.49</b> (55)	9.44 (16)	19.94 (6)	23.08 (2)	31.44 (3)	44.08 (2)	53.45 (1)	86.91 (2)
GAGurine	44.30 (100)	<b>44.26</b> (99)	<b>44.25</b> (99)	44.86 (87)	46.20 (50)	49.87 (33)	52.88 (22)	57.06 (12)	65.89 (6)	103.32 (2)
geyser	<b>77.81</b> (100)	78.15 (100)	<b>78.12</b> (100)	78.45 (100)	78.40 (92)	78.28 (78)	78.54 (59)	80.55 (32)	85.15 (16)	99.92 (1)
gilgais	<b>32.96</b> (100)	<b>33.12</b> (99)	33.27 (96)	33.42 (81)	35.08 (43)	36.62 (25)	37.94 (14)	48.17 (5)	94.65 (7)	104.12 (0)
topo	47.49 (100)	48.93 (100)	48.74 (98)	48.17 (94)	<b>41.65</b> (57)	<b>45.24</b> (38)	51.19 (26)	53.68 (16)	58.21 (12)	80.57 (6)
BostonHousing	<b>34.54</b> (100)	34.68 (99)	34.70 (97)	<b>34.09</b> (80)	35.27 (35)	37.65 (20)	41.31 (13)	55.04 (7)	73.39 (7)	112.22 (12)
CobarOre	<b>0.50</b> (100)	<b>0.05</b> (38)	8.75 (36)	12.47 (26)	23.84 (18)	35.82 (17)	47.35 (15)	66.15 (14)	84.51 (12)	106.89 (6)
engel	43.57 (100)	43.50 (100)	<b>43.47</b> (99)	<b>43.44</b> (89)	57.36 (39)	43.98 (19)	46.31 (11)	53.15 (4)	69.43 (5)	100.48 (0)
mcycle	<b>63.95</b> (100)	<b>63.88</b> (99)	64.26 (99)	64.90 (98)	65.89 (88)	67.29 (70)	70.11 (51)	74.78 (26)	86.49 (14)	109.79 (2)
BigMac2003	<b>49.94</b> (100)	<b>49.97</b> (98)	50.00 (96)	50.27 (85)	51.16 (56)	51.44 (36)	53.63 (28)	77.40 (18)	106.38 (14)	136.76 (4)
UN3	71.27 (100)	<b>70.94</b> (100)	<b>71.03</b> (100)	71.49 (99)	71.37 (87)	71.53 (65)	72.68 (50)	76.72 (27)	84.50 (13)	109.59 (0)
cpus	<b>11.31</b> (100)	<b>13.32</b> (28)	15.57 (21)	20.16 (15)	25.88 (8)	35.66 (6)	55.27 (6)	65.05 (0)	65.05 (0)	65.02 (0)

observe first that increasing  $\epsilon$  does promote data sparsity. Second, the empirical loss tends to increase with  $\epsilon$ , as expected, but we do not lose much by requiring sparsity (we can even gain a little).

## 8. Conclusion

This paper introduces a novel and unifying framework concerning  $\epsilon$ -insensitive losses, which offers a variety of opportunities for traditional, quantile and expectile regression, for uni and multivariate outcomes. Estimators are based on vector-valued RKHSs (Micchelli and Pontil, 2005) and benefit from theoretical guarantees concerning generalization as well as an efficient learning algorithm. Future directions of research include improving further the learning procedure and deriving sharper generalization bounds.

## References

- M.S. Anderson, J. Dahl, and L. Vandenberghe. CVXOPT: A Python package for convex optimization, version 1.1.5., 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.
- 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.
- H. H. Bauschke and P. L. Combettes. *Convex Analysis and Monotone Operator Theory in Hilbert Spaces*. Springer, 2011.
- B.E. Boser, I.M. Guyon, and V.N. Vapnik. A Training Algorithm for Optimal Margin Classifiers. In *Conference on Learning Theory*, 1992.
- S. Boucheron, O. Bousquet, and G. Lugosi. Theory of Classification: a Survey of Some Recent Advances. *ESAIM: Probability and Statistics*, 9:323–375, 2005.

- C. Brouard, M. Szafranski, and F. d'Alché Buc. Input Output Kernel Regression: Supervised and Semi-Supervised Structured Output Prediction with Operator-Valued Kernels. *Journal of Machine Learning Research*, 17(176):1–48, 2016.
- H. Drucker, C.J.C. Burges, L. Kaufman, A.J. Smola, and V. Vapnik. Support Vector Regression Machines. In *Advances in Neural Information Processing Systems 9*, 1997.
- M. Farooq and I. Steinwart. An SVM-like Approach for Expectile Regression. *arXiv:1507.03887 [stat]*, 2015.
- M. Farooq and I. Steinwart. Learning Rates for Kernel-Based Expectile Regression. *arXiv:1702.07552 [cs, stat]*, 2017.
- O. Fercoq and P. Bianchi. A Coordinate Descent Primal-Dual Algorithm with Large Step Size and Possibly Non Separable Functions. *arXiv:1508.04625 [math]*, 2015.
- T. Hastie, R. Tibshirani, and J.H. Friedman. *The elements of statistical learning: data mining, inference, and prediction*. Springer, 2009.
- T. Hastie, R. Tibshirani, and M. Wainwright. *Statistical Learning with Sparsity: The Lasso and Generalizations*. Taylor & Francis, 2015.
- R. Koenker. *Quantile Regression*. Cambridge University Press, Cambridge, New York, 2005.
- N. Lim, F. d'Alché Buc, C. Auliac, and G. Michailidis. Operator-valued kernel-based vector autoregressive models for network inference. *Machine Learning*, 99(3):489–513, 2014.
- C.A. Micchelli and M.A. Pontil. On Learning Vector-Valued Functions. *Neural Computation*, 17:177–204, 2005.
- H.Q. Minh, L. Bazzani, and V. Murino. A Unifying Framework in Vector-valued Reproducing Kernel Hilbert Spaces for Manifold Regularization and Co-Regularized Multi-view Learning. *Journal of Machine Learning Research*, 17(25):1–72, 2016.
- M. Mohri, A. Rostamizadeh, and A. Talwalkar. *Foundations of Machine Learning*. The MIT Press, 2012.
- E. Ndiaye, O. Fercoq, A. Gramfort, and J. Salmon. GAP Safe screening rules for sparse multi-task and multi-class models. In *Advances in Neural Information Processing Systems 28*, 2015.
- E. Ndiaye, O. Fercoq, A. Gramfort, and J. Salmon. GAP Safe Screening Rules for Sparse-Group Lasso. In *Advances in Neural Information Processing Systems 29*, 2016.
- W.K. Newey and J.L. Powell. Asymmetric Least Squares Estimation and Testing. *Econometrica*, 55(4):819–847, 1987.
- J. Park and J. Kim. Quantile regression with an epsilon-insensitive loss in a reproducing kernel Hilbert space. *Statistics & Probability Letters*, 81(1):62–70, 2011.

- J.C. Platt. *Fast training of support vector machines using sequential minimal optimization*. Advances in Kernel Methods. MIT Press, Cambridge, MA, USA, 1999.
- R.T. Rockafellar. *Convex Analysis*. Princeton University Press, 1970.
- 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.
- B. Schölkopf, R. Herbrich, and A.J. Smola. A Generalized Representer Theorem. In *Computational Learning Theory*, pages 416–426. Springer, Berlin, Heidelberg, 2001.
- S. Shalev-Shwartz and T. Zhang. Stochastic Dual Coordinate Ascent Methods for Regularized Loss Minimization. *Journal of Machine Learning Research*, 14:567–599, 2013.
- A. Shibagaki, M. Karasuyama, K. Hatano, and I. Takeuchi. Simultaneous Safe Screening of Features and Samples in Doubly Sparse Modeling. In *Proceedings of the 33rd International Conference on International Conference on Machine Learning*, 2016.
- I. Steinwart and A. Christmann. *Support Vector Machines*. Springer, New York, NY, 2008.
- I. Takeuchi, Q.V. Le, T.D. Sears, and A.J. Smola. Nonparametric Quantile Estimation. *Journal of Machine Learning Research*, 7:1231–1264, 2006.
- I. Takeuchi, T. Hongo, M. Sugiyama, and S. Nakajima. Parametric Task Learning. In *Advances in Neural Information Processing Systems 26*, pages 1358–1366, 2013.
- R. Tibshirani. Regression Shrinkage and Selection via the Lasso. *Journal of the Royal Statistical Society. Series B (Methodological)*, 58(1):267–288, 1996.
- A.N. Tikhonov and V.Y. Arsenin. *Solutions of ill-posed problems*. Winston, Washington, DC, 1977.
- M.E. Tipping. Sparse Bayesian Learning and the Relevance Vector Machine. *Journal of Machine Learning Research*, 1:211–244, 2001.
- V. Vapnik. *The Nature of Statistical Learning Theory*. Information Science and Statistics. Springer, 2010.
- Y. Yang, T. Zhang, and H. Zou. Flexible Expectile Regression in Reproducing Kernel Hilbert Space. *arXiv:1508.05987 [stat]*, 2015.
- C. Zhang, Y. Liu, and Y. Wu. On Quantile Regression in Reproducing Kernel Hilbert Spaces with the Data Sparsity Constraint. *Journal of Machine Learning Research*, 17(1):1374–1418, 2016.