



**HAL**  
open science

# Sharp Oracle Inequalities for Low-complexity Priors

Duy Tung Luu, Jalal M. Fadili, Christophe Chesneau

► **To cite this version:**

Duy Tung Luu, Jalal M. Fadili, Christophe Chesneau. Sharp Oracle Inequalities for Low-complexity Priors. *Annals of the Institute of Statistical Mathematics*, In press. hal-01422476v5

**HAL Id: hal-01422476**

**<https://hal.science/hal-01422476v5>**

Submitted on 27 Sep 2018

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

---

# Sharp Oracle Inequalities for Low-complexity Priors

Tung Duy Luu · Jalal Fadili  
(corresponding author) · Christophe  
Chesneau

**Abstract** In this paper, we consider a high-dimensional statistical estimation problem in which the number of parameters is comparable or larger than the sample size. We present a unified analysis of the performance guarantees of exponential weighted aggregation and penalized estimators with a general class of data losses and priors which encourage objects which conform to some notion of simplicity/complexity. More precisely, we show that these two estimators satisfy sharp oracle inequalities for prediction ensuring their good theoretical performances. We also highlight the differences between them. When the noise is random, we provide oracle inequalities in probability using concentration inequalities. These results are then applied to several instances including the Lasso, the group Lasso, their analysis-type counterparts, the  $\ell_\infty$  and the nuclear norm penalties. All our estimators can be efficiently implemented using proximal splitting algorithms.

**Keywords** High-dimensional estimation · Exponential weighted aggregation · Penalized estimation · Oracle inequality · Low complexity models

## 1 Introduction

### 1.1 Problem statement

Our statistical context is the following. Let  $\mathbf{y} = (\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_n)$  be  $n$  observations with common marginal distribution, and  $\mathbf{X} \in \mathbb{R}^{n \times p}$  a deterministic design matrix. The goal is to estimate a parameter vector  $\boldsymbol{\theta} \in \mathbb{R}^p$  of the observations marginal distribution based on the data  $\mathbf{y}$  and  $\mathbf{X}$ .

Let  $F : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$  be a loss function supposed to be smooth and convex that assigns to each  $\boldsymbol{\theta} \in \mathbb{R}^p$  a cost  $F(\mathbf{X}\boldsymbol{\theta}, \mathbf{y})$ . Let  $\boldsymbol{\theta}_0 \in \text{Argmin}_{\boldsymbol{\theta} \in \mathbb{R}^p} \mathbb{E} [F(\mathbf{X}\boldsymbol{\theta}, \mathbf{y})]$  be any minimizer of the population risk. We regard  $\boldsymbol{\theta}_0$  as the true parameter. A usual instance of this statistical setting is the standard linear regression

---

Tung Duy Luu, Jalal Fadili  
Normandie Univ, ENSICAEN, CNRS, GREYC, 6, Bd du Maréchal Juin, 14050 Caen, France  
E-mail: {duy-tung.luu, Jalal.Fadili}@ensicaen.fr

Christophe Chesneau  
Normandie Univ, Université de Caen Normandie, CNRS, LMNO, BP 5186, 14032 Caen, France  
E-mail: christophe.chesneau@unicaen.fr

model based on  $n$  pairs  $(\mathbf{y}_i, \mathbf{X}_i)$  of response-covariate that are linked linearly  $\mathbf{y} = \mathbf{X}\boldsymbol{\theta}_0 + \boldsymbol{\xi}$ , and  $F(\mathbf{u}, \mathbf{y}) = \frac{1}{2}\|\mathbf{y} - \mathbf{u}\|_2^2$ .

Our goal is to provide general oracle inequalities in prediction for two estimators of  $\boldsymbol{\theta}_0$ : the penalized estimator and exponential weighted aggregation. In the setting where " $p$  larger than  $n$  (possibly much larger)", the estimation problem is ill-posed since the rectangular matrix  $\mathbf{X}$  has a kernel of dimension at least  $p - n$ . To circumvent this difficulty, we will exploit the prior that  $\boldsymbol{\theta}_0$  has some low-complexity structure (among which sparsity and low-rank are the most popular). That is, even if the ambient dimension  $p$  of  $\boldsymbol{\theta}_0$  is very large, its intrinsic dimension is much smaller than the sample size  $n$ . This makes it possible to build estimates  $\mathbf{X}\widehat{\boldsymbol{\theta}}$  with good provable performance guarantees under appropriate conditions. There has been a flurry of research on the use of low-complexity regularization in ill-posed recovery problems in various areas including statistics and machine learning.

## 1.2 Penalized Estimators

Regularization is now a central theme in many fields including statistics, machine learning and inverse problems. It allows one to impose on the set of candidate solutions some prior structure on the object to be estimated. This regularization ranges from squared Euclidean or Hilbertian norms to non-Hilbertian norms (e.g.  $\ell_1$  norm for sparse objects, or nuclear norm for low-rank matrices) that have sparked considerable interest in the recent years. In this paper, we consider the class of estimators obtained by solving the convex optimization problem<sup>1</sup>

$$\widehat{\boldsymbol{\theta}}_n^{\text{PEN}} \in \underset{\boldsymbol{\theta} \in \mathbb{R}^p}{\text{Argmin}} \{V_n(\boldsymbol{\theta}) \stackrel{\text{def}}{=} \frac{1}{n}F(\mathbf{X}\boldsymbol{\theta}, \mathbf{y}) + \lambda_n J(\boldsymbol{\theta})\}, \quad (1)$$

where the regularizing penalty  $J$  is a proper closed convex function that promotes some specific notion of simplicity/low-complexity, and  $\lambda_n > 0$  is the regularization parameter. A prominent member covered by (1) is the Lasso [Chen et al \(1999\)](#); [Tibshirani \(1996\)](#); [Osborne et al \(2000\)](#); [Donoho \(2006\)](#); [Candès and Plan \(2009\)](#); [Bickel et al \(2009\)](#); [Bühlmann and van de Geer \(2011\)](#); [Koltchinskii \(2008\)](#) and its variants such the analysis/fused Lasso [Rudin et al \(1992\)](#); [Tibshirani et al \(2005\)](#), SLOPE [Bogdan et al \(2014\)](#); [Su and Candès \(2015\)](#) or group Lasso [Bakin \(1999\)](#); [Yuan and Lin \(2006\)](#); [Bach \(2008\)](#); [Wei and Huang \(2010\)](#). Another example is the nuclear norm minimization for low rank matrix recovery motivated by various applications including robust PCA, phase retrieval, control and computer vision [Recht et al \(2010\)](#); [Candès and Recht \(2009\)](#); [Fazel et al \(2001\)](#); [Candès et al \(2013\)](#). See [Negahban et al \(2012\)](#); [Bühlmann and van de Geer \(2011\)](#); [van de Geer \(2014\)](#); [Vaier et al \(2015b\)](#) for generalizations and comprehensive reviews.

<sup>1</sup> To avoid trivialities, the set of minimizers is assumed non-empty, which holds for instance if  $J$  is also coercive.

### 1.3 Exponential Weighted Aggregation (EWA)

An alternative to the penalized estimator (1) is the aggregation by exponential weighting, which consists in substituting averaging for minimization. The aggregators are defined via the probability density function

$$\mu_n(\boldsymbol{\theta}) = \frac{\exp(-V_n(\boldsymbol{\theta})/\beta)}{\int_{\Theta} \exp(-V_n(\boldsymbol{\omega})/\beta) d\boldsymbol{\omega}}, \quad (2)$$

where  $\beta > 0$  is called temperature parameter. If all  $\boldsymbol{\theta}$  are candidates to estimate the true vector  $\boldsymbol{\theta}_0$ , then  $\Theta = \mathbb{R}^p$ . The aggregate is thus defined by

$$\widehat{\boldsymbol{\theta}}_n^{\text{EWA}} = \int_{\mathbb{R}^p} \boldsymbol{\theta} \mu_n(\boldsymbol{\theta}) d\boldsymbol{\theta}. \quad (3)$$

Aggregation by exponential weighting has been widely considered in the statistical and machine learning literatures, see e.g., Dalalyan and Tsybakov (2007, 2008, 2009, 2012); Nemirovski (2000); Yang (2004); Rigollet and Tsybakov (2007); Lecué (2007); Guedj and Alquier (2013); Duy Luu et al (2016) to name a few. The technique used in these papers were initiated by Leung and Barron (2006) (use of Stein's identity to study an early version of EWA) and Catoni (2003, 2007) (PAC-Bayesian theory).  $\widehat{\boldsymbol{\theta}}_n^{\text{EWA}}$  can also be interpreted as the posterior conditional mean in the Bayesian sense if  $F/(n\beta)$  is the negative-loglikelihood associated to the noise  $\boldsymbol{\xi}$  with the prior density  $\pi(\boldsymbol{\theta}) \propto \exp(-\lambda_n J(\boldsymbol{\theta})/\beta)$ .

### 1.4 Oracle inequalities

Oracle inequalities, which are at the heart of our work, quantify the quality of an estimator compared to the best possible one among a family of estimators. These inequalities are well adapted in the scenario where the prior penalty promotes some notion of low-complexity (e.g. sparsity, low rank, etc.). Given two vectors  $\boldsymbol{\theta}_1$  and  $\boldsymbol{\theta}_2$ , let  $R_n(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2)$  be a nonnegative error measure between their predictions, respectively  $\mathbf{X}\boldsymbol{\theta}_1$  and  $\mathbf{X}\boldsymbol{\theta}_2$ . A popular example is the averaged prediction squared error  $\frac{1}{n} \|\mathbf{X}\boldsymbol{\theta}_1 - \mathbf{X}\boldsymbol{\theta}_2\|_2^2$ , where  $\|\cdot\|_2$  is the  $\ell_2$  norm.  $R_n$  will serve as a measure of the performance of the estimators  $\widehat{\boldsymbol{\theta}}_n^{\text{EWA}}$  and  $\widehat{\boldsymbol{\theta}}_n^{\text{PEN}}$ . More precisely, we aim to prove that  $\widehat{\boldsymbol{\theta}}_n^{\text{EWA}}$  and  $\widehat{\boldsymbol{\theta}}_n^{\text{PEN}}$  mimic as much as possible the best possible model. This idea is materialized in the following type of inequalities (stated here for EWA)

$$R_n(\widehat{\boldsymbol{\theta}}_n^{\text{EWA}}, \boldsymbol{\theta}_0) \leq C \inf_{\boldsymbol{\theta} \in \mathbb{R}^p} (R_n(\boldsymbol{\theta}, \boldsymbol{\theta}_0) + \Delta_{n,p,\lambda_n,\beta}(\boldsymbol{\theta})), \quad (4)$$

where  $C \geq 1$  is the leading constant of the oracle inequality and the remainder term  $\Delta_{n,\lambda_n,\beta}(\boldsymbol{\theta})$  depends on the performance of the estimator, the complexity of  $\boldsymbol{\theta}$ , the sample size  $n$ , the dimension  $p$ , and the regularization and temperature parameters  $(\lambda_n, \beta)$ . An estimator with good oracle properties would

correspond to  $C$  close to 1 (ideally,  $C = 1$ , in which case the inequality is said “sharp”), and  $\Delta_{n,p,\lambda_n,\beta}(\boldsymbol{\theta})$  is small and decreases rapidly to 0 as  $n \rightarrow +\infty$ .

### 1.5 Contributions

We provide a unified analysis where we capture the essential ingredients behind the low-complexity priors promoted by  $J$ , relying on sophisticated arguments from convex analysis and our previous work [Fadili et al \(2013\)](#); [Vaïter et al \(2015a, 2018, 2017, 2015b\)](#). Our main contributions are summarized as follows:

- We show that the EWA estimator  $\hat{\boldsymbol{\theta}}_n^{\text{EWA}}$  in (2) and the penalized estimator  $\hat{\boldsymbol{\theta}}_n^{\text{PEN}}$  in (1) satisfy (deterministic) sharp oracle inequalities for prediction with optimal remainder term, for general data losses  $F$  beyond the usual quadratic one, and  $J$  is a proper finite-valued sublinear function (i.e.  $J$  is finite-valued convex and positively homogeneous). We also highlight the differences between the two estimators in terms of the corresponding bounds.
- When the observations are random, we prove oracle inequalities in probability. The theory is non-asymptotic in nature, as it yields explicit bounds that hold with high probability for finite sample sizes, and reveals the dependence on dimension and other structural parameters of the model.
- For the standard linear model with Gaussian or sub-Gaussian noise, and a quadratic loss, we deliver refined versions of these oracle inequalities in probability. We underscore the role of the Gaussian width, a concept that captures important geometric characteristics of sets in  $\mathbb{R}^n$ .
- These results yield naturally a large number of corollaries when specialized to penalties routinely used in the literature, among which the Lasso, the group Lasso, their analysis-type counterparts (fused (group) Lasso), the  $\ell_\infty$  and the nuclear norms. Some of these corollaries are known and others novel.

The estimators  $\hat{\boldsymbol{\theta}}_n^{\text{EWA}}$  and  $\hat{\boldsymbol{\theta}}_n^{\text{PEN}}$  can be easily implemented thanks to the framework of proximal splitting methods, and more precisely forward-backward type splitting. While the latter is well-known to solve (1) [Vaïter et al \(2015b\)](#), its application within a proximal Langevin Monte-Carlo algorithm to compute  $\hat{\boldsymbol{\theta}}_n^{\text{EWA}}$  with provable guarantees has been recently developed by the authors in [Duy Luu et al \(2016\)](#) to sample from log-semiconcave densities<sup>2</sup>, see also [Durmus et al \(2016\)](#) for log-concave densities.

### 1.6 Relation to previous work

Our oracle inequality for  $\hat{\boldsymbol{\theta}}_n^{\text{EWA}}$  extends the work of [Dalalyan et al \(2018\)](#) with an unprecedented level of generality, far beyond the Lasso and the nuclear

<sup>2</sup> In a forthcoming paper, this framework was extended to cover the even more general class of prox-regular functions.

norm. Our prediction sharp oracle inequality for  $\widehat{\boldsymbol{\theta}}_n^{\text{PEN}}$  specializes to that of Sun and Zhang (2012) in the case of the Lasso (see also the discussion in Dalalyan et al (2017) and references therein) and that of Koltchinskii et al (2011) for the case of the nuclear norm. Our work also goes much beyond that in van de Geer (2014) on weakly decomposable priors, where we show in particular that there is no need to impose decomposability on the regularizer, since it is rather an intrinsic property of it.

## 1.7 Paper organization

Section 2 states our main assumptions on the data loss and the prior penalty. All the concepts and notions are exemplified on some penalties some of which are popular in the literature. In Section 3, we prove our main oracle inequalities, and their versions in probability. We then tackle the case of linear regression with quadratic data loss in Section 4. Concepts from convex analysis that are essential to this work are gathered in Section A. A key intermediate result in the proof of our main results is established in Section B with an elegant argument relying on Moreau-Yosida regularization.

## 1.8 Notations

*Vectors and matrices* For a  $d$ -dimensional Euclidean space  $\mathbb{R}^d$ , we endow it with its usual inner product  $\langle \cdot, \cdot \rangle$  and associated norm  $\|\cdot\|_2$ .  $\mathbf{Id}_d$  is the identity matrix on  $\mathbb{R}^d$ . For  $p \geq 1$ ,  $\|\cdot\|_p$  will denote the  $\ell_p$  norm of a vector with the usual adaptation for  $p = +\infty$ .

In the following, if  $T$  is a vector space,  $P_T$  denotes the orthogonal projector on  $T$ , and

$$\boldsymbol{\theta}_T = P_T \boldsymbol{\theta} \quad \text{and} \quad \mathbf{X}_T = \mathbf{X} P_T.$$

For a finite set  $\mathcal{C}$  we denote  $|\mathcal{C}|$  its cardinality. For  $I \subset \{1, \dots, p\}$ , we denote by  $I^c$  its complement.  $\boldsymbol{\theta}_I$  is the subvector whose entries are those of  $\boldsymbol{\theta}$  restricted to the indices in  $I$ , and  $\mathbf{X}_I$  the submatrix whose columns are those of  $\mathbf{X}$  indexed by  $I$ . For any matrix  $\mathbf{X}$ ,  $\mathbf{X}^\top$  denotes its transpose and  $\mathbf{X}^+$  its Moore-Penrose pseudo-inverse. For a linear operator  $\mathbf{A}$ ,  $\mathbf{A}^*$  is its adjoint.

*Sets* For a nonempty set  $\mathcal{C} \in \mathbb{R}^p$ , we denote  $\overline{\text{conv}}(\mathcal{C})$  the closure of its convex hull, and  $\iota_{\mathcal{C}}$  its indicator function, i.e.  $\iota_{\mathcal{C}}(\boldsymbol{\theta}) = 0$  if  $\boldsymbol{\theta} \in \mathcal{C}$  and  $+\infty$  otherwise. For a nonempty convex set  $\mathcal{C}$ , its *affine hull*  $\text{aff}(\mathcal{C})$  is the smallest affine manifold containing it. It is a translate of its *parallel subspace*  $\text{par}(\mathcal{C})$ , i.e.  $\text{par}(\mathcal{C}) = \text{aff}(\mathcal{C}) - \boldsymbol{\theta} = \mathbb{R}(\mathcal{C} - \mathcal{C})$ ; for any  $\boldsymbol{\theta} \in \mathcal{C}$ . The *relative interior*  $\text{ri}(\mathcal{C})$  of a convex set  $\mathcal{C}$  is the interior of  $\mathcal{C}$  for the topology relative to its affine full.

*Functions* A function  $f : \mathbb{R}^p \rightarrow \mathbb{R} \cup \{+\infty\}$  is closed (or lower semicontinuous (lsc)) if so is its epigraph. It is coercive if  $\lim_{\|\boldsymbol{\theta}\|_2 \rightarrow +\infty} f(\boldsymbol{\theta}) = +\infty$ , and strongly coercive if  $\lim_{\|\boldsymbol{\theta}\|_2 \rightarrow +\infty} f(\boldsymbol{\theta})/\|\boldsymbol{x}\|_2 = +\infty$ . The effective domain of  $f$  is  $\text{dom}(f) = \{\boldsymbol{\theta} \in \mathbb{R}^p : f(\boldsymbol{\theta}) < +\infty\}$  and  $f$  is proper if  $\text{dom}(f) \neq \emptyset$  as is the case when it is finite-valued. A function is said sublinear if it is convex and positively homogeneous. The Legendre-Fenchel conjugate of  $f$  is  $f^*(\boldsymbol{z}) = \sup_{\boldsymbol{\theta} \in \mathbb{R}^p} \langle \boldsymbol{z}, \boldsymbol{\theta} \rangle - f(\boldsymbol{\theta})$ . For  $f$  proper, the functions  $(f, f^*)$  obey the Fenchel-Young inequality

$$f(\boldsymbol{\theta}) + f^*(\boldsymbol{z}) \geq \langle \boldsymbol{z}, \boldsymbol{\theta} \rangle, \quad \forall (\boldsymbol{\theta}, \boldsymbol{z}) \in \mathbb{R}^p \times \mathbb{R}^p. \quad (5)$$

When  $f$  is a proper lower semicontinuous and convex function,  $(f, f^*)$  is actually the best pair for which this inequality cannot be tightened. For a function  $g$  on  $\mathbb{R}_+$ , the function  $g^+ : a \in \mathbb{R}_+ \mapsto g^+(a) = \sup_{t \geq 0} at - g(t)$  is called the monotone conjugate of  $g$ . The pair  $(g, g^+)$  obviously obeys (5) on  $\mathbb{R}_+ \times \mathbb{R}_+$ .

For a  $C^1$ -smooth function  $f$ ,  $\nabla f(\boldsymbol{\theta})$  is its (Euclidean) gradient. For a bivariate function  $g : (\boldsymbol{\eta}, \boldsymbol{y}) \in \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$  that is  $C^2$  with respect to the first variable  $\boldsymbol{\eta}$ , for any  $\boldsymbol{y}$ , we will denote  $\nabla g(\boldsymbol{\eta}, \boldsymbol{y})$  the gradient of  $g$  at  $\boldsymbol{\eta}$  with respect to the first variable.

The *subdifferential*  $\partial f(\boldsymbol{\theta})$  of a convex function  $f$  at  $\boldsymbol{\theta}$  is the set

$$\partial f(\boldsymbol{\theta}) = \{\boldsymbol{\eta} \in \mathbb{R}^p : f(\boldsymbol{\theta}') \geq f(\boldsymbol{\theta}) + \langle \boldsymbol{\eta}, \boldsymbol{\theta}' - \boldsymbol{\theta} \rangle, \quad \forall \boldsymbol{\theta}' \in \text{dom}(f)\}.$$

An element of  $\partial f(\boldsymbol{\theta})$  is a subgradient. If the convex function  $f$  is differentiable at  $\boldsymbol{\theta}$ , then its only subgradient is its gradient, i.e.  $\partial f(\boldsymbol{\theta}) = \{\nabla f(\boldsymbol{\theta})\}$ .

The *Bregman divergence* associated to a convex function  $f$  at  $\boldsymbol{\theta}$  with respect to  $\boldsymbol{\eta} \in \partial f(\boldsymbol{\theta}) \neq \emptyset$  is

$$D_f^\boldsymbol{\eta}(\bar{\boldsymbol{\theta}}, \boldsymbol{\theta}) = f(\bar{\boldsymbol{\theta}}) - f(\boldsymbol{\theta}) - \langle \boldsymbol{\eta}, \bar{\boldsymbol{\theta}} - \boldsymbol{\theta} \rangle.$$

The Bregman divergence is in general nonsymmetric. It is also nonnegative by convexity. When  $f$  is differentiable at  $\bar{\boldsymbol{\theta}}$ , we simply write  $D_f(\boldsymbol{\theta}, \bar{\boldsymbol{\theta}})$  (which is, in this case, also known as the Taylor distance).

## 2 Estimation with low-complexity penalties

The estimators  $\hat{\boldsymbol{\theta}}_n^{\text{PEN}}$  and  $\hat{\boldsymbol{\theta}}_n^{\text{EWA}}$  in (1) and (3) require two essential ingredients: the data loss term  $F$  and the prior penalty  $J$ . We here specify the class of such functions covered in our work, and provide illustrating examples.

### 2.1 Data loss

The class of loss functions  $F$  that we consider obey the following assumptions:

**(H.1)**  $F(\cdot, \mathbf{y}) : \mathbb{R}^n \rightarrow \mathbb{R}$  is  $C^1(\mathbb{R}^n)$  and uniformly convex for all  $\mathbf{y}$  of modulus  $\varphi$ , i.e.

$$F(\mathbf{v}, \mathbf{y}) \geq F(\mathbf{u}, \mathbf{y}) + \langle \nabla F(\mathbf{u}, \mathbf{y}), \mathbf{v} - \mathbf{u} \rangle + \varphi(\|\mathbf{v} - \mathbf{u}\|_2),$$

where  $\varphi : \mathbb{R}_+ \rightarrow \mathbb{R}_+$  is a convex non-decreasing function that vanishes only at 0.

**(H.2)** For any  $\bar{\boldsymbol{\theta}} \in \mathbb{R}^p$  and  $\mathbf{y} \in \mathbb{R}^n$ ,  $\int_{\mathbb{R}^p} \exp(-F(\mathbf{X}\boldsymbol{\theta}, \mathbf{y})/(n\beta)) |\langle \nabla F(\mathbf{X}\boldsymbol{\theta}, \mathbf{y}), \mathbf{X}(\bar{\boldsymbol{\theta}} - \boldsymbol{\theta}) \rangle| d\boldsymbol{\theta} < +\infty$ .

Recall that by Lemma 2, the monotone conjugate  $\varphi^+$  of  $\varphi$  is a proper, closed, convex, strongly coercive and non-decreasing function on  $\mathbb{R}_+$  that vanishes at 0. Moreover,  $\varphi^{++} = \varphi$ . The function  $\varphi^+$  is finite-valued on  $\mathbb{R}_+$  if  $\varphi$  is strongly coercive, and it vanishes only at 0 under e.g. Lemma 2(iii).

The class of data loss functions in (H.1) is fairly general. It is reminiscent of the negative log-likelihood in the regular exponential family. For the moment assumption (H.2) to be satisfied, it is sufficient that

$$\int_{\mathbb{R}^p} \exp(-\varphi(\|\mathbf{X}\boldsymbol{\theta}\|_2)/(n\beta)) \|\nabla F(\mathbf{X}\boldsymbol{\theta} + \mathbf{u}^*, \mathbf{y})\|_2 \|\mathbf{X}\boldsymbol{\theta} + (\mathbf{u}^* - \mathbf{X}\bar{\boldsymbol{\theta}})\|_2 d\boldsymbol{\theta} < +\infty,$$

where  $\mathbf{u}^*$  be a minimizer of  $F(\cdot, \mathbf{y})$ , which is unique by uniform convexity. We here provide an example.

*Example 1* Consider the case where<sup>3</sup>  $\varphi(t) = t^q/q$ ,  $q \in ]1, +\infty[$ , or equivalently  $\varphi^+(t) = t^{q^*}/q^*$  where  $1/q + 1/q^* = 1$ . For  $q = q^* = 2$ , (H.1) amounts to saying that  $F(\cdot, \mathbf{y})$  is strongly convex for all  $\mathbf{y}$ . In particular, (Bauschke and Combettes 2011, Proposition 10.13) shows that  $F(\mathbf{u}, \mathbf{y}) = \|\mathbf{u} - \mathbf{y}\|_2^q/q$  is uniformly convex for  $q \in [2, +\infty[$  with modulus  $\varphi(t) = C_q t^q/q$ , where  $C_q > 0$  is a constant that depends solely on  $q$ .

For (H.2) to be verified, it is sufficient that

$$\int_{\mathbb{R}^p} \exp(-\|\mathbf{X}\boldsymbol{\theta}\|_2^q/(qn\beta)) \|\nabla F(\mathbf{X}\boldsymbol{\theta} + \mathbf{u}^*, \mathbf{y})\|_2 \|(\mathbf{X}\boldsymbol{\theta} + \mathbf{u}^*) - \mathbf{X}\bar{\boldsymbol{\theta}}\|_2 d\boldsymbol{\theta} < +\infty.$$

In particular, taking  $F(\mathbf{u}, \mathbf{y}) = \|\mathbf{u} - \mathbf{y}\|_2^q/q$ ,  $q \in [2, +\infty[$ , we have  $\|\nabla F(\mathbf{u}, \mathbf{y})\|_2 = \|\mathbf{u} - \mathbf{y}\|_2^{q-1}$ , and thus (H.2) holds since

$$\int_{\mathbb{R}^p} \exp(-\|\mathbf{X}\boldsymbol{\theta}\|_2^q/(qn\beta)) \|\mathbf{y} - (\mathbf{X}\boldsymbol{\theta} + \mathbf{u}^*)\|_2^{q-1} \|\mathbf{X}\bar{\boldsymbol{\theta}} - (\mathbf{X}\boldsymbol{\theta} + \mathbf{u}^*)\|_2 d\boldsymbol{\theta} < +\infty.$$

## 2.2 Prior penalty

Recall the main definitions and results from convex analysis that are collected in Section A. Our main assumption on  $J$  is the following.

**(H.3)**  $J : \mathbb{R}^p \rightarrow \mathbb{R}$  is the gauge of a non-empty convex compact set containing the origin as an interior point.

<sup>3</sup> We consider a scaled version of  $\varphi$  for simplicity, but the same conclusions remain valid if we take  $\varphi(t) = Ct^q/q$ , with  $C > 0$ .



By Lemma 4, this assumption is equivalent to saying that  $J \stackrel{\text{def}}{=} \gamma_{\mathcal{C}}$  is proper, convex, positively homogeneous, finite-valued and coercive. In turn,  $J$  is locally Lipschitz continuous on  $\mathbb{R}^p$ . Observe also that by virtue of Lemma 5 and Lemma 3, the polar gauge  $J^\circ \stackrel{\text{def}}{=} \gamma_{\mathcal{C}^\circ}$  enjoys the same properties as  $J$  in (H.3).

### 2.3 Decomposability of the prior penalty

We are now in position to provide an important characterization of the subdifferential mapping of a function  $J$  satisfying (H.3). This characterization will play a pivotal role in our proof of the oracle inequality.

We start by defining some essential geometrical objects that were introduced in Vaïter et al (2015a).

**Definition 1 (Model Subspace)** Let  $\boldsymbol{\theta} \in \mathbb{R}^p$ . We denote by  $e_{\boldsymbol{\theta}}$  as

$$e_{\boldsymbol{\theta}} = P_{\text{aff}(\partial J(\boldsymbol{\theta}))}(0).$$

We denote

$$S_{\boldsymbol{\theta}} = \text{par}(\partial J(\boldsymbol{\theta})) \quad \text{and} \quad T_{\boldsymbol{\theta}} = S_{\boldsymbol{\theta}}^\perp.$$

$T_{\boldsymbol{\theta}}$  is coined the *model subspace* of  $\boldsymbol{\theta}$  associated to  $J$ .

It can be shown, see (Vaïter et al 2015a, Proposition 5), that  $\boldsymbol{\theta} \in T_{\boldsymbol{\theta}}$ , hence the name model subspace. When  $J$  is differentiable at  $\boldsymbol{\theta}$ , we have  $e_{\boldsymbol{\theta}} = \nabla J(\boldsymbol{\theta})$  and  $T_{\boldsymbol{\theta}} = \mathbb{R}^p$ . When  $J$  is the  $\ell_1$ -norm (Lasso), the vector  $e_{\boldsymbol{\theta}}$  is nothing but the sign of  $\boldsymbol{\theta}$ . Thus,  $e_{\boldsymbol{\theta}}$  can be viewed as a generalization of the sign vector. Observe also that  $e_{\boldsymbol{\theta}} = P_{T_{\boldsymbol{\theta}}}(\partial J(\boldsymbol{\theta}))$ , and thus  $e_{\boldsymbol{\theta}} \in T_{\boldsymbol{\theta}} \cap \text{aff}(\partial J(\boldsymbol{\theta}))$ . However, in general,  $e_{\boldsymbol{\theta}} \notin \partial J(\boldsymbol{\theta})$ .

We now provide a fundamental equivalent description of the subdifferential of  $J$  at  $\boldsymbol{\theta}$  in terms of  $e_{\boldsymbol{\theta}}$ ,  $T_{\boldsymbol{\theta}}$ ,  $S_{\boldsymbol{\theta}}$  and the polar gauge  $J^\circ$ .

**Theorem 1** Let  $J$  satisfy (H.3). Let  $\boldsymbol{\theta} \in \mathbb{R}^p$  and  $f_{\boldsymbol{\theta}} \in \text{ri}(\partial J(\boldsymbol{\theta}))$ .

(i) The subdifferential of  $J$  at  $\boldsymbol{\theta}$  reads

$$\begin{aligned} \partial J(\boldsymbol{\theta}) &= \text{aff}(\partial J(\boldsymbol{\theta})) \cap \mathcal{C}^\circ \\ &= \{ \boldsymbol{\eta} \in \mathbb{R}^n : \\ &\quad \boldsymbol{\eta}_{T_{\boldsymbol{\theta}}} = e_{\boldsymbol{\theta}} \text{ and } \inf_{\tau \geq 0} \max (J^\circ(\tau e_{\boldsymbol{\theta}} + \boldsymbol{\eta}_{S_{\boldsymbol{\theta}}} + (\tau - 1) P_{S_{\boldsymbol{\theta}}} f_{\boldsymbol{\theta}}), \tau) \leq 1 \}. \end{aligned}$$

(ii) For any  $\boldsymbol{\omega} \in \mathbb{R}^p$ ,  $\exists \boldsymbol{\eta} \in \partial J(\boldsymbol{\theta})$  such that

$$J(\boldsymbol{\omega}_{S_{\boldsymbol{\theta}}}) = \langle \boldsymbol{\eta}_{S_{\boldsymbol{\theta}}}, \boldsymbol{\omega}_{S_{\boldsymbol{\theta}}} \rangle.$$

*Proof* (i) This follows by piecing together (Vaïter et al 2015a, Theorem 1, Proposition 4 and Proposition 5(iii)).

(ii) From (Vaier et al 2015a, Proposition 5(iv)), we have

$$\sigma_{\partial J(\theta) - f_\theta}(\omega) = J(\omega_{S_\theta}) - \langle P_{S_\theta} f_\theta, \omega_{S_\theta} \rangle.$$

Thus there exists a supporting point  $\mathbf{v} \in \partial J(\theta) - f_\theta \subset S_\theta$  with normal vector  $\omega$  (Bauschke and Combettes 2011, Corollary 7.6(iii)), i.e.

$$\sigma_{\partial J(\theta) - f_\theta}(\omega) = \langle \mathbf{v}, \omega_{S_\theta} \rangle.$$

Taking  $\boldsymbol{\eta} = \mathbf{v} + f_\theta$  concludes the proof.

*Remark 1* The coercivity assumption in (H.3) is not needed for Theorem 1 to hold.

The decomposability of described in Theorem 1(i) depends on the particular choice of the mapping  $\theta \mapsto f_\theta \in \text{ri}(\partial J(\theta))$ . An interesting situation is encountered when  $e_\theta \in \text{ri}(J(\theta))$ , so that one can choose  $f_\theta = e_\theta$ . Strong gauges, see (Vaier et al 2015a, Definition 6), are precisely a class of gauges for which this situation occurs, and in this case, Theorem 1(i) has the simpler form

$$\partial J(\theta) = \text{aff}(\partial J(\theta)) \cap \mathcal{C}^\circ = \{ \boldsymbol{\eta} \in \mathbb{R}^n : \boldsymbol{\eta}_{T_\theta} = e_\theta \text{ and } J^\circ(\boldsymbol{\eta}_{S_\theta}) \leq 1 \}. \quad (6)$$

The Lasso, group Lasso and nuclear norms are typical examples of (symmetric) strong gauges. However, analysis sparsity penalties (e.g. the fused Lasso) or the  $\ell_\infty$ -penalty are not strong gauges, though they obviously satisfy (H.3). See the next section for a detailed discussion.

## 2.4 Calculus with the prior family

The family of penalties complying with (H.3) form a robust class enjoying important calculus rules. In particular it is closed under the sum and composition with an injective linear operator as we now prove.

**Lemma 1** *The set of functions satisfying (H.3) is closed under addition<sup>4</sup> and pre-composition by an injective linear operator. More precisely, the following holds:*

- (i) Let  $J$  and  $G$  be two gauges satisfying (H.3). Then  $H \stackrel{\text{def}}{=} J + G$  also obeys (H.3). Moreover,
  - (a)  $T_\theta^H = T_\theta^J \cap T_\theta^G$  and  $e_\theta^H = P_{T_\theta^H}(e_\theta^J + e_\theta^G)$ , where  $T_\theta^J$  and  $e_\theta^J$  (resp.  $T_\theta^G$  and  $e_\theta^G$ ) are the model subspace and vector at  $\theta$  associated to  $J$  (resp.  $G$ );
  - (b)  $H^\circ(\omega) = \max_{\rho \in [0,1]} \overline{\text{conv}}(\inf(\rho J^\circ(\omega), (1-\rho)G^\circ(\omega)))$ .
- (ii) Let  $J$  be a gauge satisfying (H.3), and  $\mathbf{D} : \mathbb{R}^q \rightarrow \mathbb{R}^p$  be surjective. Then  $H \stackrel{\text{def}}{=} J \circ \mathbf{D}^\top$  also fulfills (H.3). Moreover,

<sup>4</sup> It is obvious that the same holds with any positive linear combination.

- (a)  $T_{\boldsymbol{\theta}}^H = \text{Ker}(\mathbf{D}_{S_u^J}^\top)$  and  $e_{\boldsymbol{\theta}}^H = P_{T_{\boldsymbol{\theta}}^H} \mathbf{D} e_{\mathbf{u}}^J$ , where  $T_{\mathbf{u}}^J$  and  $e_{\mathbf{u}}^J$  are the model subspace and vector at  $\mathbf{u} \stackrel{\text{def}}{=} \mathbf{D}^\top \boldsymbol{\theta}$  associated to  $J$ ;
- (b)  $H^\circ(\boldsymbol{\omega}) = J^\circ(\mathbf{D}^+ \boldsymbol{\omega})$ , where  $\mathbf{D}^+ = \mathbf{D}^\top (\mathbf{D} \mathbf{D}^\top)^{-1}$ .

The outcome of Lemma 1 is naturally expected. For instance, assertion (i) states that combining several penalties/priors will promote objects living on the intersection of the respective low-complexity models. Similarly, for (ii), one promotes low-complexity in the image of the analysis operator  $\mathbf{D}^\top$ . It then follows that one has not to deploy an ad hoc analysis when linearly pre-composing or combining (or both) several penalties (e.g.  $\ell_1$ +nuclear norms for recovering sparse and low-rank matrices) since our unified analysis in Section 3 will apply to them just as well.

*Proof* (i) Convexity, positive homogeneity, coercivity and finite-valuedness are straightforward.

- (a) This is (Vaiter et al 2015a, Proposition 8(i)-(ii)).
- (b) We have from Lemma 5 and calculus rules on support functions,

$$\begin{aligned}
H^\circ(\boldsymbol{\omega}) &= \sigma_{J(\boldsymbol{\theta})+G(\boldsymbol{\theta})\leq 1}(\boldsymbol{\omega}) = \sup_{J(\boldsymbol{\theta})+G(\boldsymbol{\theta})\leq 1} \langle \boldsymbol{\omega}, \boldsymbol{\theta} \rangle \\
&= \max_{\rho \in [0,1]} \sup_{J(\boldsymbol{\theta})\leq \rho, G(\boldsymbol{\theta})\leq 1-\rho} \langle \boldsymbol{\omega}, \boldsymbol{\theta} \rangle \\
&= \max_{\rho \in [0,1]} \overline{\text{conv}} \left( \inf \left( \sigma_{J(\boldsymbol{\theta})\leq \rho}(\boldsymbol{\omega}), \sigma_{G(\boldsymbol{\theta})\leq 1-\rho}(\boldsymbol{\omega}) \right) \right) \\
&= \max_{\rho \in [0,1]} \overline{\text{conv}} \left( \inf \left( \rho \sigma_{J(\boldsymbol{\theta})\leq 1}(\boldsymbol{\omega}), (1-\rho) \sigma_{G(\boldsymbol{\theta})\leq 1}(\boldsymbol{\omega}) \right) \right) \\
&= \max_{\rho \in [0,1]} \overline{\text{conv}} \left( \inf \left( \rho J^\circ(\boldsymbol{\omega}), (1-\rho) G^\circ(\boldsymbol{\omega}) \right) \right),
\end{aligned}$$

where we used (Hiriart-Urruty and Lemaréchal 2001, Theorem V.3.3.3) in third row, positive homogeneity in the fourth, and Lemma 5 in the fifth.

- (ii) Again, convexity, positive homogeneity and finite-valuedness are immediate. Coercivity holds by injectivity of  $\mathbf{D}^\top$ .
- (a) This is (Vaiter et al 2015a, Proposition 10(i)-(ii)).
- (b) Denote  $J = \gamma_{\mathcal{C}}$ . We have

$$\begin{aligned}
H^\circ(\boldsymbol{\omega}) &= \sup_{\mathbf{D}^\top \boldsymbol{\theta} \in \mathcal{C}} \langle \boldsymbol{\omega}, \boldsymbol{\theta} \rangle \\
(\mathbf{D}^\top \text{ is injective}) &= \sup_{\mathbf{D}^\top \boldsymbol{\theta} \in \mathcal{C}} \langle \mathbf{D}^+ \boldsymbol{\omega}, \mathbf{D}^\top \boldsymbol{\theta} \rangle \\
&= \sup_{\mathbf{u} \in \mathcal{C} \cap \text{Span}(\mathbf{D}^\top)} \langle \mathbf{D}^+ \boldsymbol{\omega}, \mathbf{u} \rangle \\
&= \overline{\text{conv}} \left( \inf \left( J^\circ(\mathbf{D}^+ \boldsymbol{\omega}), \iota_{\text{Ker}(\mathbf{D})}(\mathbf{D}^+ \boldsymbol{\omega}) \right) \right) \\
&= J^\circ(\mathbf{D}^+ \boldsymbol{\omega}).
\end{aligned}$$

where in the last equality, we used the fact that  $\mathbf{D}^+ \boldsymbol{\omega} \in \text{Span}(\mathbf{D}^\top) = \text{Ker}(\mathbf{D})^\perp$ , and thus  $\iota_{\text{Ker}(\mathbf{D})}(\mathbf{D}^+ \boldsymbol{\omega}) = +\infty$  unless  $\boldsymbol{\omega} = 0$ , and  $J^\circ$  is

continuous and convex by [\(H.3\)](#) and [Lemma 5](#). In the fourth equality, we invoked [\(Hiriart-Urruty and Lemaréchal 2001, Theorem V.3.3.3\)](#) and [Lemma 5](#).

## 2.5 Examples

### 2.5.1 Lasso

The Lasso regularization is used to promote the sparsity of the minimizers, see [Bühlmann and van de Geer \(2011\)](#) for a comprehensive review. It corresponds to choosing  $J$  as the  $\ell_1$ -norm

$$J(\boldsymbol{\theta}) = \|\boldsymbol{\theta}\|_1 = \sum_{i=1}^p |\boldsymbol{\theta}_i|. \quad (7)$$

It is also referred to as  $\ell_1$ -synthesis in the signal processing community, in contrast to the more general  $\ell_1$ -analysis sparsity penalty detailed below.

We denote  $(\mathbf{a}_i)_{1 \leq i \leq p}$  the canonical basis of  $\mathbb{R}^p$  and  $\text{supp}(\boldsymbol{\theta}) \stackrel{\text{def}}{=} \{i \in \{1, \dots, p\} : \boldsymbol{\theta}_i \neq 0\}$ . Then,

$$T_{\boldsymbol{\theta}} = \text{Span}\{(\mathbf{a}_i)_{i \in \text{supp}(\boldsymbol{\theta})}\}, \quad (e_{\boldsymbol{\theta}})_i = \begin{cases} \text{sign}(\boldsymbol{\theta}_i) & \text{if } i \in \text{supp}(\boldsymbol{\theta}) \\ 0 & \text{otherwise} \end{cases}, \quad \text{and } J^\circ = \|\cdot\|_\infty. \quad (8)$$

### 2.5.2 Group Lasso

The group Lasso has been advocated to promote sparsity by groups, i.e. it drives all the coefficients in one group to zero together hence leading to group selection, see [Bakin \(1999\)](#); [Yuan and Lin \(2006\)](#); [Bach \(2008\)](#); [Wei and Huang \(2010\)](#) to cite a few. The group Lasso penalty with  $L$  groups reads

$$J(\boldsymbol{\theta}) = \|\boldsymbol{\theta}\|_{1,2} \stackrel{\text{def}}{=} \sum_{i=1}^L \|\boldsymbol{\theta}_{b_i}\|_2, \quad (9)$$

where  $\bigcup_{i=1}^L b_i = \{1, \dots, p\}$ ,  $b_i, b_j \subset \{1, \dots, p\}$ , and  $b_i \cap b_j = \emptyset$  whenever  $i \neq j$ . Define the group support as  $\text{supp}_{\mathcal{B}}(\boldsymbol{\theta}) \stackrel{\text{def}}{=} \{i \in \{1, \dots, L\} : \boldsymbol{\theta}_{b_i} \neq 0\}$ . Thus, one has

$$T_{\boldsymbol{\theta}} = \text{Span}\{(a_j)_{\{j : \exists i \in \text{supp}_{\mathcal{B}}(\boldsymbol{\theta}), j \in b_i\}}\},$$

$$(e_{\boldsymbol{\theta}})_{b_i} = \begin{cases} \frac{\boldsymbol{\theta}_{b_i}}{\|\boldsymbol{\theta}_{b_i}\|_2} & \text{if } i \in \text{supp}_{\mathcal{B}}(\boldsymbol{\theta}) \\ 0 & \text{otherwise} \end{cases}, \quad \text{and } J^\circ(\boldsymbol{\omega}) = \max_{i \in \{1, \dots, L\}} \|\boldsymbol{\omega}_{b_i}\|_2. \quad (10)$$

### 2.5.3 Analysis (group) Lasso

One can push the structured sparsity idea one step further by promoting group/block sparsity through a linear operator, i.e. analysis-type sparsity. Given a linear operator  $\mathbf{D} : \mathbb{R}^q \rightarrow \mathbb{R}^p$  (seen as a matrix), the analysis group sparsity penalty is

$$J(\boldsymbol{\theta}) = \|\mathbf{D}^\top \boldsymbol{\theta}\|_{1,2}. \quad (11)$$

This encompasses the 2-D isotropic total variation [Rudin et al \(1992\)](#). For when all groups of cardinality one, we have the analysis- $\ell_1$  penalty (a.k.a. general Lasso), which encapsulates several important penalties including that of the 1-D total variation [Rudin et al \(1992\)](#), and the fused Lasso [Tibshirani et al \(2005\)](#). The overlapping group Lasso [Jacob et al \(2009\)](#) is also a special case of (9) by taking  $\mathbf{D}^\top$  to be an operator that extract the blocks [Peyré et al \(2011\)](#); [Chen et al \(2010\)](#) (in which case  $\mathbf{D}$  has even orthogonal rows).

Let  $A_\theta = \bigcup_{i \in \text{supp}_B(\mathbf{D}^\top \boldsymbol{\theta})} b_i$  and  $A_\theta^c$  its complement. From Lemma 1(ii) and (10), we get

$$\begin{aligned} T_\theta &= \text{Ker}(\mathbf{D}_{A_\theta^c}^\top), \quad e_\theta = P_{T_\theta} \mathbf{D} e_{\mathbf{D}^\top \boldsymbol{\theta}}^{\|\cdot\|_{1,2}} \\ \text{where } \left( e_{\mathbf{D}^\top \boldsymbol{\theta}}^{\|\cdot\|_{1,2}} \right)_{b_i} &= \begin{cases} \frac{(\mathbf{D}^\top \boldsymbol{\theta})_{b_i}}{\|(\mathbf{D}^\top \boldsymbol{\theta})_{b_i}\|_2} & \text{if } i \in \text{supp}_B(\mathbf{D}^\top \boldsymbol{\theta}) \\ 0 & \text{otherwise.} \end{cases} \end{aligned} \quad (12)$$

If, in addition,  $\mathbf{D}$  is surjective, then by virtue of Lemma 1(ii) we also have

$$J^\circ(\boldsymbol{\omega}) = \|\mathbf{D}^+ \boldsymbol{\omega}\|_{\infty,2} \stackrel{\text{def}}{=} \max_{i \in \{1, \dots, L\}} \|(\mathbf{D}^+ \boldsymbol{\omega})_{b_i}\|_2. \quad (13)$$

### 2.5.4 Anti-sparsity

If the vector to be estimated is expected to be flat (anti-sparse), this can be captured using the  $\ell_\infty$  norm (a.k.a. Tchebychev norm) as prior

$$J(\boldsymbol{\theta}) = \|\boldsymbol{\theta}\|_\infty = \max_{i \in \{1, \dots, p\}} |\boldsymbol{\theta}_i|. \quad (14)$$

The  $\ell_\infty$  regularization has found applications in several fields [Jégou et al \(2012\)](#); [Lyubarskii and Vershynin \(2010\)](#); [Studer et al \(2012\)](#). Suppose that  $\boldsymbol{\theta} \neq 0$ , and define the saturation support of  $\boldsymbol{\theta}$  as  $I_\theta^{\text{sat}} \stackrel{\text{def}}{=} \{i \in \{1, \dots, p\} : |\boldsymbol{\theta}_i| = \|\boldsymbol{\theta}\|_\infty\} \neq \emptyset$ . From [\(Vaiter et al 2015a, Proposition 14\)](#), we have

$$\begin{aligned} T_\theta &= \{\bar{\boldsymbol{\theta}} \in \mathbb{R}^p : \bar{\boldsymbol{\theta}}_{I_\theta^{\text{sat}}} \in \mathbb{R} \text{sign}(\boldsymbol{\theta}_{I_\theta^{\text{sat}}})\}, \\ (e_\theta)_i &= \begin{cases} \text{sign}(\boldsymbol{\theta}_i)/|I_\theta^{\text{sat}}| & \text{if } i \in I_\theta^{\text{sat}} \\ 0 & \text{otherwise} \end{cases}, \quad \text{and } J^\circ = \|\cdot\|_1. \end{aligned} \quad (15)$$

### 2.5.5 Nuclear norm

The natural extension of low-complexity priors to matrices  $\boldsymbol{\theta} \in \mathbb{R}^{p_1 \times p_2}$  is to penalize the singular values of the matrix. Let  $\text{rank}(\boldsymbol{\theta}) = r$ , and  $\boldsymbol{\theta} = \mathbf{U} \text{diag}(\boldsymbol{\lambda}(\boldsymbol{\theta})) \mathbf{V}^\top$  be a reduced rank- $r$  SVD decomposition, where  $\mathbf{U} \in \mathbb{R}^{p_1 \times r}$  and  $\mathbf{V} \in \mathbb{R}^{p_2 \times r}$  have orthonormal columns, and  $\boldsymbol{\lambda}(\boldsymbol{\theta}) \in (\mathbb{R}_+ \setminus \{0\})^r$  is the vector of singular values  $(\lambda_1(\boldsymbol{\theta}), \dots, \lambda_r(\boldsymbol{\theta}))$  in non-increasing order. The nuclear norm of  $\boldsymbol{\theta}$  is

$$J(\boldsymbol{\theta}) = \|\boldsymbol{\theta}\|_* = \|\boldsymbol{\lambda}(\boldsymbol{\theta})\|_1. \quad (16)$$

This penalty is the best convex surrogate to enforce a low-rank prior. It has been widely used for various applications [Recht et al \(2010\)](#); [Candès and Recht \(2009\)](#); [Candès et al \(2011\)](#); [Fazel et al \(2001\)](#); [Candès et al \(2013\)](#).

Following e.g. ([Vaiter et al 2017](#), Example 21), we have

$$T_\boldsymbol{\theta} = \{\mathbf{U}\mathbf{A}^\top + \mathbf{B}\mathbf{V}^\top : \mathbf{A} \in \mathbb{R}^{p_2 \times r}, \mathbf{B} \in \mathbb{R}^{p_1 \times r}\}, \quad e_\boldsymbol{\theta} = \mathbf{U}\mathbf{V}^\top \quad (17)$$

and  $J^\circ(\boldsymbol{\omega}) = \|\boldsymbol{\omega}\|_{2 \rightarrow 2} = \|\boldsymbol{\lambda}(\boldsymbol{\omega})\|_\infty$ .

## 3 Oracle inequalities for a general loss

Before delving into the details, in the sequel, we will need a bit of notations.

We recall  $T_\boldsymbol{\theta}$  and  $e_\boldsymbol{\theta}$  the model subspace and vector associated to  $\boldsymbol{\theta}$  (see Definition 1). Denote  $S_\boldsymbol{\theta} = T_\boldsymbol{\theta}^\perp$ . Given two coercive finite-valued gauges  $J_1 = \gamma_{\mathcal{C}_1}$  and  $J_2 = \gamma_{\mathcal{C}_2}$ , and a linear operator  $\mathbf{A}$ , we define  $\|\mathbf{A}\|_{J_1 \rightarrow J_2}$  the *operator bound* as

$$\|\mathbf{A}\|_{J_1 \rightarrow J_2} = \sup_{\boldsymbol{\theta} \in \mathcal{C}_1} J_2(\mathbf{A}\boldsymbol{\theta}).$$

Note that  $\|\mathbf{A}\|_{J_1 \rightarrow J_2}$  is bounded (this follows from Lemma 4(v)). Furthermore, we have from Lemma 5 that

$$\begin{aligned} \|\mathbf{A}\|_{J_1 \rightarrow J_2} &= \sup_{\boldsymbol{\theta} \in \mathcal{C}_1} \sup_{\boldsymbol{\omega} \in \mathcal{C}_2^\circ} \langle \mathbf{A}^\top \boldsymbol{\omega}, \boldsymbol{\theta} \rangle = \sup_{\boldsymbol{\omega} \in \mathcal{C}_2^\circ} \sup_{\boldsymbol{\theta} \in \mathcal{C}_1} \langle \mathbf{A}^\top \boldsymbol{\omega}, \boldsymbol{\theta} \rangle = \sup_{\boldsymbol{\omega} \in \mathcal{C}_2^\circ} J_1^\circ(\mathbf{A}^\top \boldsymbol{\omega}) \\ &= \|\mathbf{A}^\top\|_{J_2^\circ \rightarrow J_1^\circ}. \end{aligned}$$

In the following, whenever it is clear from the context, to lighten notation when  $J_i$  is a norm, we write the subscript of the norm instead of  $J_i$  (e.g.  $p$  for the  $\ell_p$  norm,  $*$  for the nuclear norm, etc.).

Our main result will involve a measure of well-conditionedness of the design matrix  $\mathbf{X}$  when restricted to some subspace  $T$ . More precisely, for  $c > 0$ , we introduce the coefficient

$$\mathcal{X}(T, c) = \inf_{\{\boldsymbol{\omega} \in \mathbb{R}^p : J(\boldsymbol{\omega}_S) < cJ(\boldsymbol{\omega}_T)\}} \frac{\|\mathbf{P}_T\|_{2 \rightarrow J} \|\mathbf{X}\boldsymbol{\omega}\|_2}{n^{1/2}(J(\boldsymbol{\omega}_T) - J(\boldsymbol{\omega}_S)/c)}. \quad (18)$$

This generalizes the compatibility factor introduced in [van de Geer and Bühlmann \(2009\)](#) for the Lasso (and used in [Dalalyan et al \(2018\)](#)). The experienced

reader may have recognized that this factor is reminiscent of the null space property and restricted injectivity that play a central role in the analysis of the performance guarantees of penalized estimators (1); see Fadili et al (2013); Vaïter et al (2015a, 2018, 2017, 2015b). One can see in particular that  $\Upsilon(T, c)$  is larger than the smallest singular value of  $\mathbf{X}_T$ .

The oracle inequalities will be provided in terms of the loss

$$R_n(\boldsymbol{\theta}, \boldsymbol{\theta}_0) = \frac{1}{n} D_F(\mathbf{X}\boldsymbol{\theta}, \mathbf{X}\boldsymbol{\theta}_0).$$

### 3.1 Oracle inequality for $\hat{\boldsymbol{\theta}}_n^{\text{EWA}}$

We are now ready to establish our first main result: an oracle inequality for the EWA estimator (3).

**Theorem 2** Consider the EWA estimator  $\hat{\boldsymbol{\theta}}_n^{\text{EWA}}$  in (3) with the density (2), where  $F$  and  $J$  satisfy Assumptions (H.1)-(H.2) and (H.3). Then, for any  $\tau > 1$  such that  $\lambda_n \geq \tau J^\circ(-\mathbf{X}^\top \nabla F(\mathbf{X}\boldsymbol{\theta}_0, \mathbf{y}))/n$ , the following holds,

$$R_n(\hat{\boldsymbol{\theta}}_n^{\text{EWA}}, \boldsymbol{\theta}_0) \leq \inf_{\boldsymbol{\theta} \in \mathbb{R}^p} \left( R_n(\boldsymbol{\theta}, \boldsymbol{\theta}_0) + \frac{1}{n} \varphi^+ \left( \frac{\lambda_n \sqrt{n} (\tau J^\circ(e_{\boldsymbol{\theta}}) + 1) \|\mathbf{P}_{T_{\boldsymbol{\theta}}}\|_{2 \rightarrow J}}{\tau \Upsilon(T_{\boldsymbol{\theta}}, \frac{\tau J^\circ(e_{\boldsymbol{\theta}}) + 1}{\tau - 1})} \right) \right) + p\beta. \quad (19)$$

*Remark 2*

1. It should be emphasized that Theorem 2 is actually a deterministic statement for a fixed choice of  $\lambda_n$ . Probabilistic analysis will be required when the result is applied to particular statistical models as we will see later. For this, we will use concentration inequalities in order to provide bounds that hold with high probability over the data.
2. The oracle inequality is sharp. The remainder in it has two terms. The first one encodes the complexity of the model promoted by  $J$ . The second one,  $p\beta$ , captures the influence of the temperature parameter. In particular, taking  $\beta$  sufficiently small of the order  $O((pn)^{-1})$ , this term becomes  $O(n^{-1})$ .
3. When  $\varphi(t) = \nu t^2/2$ , i.e.  $F(\cdot, \mathbf{y})$  is  $\nu$ -strongly convex, then  $\varphi^+(t) = t^2/(2\nu)$ , and the remainder term becomes

$$\frac{\lambda_n^2 (\tau J^\circ(e_{\boldsymbol{\theta}}) + 1)^2 \|\mathbf{P}_{T_{\boldsymbol{\theta}}}\|_{2 \rightarrow J}^2}{2\tau^2 \nu \Upsilon(T_{\boldsymbol{\theta}}, \frac{\tau J^\circ(e_{\boldsymbol{\theta}}) + 1}{\tau - 1})^2}. \quad (20)$$

If, moreover,  $\nabla F$  is also  $\kappa$ -Lipschitz continuous, then it can be shown that  $R_n(\boldsymbol{\theta}, \boldsymbol{\theta}_0)$  is equivalent to a quadratic loss. This means that the oracle inequality in Theorem 2 can be stated in terms of the quadratic prediction

error. However, the inequality is not anymore sharp in this case as a constant factor equal to the condition number  $\kappa/\nu \geq 1$  naturally multiplies the right-hand side.

4. If  $J$  is such that  $e_{\boldsymbol{\theta}} \in \partial J(\boldsymbol{\theta}) \subset \mathcal{C}^\circ$  (typically for a strong gauge by (6)), then  $J^\circ(e_{\boldsymbol{\theta}}) \leq 1$  (in fact an equality if  $\boldsymbol{\theta} \neq 0$ ). Thus the term  $J^\circ(e_{\boldsymbol{\theta}})$  can be omitted in (19).
5. A close inspection of the proof of Theorem 2 reveals that the term  $p\beta$  can be improved to the smaller bound

$$p\beta + \left( V_n(\widehat{\boldsymbol{\theta}}_n^{\text{EWA}}) - \mathbb{E}_{\mu_n} [V_n(\boldsymbol{\theta})] \right),$$

where the upper-bound is a consequence of Jensen inequality.

*Proof* By convexity of  $J$  and assumption (H.1), we have for any  $\boldsymbol{\eta} \in \partial V_n(\boldsymbol{\theta})$  and any  $\bar{\boldsymbol{\theta}} \in \mathbb{R}^p$ ,

$$D_{V_n}^{\boldsymbol{\eta}}(\bar{\boldsymbol{\theta}}, \boldsymbol{\theta}) \geq \frac{1}{n} \varphi(\|\mathbf{X}\bar{\boldsymbol{\theta}} - \mathbf{X}\boldsymbol{\theta}\|_2).$$

Since  $\varphi$  is non-decreasing and convex,  $\varphi \circ \|\cdot\|_2$  is a convex function. Thus, taking the expectation w.r.t. to  $\mu_n$  on both sides and using Jensen inequality, we get

$$\begin{aligned} V_n(\bar{\boldsymbol{\theta}}) &\geq \mathbb{E}_{\mu_n} [V_n(\boldsymbol{\theta})] + \mathbb{E}_{\mu_n} [\langle \boldsymbol{\eta}, \bar{\boldsymbol{\theta}} - \boldsymbol{\theta} \rangle] + \frac{1}{n} \mathbb{E}_{\mu_n} [\varphi(\|\mathbf{X}\bar{\boldsymbol{\theta}} - \mathbf{X}\boldsymbol{\theta}\|_2)] \\ &\geq V_n(\widehat{\boldsymbol{\theta}}_n^{\text{EWA}}) + \mathbb{E}_{\mu_n} [\langle \boldsymbol{\eta}, \bar{\boldsymbol{\theta}} - \boldsymbol{\theta} \rangle] + \frac{1}{n} \varphi(\|\mathbf{X}\bar{\boldsymbol{\theta}} - \mathbf{X}\widehat{\boldsymbol{\theta}}_n^{\text{EWA}}\|_2). \end{aligned}$$

This holds for any  $\boldsymbol{\eta} \in \partial V_n(\boldsymbol{\theta})$ , and in particular at the minimal selection  $(\partial V_n(\boldsymbol{\theta}))^0$  (see Section B for details). It then follows from the pillar result in Proposition 5<sup>5</sup> that

$$\mathbb{E}_{\mu_n} [\langle (\partial V_n(\boldsymbol{\theta}))^0, \bar{\boldsymbol{\theta}} - \boldsymbol{\theta} \rangle] = -p\beta.$$

We thus deduce the inequality

$$V_n(\widehat{\boldsymbol{\theta}}_n^{\text{EWA}}) - V_n(\bar{\boldsymbol{\theta}}) \leq p\beta - \frac{1}{n} \varphi(\|\mathbf{X}\widehat{\boldsymbol{\theta}}_n^{\text{EWA}} - \mathbf{X}\bar{\boldsymbol{\theta}}\|_2), \quad \forall \bar{\boldsymbol{\theta}} \in \mathbb{R}^p. \quad (21)$$

By definition of the Bregman divergence, we have

$$\begin{aligned} &R_n(\widehat{\boldsymbol{\theta}}_n^{\text{EWA}}, \boldsymbol{\theta}_0) - R_n(\boldsymbol{\theta}, \boldsymbol{\theta}_0) \\ &= \frac{1}{n} \left( F(\mathbf{X}\widehat{\boldsymbol{\theta}}_n^{\text{EWA}}, \mathbf{y}) - F(\mathbf{X}\boldsymbol{\theta}, \mathbf{y}) + \langle -\mathbf{X}^\top \nabla F(\mathbf{X}\boldsymbol{\theta}_0, \mathbf{y}), \widehat{\boldsymbol{\theta}}_n^{\text{EWA}} - \boldsymbol{\theta} \rangle \right) \\ &= \left( V_n(\widehat{\boldsymbol{\theta}}_n^{\text{EWA}}) - V_n(\boldsymbol{\theta}) \right) + \frac{1}{n} \langle -\mathbf{X}^\top \nabla F(\mathbf{X}\boldsymbol{\theta}_0, \mathbf{y}), \widehat{\boldsymbol{\theta}}_n^{\text{EWA}} - \boldsymbol{\theta} \rangle \\ &\quad - \lambda_n (J(\widehat{\boldsymbol{\theta}}_n^{\text{EWA}}) - J(\boldsymbol{\theta})). \end{aligned}$$

<sup>5</sup> In the appendix, we provide a self-contained proof based on a novel Moreau-Yosida regularization argument. In (Dalalyan et al 2018, Corollary 1 and 2), an alternative proof is given using an absolute continuity argument since  $\mu_n$  is locally Lipschitz, hence a Sobolev function.



By virtue of the duality inequality (42), we have

$$\begin{aligned}
& R_n(\widehat{\boldsymbol{\theta}}_n^{\text{EWA}}, \boldsymbol{\theta}_0) - R_n(\boldsymbol{\theta}, \boldsymbol{\theta}_0) \\
& \leq \left( V_n(\widehat{\boldsymbol{\theta}}_n^{\text{EWA}}) - V_n(\boldsymbol{\theta}) \right) + \frac{1}{n} J^\circ(-\mathbf{X}^\top \nabla F(\mathbf{X}\boldsymbol{\theta}_0, \mathbf{y})) J(\widehat{\boldsymbol{\theta}}_n^{\text{EWA}} - \boldsymbol{\theta}) \\
& \quad - \lambda_n (J(\widehat{\boldsymbol{\theta}}_n^{\text{EWA}}) - J(\boldsymbol{\theta})) \\
& \leq \left( V_n(\widehat{\boldsymbol{\theta}}_n^{\text{EWA}}) - V_n(\boldsymbol{\theta}) \right) + \frac{\lambda_n}{\tau} \left( J(\widehat{\boldsymbol{\theta}}_n^{\text{EWA}} - \boldsymbol{\theta}) - \tau (J(\widehat{\boldsymbol{\theta}}_n^{\text{EWA}}) - J(\boldsymbol{\theta})) \right).
\end{aligned}$$

Denote  $\boldsymbol{\omega} = \widehat{\boldsymbol{\theta}}_n^{\text{EWA}} - \boldsymbol{\theta}$ . By virtue of (H.3), Theorem 1 and (42), we obtain

$$\begin{aligned}
J(\boldsymbol{\omega}) - \tau (J(\widehat{\boldsymbol{\theta}}_n^{\text{EWA}}) - J(\boldsymbol{\theta})) & \leq J(\boldsymbol{\omega}_{T_\theta}) + J(\boldsymbol{\omega}_{S_\theta}) - \tau \langle e_\theta, \boldsymbol{\omega}_{T_\theta} \rangle - \tau J(\boldsymbol{\omega}_{S_\theta}) \\
& \leq J(\boldsymbol{\omega}_{T_\theta}) + J(\boldsymbol{\omega}_{S_\theta}) + \tau J^\circ(e_\theta) J(\boldsymbol{\omega}_{T_\theta}) - \tau J(\boldsymbol{\omega}_{S_\theta}) \\
& = (\tau J^\circ(e_\theta) + 1) J(\boldsymbol{\omega}_{T_\theta}) - (\tau - 1) J(\boldsymbol{\omega}_{S_\theta}) \\
& \leq (\tau J^\circ(e_\theta) + 1) \left( J(\boldsymbol{\omega}_{T_\theta}) - \frac{\tau - 1}{\tau J^\circ(e_\theta) + 1} J(\boldsymbol{\omega}_{S_\theta}) \right).
\end{aligned}$$

This inequality together with (21) (applied with  $\bar{\boldsymbol{\theta}} = \boldsymbol{\theta}$ ) and (18) yield

$$\begin{aligned}
& R_n(\widehat{\boldsymbol{\theta}}_n^{\text{EWA}}, \boldsymbol{\theta}_0) - R_n(\boldsymbol{\theta}, \boldsymbol{\theta}_0) \\
& \leq p\beta - \frac{1}{n} \varphi(\|\mathbf{X}\boldsymbol{\omega}\|_2) + \frac{\lambda_n (\tau J^\circ(e_\theta) + 1) \|P_{T_\theta}\|_{2 \rightarrow J} \|\mathbf{X}\boldsymbol{\omega}\|_2}{n^{1/2} \tau \Upsilon\left(T_\theta, \frac{\tau J^\circ(e_\theta) + 1}{\tau - 1}\right)} \\
& \leq p\beta + \frac{1}{n} \varphi^+ \left( \frac{\lambda_n \sqrt{n} (\tau J^\circ(e_\theta) + 1) \|P_{T_\theta}\|_{2 \rightarrow J}}{\tau \Upsilon\left(T_\theta, \frac{\tau J^\circ(e_\theta) + 1}{\tau - 1}\right)} \right),
\end{aligned}$$

where we applied Fenchel-Young inequality (5) to get the last bound. Taking the infimum over  $\boldsymbol{\theta} \in \mathbb{R}^p$  yields the desired result.

*Stratifiable functions* Theorem 2 has a nice instantiation when  $\mathbb{R}^p$  can be partitioned into a collection of subsets  $\{\mathcal{M}_i\}_i$  that form a stratification of  $\mathbb{R}^p$ . That is,  $\mathbb{R}^p$  is a finite disjoint union  $\cup_i \mathcal{M}_i$  such that the partitioning sets  $\mathcal{M}_i$  (called strata) must fit nicely together and the stratification is endowed with a partial ordering for the closure operation. For example, it is known that a polyhedral function has a polyhedral stratification, and more generally, semialgebraic functions induce stratifications into finite disjoint unions of manifolds; see, e.g., Coste (2002). Another example is that of partly smooth convex functions thoroughly studied in Vaïter et al (2015a, 2018, 2017, 2015b) for various statistical and inverse problems. These functions induce a stratification into strata that are  $C^2$ -smooth submanifolds of  $\mathbb{R}^p$ . It turns out that all popular penalty functions discussed in this paper are partly smooth (see Vaïter et al

(2017, 2015b)). Let's denote  $\mathcal{M}$  the set of strata associated to  $J$ . With this notation at hand, the oracle inequality (19) now reads

$$R_n(\widehat{\boldsymbol{\theta}}_n^{\text{EWA}}, \boldsymbol{\theta}_0) \leq \inf_{\substack{\mathcal{M} \in \mathcal{M} \\ \boldsymbol{\theta} \in \mathcal{M}}} \left( R_n(\boldsymbol{\theta}, \boldsymbol{\theta}_0) + \frac{1}{n} \varphi^+ \left( \frac{\lambda_n \sqrt{n} (\tau^{J^\circ}(e_{\boldsymbol{\theta}}) + 1) \|P_{T_{\boldsymbol{\theta}}}\|_{2 \rightarrow J}}{\tau \mathcal{Y}(T_{\boldsymbol{\theta}}, \frac{\tau^{J^\circ}(e_{\boldsymbol{\theta}}) + 1}{\tau - 1})} \right) \right) + p\beta. \quad (22)$$

### 3.2 Oracle inequality for $\widehat{\boldsymbol{\theta}}_n^{\text{PEN}}$

The next result establishes that  $\widehat{\boldsymbol{\theta}}_n^{\text{PEN}}$  satisfies a sharp prediction oracle inequality that we will compare to (19).

**Theorem 3** Consider the penalized estimator  $\widehat{\boldsymbol{\theta}}_n^{\text{PEN}}$  in (1), where  $F$  and  $J$  satisfy Assumptions (H.1) and (H.3). Then, for any  $\tau > 1$  such that  $\lambda_n \geq \tau J^\circ(-\mathbf{X}^\top \nabla F(\mathbf{X}\boldsymbol{\theta}_0, \mathbf{y}))/n$ , the following holds,

$$R_n(\widehat{\boldsymbol{\theta}}_n^{\text{PEN}}, \boldsymbol{\theta}_0) \leq \inf_{\boldsymbol{\theta} \in \mathbb{R}^p} \left( R_n(\boldsymbol{\theta}, \boldsymbol{\theta}_0) + \frac{1}{n} \varphi^+ \left( \frac{\lambda_n \sqrt{n} (\tau^{J^\circ}(e_{\boldsymbol{\theta}}) + 1) \|P_{T_{\boldsymbol{\theta}}}\|_{2 \rightarrow J}}{\tau \mathcal{Y}(T_{\boldsymbol{\theta}}, \frac{\tau^{J^\circ}(e_{\boldsymbol{\theta}}) + 1}{\tau - 1})} \right) \right). \quad (23)$$

*Proof* The proof follows the same lines as that of Theorem 2 except that we use the fact that  $\widehat{\boldsymbol{\theta}}_n^{\text{PEN}}$  is a global minimizer of  $V_n$ , i.e.  $0 \in \partial V_n(\widehat{\boldsymbol{\theta}}_n^{\text{PEN}})$ . Indeed, we have for any  $\boldsymbol{\theta} \in \mathbb{R}^p$

$$V_n(\boldsymbol{\theta}) \geq V_n(\widehat{\boldsymbol{\theta}}_n^{\text{PEN}}) + \frac{1}{n} \varphi(\|\mathbf{X}\boldsymbol{\theta} - \mathbf{X}\widehat{\boldsymbol{\theta}}_n^{\text{PEN}}\|_2). \quad (24)$$

Continuing exactly as just after (21), replacing  $\widehat{\boldsymbol{\theta}}_n^{\text{EWA}}$  with  $\widehat{\boldsymbol{\theta}}_n^{\text{PEN}}$  and invoking (24) instead of (21), we arrive at the claimed result.

*Remark 3*

1. Observe that the penalized estimator  $\widehat{\boldsymbol{\theta}}_n^{\text{PEN}}$  does not require the moment assumption (H.2) for (23) to hold. The convexity assumption on  $\varphi$  in (H.1), which was important to apply Jensen's inequality in the proof of (19), is not needed either to get (23).
2. As we remarked for Theorem 2, Theorem 3 is also a deterministic statement for a fixed choice of  $\lambda_n$  that holds for any minimizer  $\widehat{\boldsymbol{\theta}}_n^{\text{PEN}}$ , which is not unique in general. The condition on  $\lambda_n$  is similar to the one in Negahban et al (2012) where authors established different guarantees for  $\widehat{\boldsymbol{\theta}}_n^{\text{PEN}}$ .

### 3.3 Discussion of $\hat{\boldsymbol{\theta}}_n^{\text{EWA}}$ vs $\hat{\boldsymbol{\theta}}_n^{\text{PEN}}$

One clearly sees that the difference between the prediction performance of  $\hat{\boldsymbol{\theta}}_n^{\text{EWA}}$  and  $\hat{\boldsymbol{\theta}}_n^{\text{PEN}}$  lies in the term  $p\beta$  (or rather its lower-bound in Remark 2-5). In particular, for  $\beta = O((pn)^{-1})$ , this term is on the order  $O(n^{-1})$ . This choice can be refined in most situations. In particular, for the case of quadratic loss, one can take  $\beta = O(\lambda_n^2 \|\text{P}_{T_{\boldsymbol{\theta}_0}}\|_{2 \rightarrow J}^2/p)$ , hence leading to remainder terms in (19) and (23) of the same order.

In view of this discussion, one may wonder what are the actual benefits of using  $\hat{\boldsymbol{\theta}}_n^{\text{EWA}}$  instead of  $\hat{\boldsymbol{\theta}}_n^{\text{PEN}}$ . Generalizing the arguments of Dalalyan et al (2018), we will show that  $\hat{\boldsymbol{\theta}}_n^{\text{EWA}}$  enjoys one main advantage compared to  $\hat{\boldsymbol{\theta}}_n^{\text{PEN}}$ . To simplify the discussion, we will focus on the case of linear regression (30) with Gaussian noise  $\boldsymbol{\xi} \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_n)$  and  $F$  is the quadratic loss.

The chief advantage of  $\hat{\boldsymbol{\theta}}_n^{\text{EWA}}$  is its stability as a function of the data and hyperparameters. It has been shown that for a large class of penalties  $J$ , including those studied here, the prediction  $\mathbf{X}\hat{\boldsymbol{\theta}}_n^{\text{PEN}}$  is a smooth function of  $\mathbf{y}$  outside a set of Lebesgue measure zero; see (Vaïter et al 2017, Theorem 2). Those authors also provided in (Vaïter et al 2017, Theorem 3) an expression of the Stein unbiased risk estimator (SURE). For instance, when  $J$  is the gauge of a polytope, the SURE is given by

$$\|\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\theta}}_n^{\text{PEN}}\|_2^2 + \sigma^2 \dim(T_{\hat{\boldsymbol{\theta}}_n^{\text{PEN}}}) - n\sigma^2.$$

The SURE was advocated as an automatic and objective way to choose  $\lambda$ . However, one can see that  $\dim(T_{\hat{\boldsymbol{\theta}}_n^{\text{PEN}}})$  is a non-smooth function of  $\lambda$ , which may lead to numerical instabilities in practice. In contrast, the SURE of  $\hat{\boldsymbol{\theta}}_n^{\text{EWA}}$ , whose closed-form is given in (Dalalyan et al 2018, (10)), is such that  $\mathbf{X}\hat{\boldsymbol{\theta}}_n^{\text{EWA}}$  is a continuous function of  $(\lambda, \beta) \in ]0, +\infty[^2$  and  $\mathbf{y} \in \mathbb{R}^n$ . This better regularity suggests that it would be wiser to use the SURE associated to  $\hat{\boldsymbol{\theta}}_n^{\text{EWA}}$  for an automatic choice of  $\lambda$ .

### 3.4 Oracle inequalities in probability

It remains to check when the event  $\mathcal{E} = \{\lambda_n \geq \tau J^\circ(-\mathbf{X}^\top \nabla F(\mathbf{X}\boldsymbol{\theta}_0, \mathbf{y}))/n\}$  holds with high probability when  $\mathbf{y}$  is random. We will use concentration inequalities in order to provide bounds that hold with high probability over the data. Toward this goal, we will need the following assumption.

- (H.4)  $\mathbf{y} = (\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_n)$  are independent random observations, and  $F(\mathbf{u}, \mathbf{y}) = \sum_{i=1}^n f_i(\mathbf{u}_i, \mathbf{y}_i)$ ,  $f_i : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ . Moreover,
- (i)  $\mathbb{E} [|f_i((\mathbf{X}\boldsymbol{\theta}_0)_i, \mathbf{y}_i)|] < +\infty, \forall 1 \leq i \leq n$ ;

- (ii)  $|f'_i((\mathbf{X}\boldsymbol{\theta}_0)_i, t)| \leq g(t)$ , where  $\mathbb{E}[g(\mathbf{y}_i)] < +\infty$ ,  $\forall 1 \leq i \leq n$ ;
- (iii) Bernstein moment condition:  $\forall 1 \leq i \leq n$  and all integers  $m \geq 2$ ,  
 $\mathbb{E}[|f'_i((\mathbf{X}\boldsymbol{\theta}_0)_i, \mathbf{y}_i)|^m] \leq m! \kappa^{m-2} \sigma_i^2 / 2$  for some constants  $\kappa > 0$ ,  $\sigma_i > 0$   
independent of  $n$ .  
Let  $\sigma^2 = \max_{1 \leq i \leq n} \sigma_i^2$ .

Observe that under [\(H.4\)](#), and by virtue of Lemma 5(iv) and ([Hiriart-Urruty and Lemaréchal 2001](#), Proposition V.3.3.4), we have

$$\begin{aligned} J^\circ(-\mathbf{X}^\top \nabla F(\mathbf{X}\boldsymbol{\theta}_0, \mathbf{y})) &= \sigma_{\mathcal{C}}(-\mathbf{X}^\top \nabla F(\mathbf{X}\boldsymbol{\theta}_0, \mathbf{y})) \\ &= \sup_{\mathbf{z} \in \mathbf{X}(\mathcal{C})} -\sum_{i=1}^n f'_i((\mathbf{X}\boldsymbol{\theta}_0)_i, \mathbf{y}_i) \mathbf{z}_i. \end{aligned} \quad (25)$$

Thus, checking the event  $\mathcal{E}$  amounts to establishing a deviation inequality for the supremum of an empirical process<sup>6</sup> above its mean under the weak Bernstein moment condition [\(H.4\)](#)(iii), which essentially requires that the  $f'_i((\mathbf{X}\boldsymbol{\theta}_0)_i, \mathbf{y}_i)$  have sub-exponential tails, We will first tackle the case where  $\mathcal{C}$  is the convex hull of a finite set (i.e.  $\mathcal{C}$  is a polytope).

### 3.4.1 Polyhedral penalty

We here suppose that  $J$  is a finite-valued gauge of  $\mathcal{C} = \overline{\text{conv}}(\mathcal{V})$ , where  $\mathcal{V}$  is finite, i.e.  $\mathcal{C}$  is a polytope with vertices ([Rockafellar 1996](#), Corollary 19.1.1). Our first oracle inequality in probability is the following.

**Proposition 1** Consider the estimators  $\hat{\boldsymbol{\theta}}_n^{\text{EWA}}$  and  $\hat{\boldsymbol{\theta}}_n^{\text{PEN}}$ , where  $F$  and  $J \stackrel{\text{def}}{=} \gamma_{\mathcal{C}}$  satisfy Assumptions [\(H.1\)](#), [\(H.2\)](#), [\(H.3\)](#) and [\(H.4\)](#), and  $\mathcal{C}$  is a polytope with vertices  $\mathcal{V}$ . Suppose that  $\text{rank}(\mathbf{X}) = n$  and let  $s(\mathbf{X}) = \max_{\mathbf{v} \in \mathcal{V}} \|\mathbf{X}\mathbf{v}\|_\infty$ . Choose

$$\lambda_n \geq \tau \sigma s(\mathbf{X}) \sqrt{\frac{2\delta \log(|\mathcal{V}|)}{n}} \left( 1 + \sqrt{2}\kappa/\sigma \sqrt{\frac{\delta \log(|\mathcal{V}|)}{n}} \right),$$

for some  $\tau > 1$  and  $\delta > 1$ . Then [\(19\)](#) and [\(23\)](#) hold with probability at least  $1 - 2|\mathcal{V}|^{1-\delta}$ .

*Proof* In view of Assumptions [\(H.1\)](#) and [\(H.4\)](#), one can differentiate under the expectation sign (Leibniz rule) to conclude that  $\mathbb{E}[F(\mathbf{X}\cdot, \mathbf{y})]$  is  $C^1$  at  $\boldsymbol{\theta}_0$  and  $\nabla \mathbb{E}[F(\mathbf{X}\boldsymbol{\theta}_0, \mathbf{y})] = \mathbf{X}^\top \mathbb{E}[\nabla F(\mathbf{X}\boldsymbol{\theta}_0, \mathbf{y})]$ . As  $\boldsymbol{\theta}_0$  minimizes the population risk, one has  $\nabla \mathbb{E}[F(\mathbf{X}\boldsymbol{\theta}_0, \mathbf{y})] = 0$ . Using the rank assumption on  $\mathbf{X}$ , we deduce that

$$\mathbb{E}[f'_i((\mathbf{X}\boldsymbol{\theta}_0)_i, \mathbf{y}_i)] = 0, \quad \forall 1 \leq i \leq n.$$

Moreover, [\(25\)](#) specializes to

$$J^\circ(-\mathbf{X}^\top \nabla F(\mathbf{X}\boldsymbol{\theta}_0, \mathbf{y})) = \sup_{\mathbf{z} \in \mathbf{X}(\mathcal{V})} -\sum_{i=1}^n f'_i((\mathbf{X}\boldsymbol{\theta}_0)_i, \mathbf{y}_i) \mathbf{z}_i.$$

<sup>6</sup> As  $\mathbf{X}(\mathcal{C})$  is compact, it has a dense countable subset.

Let  $t' = \lambda_n n / \tau$  and  $t = t' / s(\mathbf{X})$ . By the union bound and (25), we have

$$\begin{aligned} \mathbb{P}\left(\mathcal{J}^\circ(-\mathbf{X}^\top \nabla F(\mathbf{X}\boldsymbol{\theta}_0, \mathbf{y})) \geq t'\right) &\leq \mathbb{P}\left(\max_{\mathbf{z} \in \mathbf{X}(\mathcal{V})} -\sum_{i=1}^n f'_i((\mathbf{X}\boldsymbol{\theta}_0)_i, \mathbf{y}_i) \mathbf{z}_i \geq t'\right) \\ &\leq |\mathcal{V}| \max_{\mathbf{z} \in \mathbf{X}(\mathcal{V})} \mathbb{P}\left(\left|\sum_{i=1}^n f'_i((\mathbf{X}\boldsymbol{\theta}_0)_i, \mathbf{y}_i) \mathbf{z}_i\right| \geq t'\right) \\ &\leq |\mathcal{V}| \mathbb{P}\left(s(\mathbf{X}) \left|\sum_{i=1}^n f'_i((\mathbf{X}\boldsymbol{\theta}_0)_i, \mathbf{y}_i)\right| \geq t'\right) \\ &= |\mathcal{V}| \mathbb{P}\left(\left|\sum_{i=1}^n f'_i((\mathbf{X}\boldsymbol{\theta}_0)_i, \mathbf{y}_i)\right| \geq t\right). \end{aligned}$$

Owing to assumption (H.4)(iii), we are in position to apply the Bernstein inequality to get

$$\mathbb{P}\left(\mathcal{J}^\circ(-\mathbf{X}^\top \nabla F(\mathbf{X}\boldsymbol{\theta}_0, \mathbf{y})) \geq t\right) \leq 2|\mathcal{V}| \exp\left(-\frac{t^2}{2(\kappa t + n\sigma^2)}\right).$$

Every  $t$  such that

$$t \geq \sqrt{\delta \log(|\mathcal{V}|)} \left( \kappa \sqrt{\delta \log(|\mathcal{V}|)} + \sqrt{\kappa^2 \delta \log(|\mathcal{V}|) + 2n\sigma^2} \right),$$

satisfies  $t^2 \geq 2\delta \log(|\mathcal{V}|)(\kappa t + n\sigma^2)$ . Applying the trivial inequality  $\sqrt{a+b} \leq \sqrt{a} + \sqrt{b}$  to the bound on  $t$ , we conclude.

*Remark 4* In the monograph (Bühlmann and van de Geer 2011, Lemma 14.12), the authors derived an exponential deviation inequality for the supremum of an empirical process with finite  $\mathcal{V}$  and possibly unbounded empirical processes under a Bernstein moment condition similar to ours (in fact ours implies theirs). The very last part of our proof can be obtained by applying their result. We detailed it here for the sake of completeness.

*Lasso* To lighten the notation, let  $I_\boldsymbol{\theta} = \text{supp}(\boldsymbol{\theta})$ . From (8), it is easy to see that

$$\|P_{T_\boldsymbol{\theta}}\|_{2 \rightarrow 1} = \sqrt{|I_\boldsymbol{\theta}|} \quad \text{and} \quad \mathcal{J}^\circ(e_\boldsymbol{\theta}) = \|\text{sign}(\boldsymbol{\theta}_{I_\boldsymbol{\theta}})\|_\infty \leq 1,$$

where last bound holds as an equality whenever  $\boldsymbol{\theta} \neq 0$ . Further the  $\ell_1$  norm is the gauge of the cross-polytope (i.e. the unit  $\ell_1$  ball). Its vertex set  $\mathcal{V}$  is the set of unit-norm one-sparse vectors  $(\pm \mathbf{a}_i)_{1 \leq i \leq p}$ , where we recall  $(\mathbf{a}_i)_{1 \leq i \leq p}$  the canonical basis. Thus

$$|\mathcal{V}| = 2p \quad \text{and} \quad s(\mathbf{X}) = \max_{\mathbf{v} \in \mathcal{V}} \|\mathbf{X}\mathbf{v}\|_\infty = \max_{1 \leq i \leq p} \|\mathbf{X}_i\|_\infty.$$

Inserting this into Proposition 1, we obtain the following corollary.

**Corollary 1** Consider the estimators  $\hat{\boldsymbol{\theta}}_n^{\text{EWA}}$  and  $\hat{\boldsymbol{\theta}}_n^{\text{PEN}}$ , where where  $J$  is the Lasso penalty and  $F$  satisfies Assumptions **(H.1)**, **(H.2)** and **(H.4)**. Suppose that  $\text{rank}(\mathbf{X}) = n$  and take

$$\lambda_n \geq \tau \sigma s(\mathbf{X}) \sqrt{\frac{2\delta \log(2p)}{n}} \left( 1 + \sqrt{2}\kappa/\sigma \sqrt{\frac{\delta \log(2p)}{n}} \right),$$

for some  $\tau > 1$  and  $\delta > 1$ . Then, with probability at least  $1 - 2(2p)^{1-\delta}$ , the following holds

$$R_n(\hat{\boldsymbol{\theta}}_n^{\text{EWA}}, \boldsymbol{\theta}_0) \leq \inf_{\substack{I \subset \{1, \dots, p\} \\ \boldsymbol{\theta}: \text{supp}(\boldsymbol{\theta})=I}} \left( R_n(\boldsymbol{\theta}, \boldsymbol{\theta}_0) + \frac{1}{n} \varphi^+ \left( \frac{\lambda_n \sqrt{n}(\tau+1)\sqrt{|I|}}{\tau \mathcal{Y}(\text{Span}\{\mathbf{a}_i\}_{i \in I}, \frac{\tau+1}{\tau-1})} \right) \right) + p\beta, \quad (26)$$

and

$$R_n(\hat{\boldsymbol{\theta}}_n^{\text{PEN}}, \boldsymbol{\theta}_0) \leq \inf_{\substack{I \subset \{1, \dots, p\} \\ \boldsymbol{\theta}: \text{supp}(\boldsymbol{\theta})=I}} \left( R_n(\boldsymbol{\theta}, \boldsymbol{\theta}_0) + \frac{1}{n} \varphi^+ \left( \frac{\lambda_n \sqrt{n}(\tau+1)\sqrt{|I|}}{\tau \mathcal{Y}(\text{Span}\{\mathbf{a}_i\}_{i \in I}, \frac{\tau+1}{\tau-1})} \right) \right). \quad (27)$$

For  $\hat{\boldsymbol{\theta}}_n^{\text{PEN}}$ , we recover a similar scaling for  $\lambda_n$  and the oracle inequality as in [van de Geer \(2008\)](#), though in the latter the oracle inequality is not sharp unlike ours. Note that the above oracle inequality extends readily to the case of analysis/fused Lasso  $\|\mathbf{D}^\top \cdot\|_1$  where  $\mathbf{D}$  is surjective. We leave the details to the interested reader (see also the analysis group Lasso example in [Section 4](#)).

*Anti-sparsity* From [Section 2.5.4](#), recall the saturation support  $I_{\boldsymbol{\theta}}^{\text{sat}}$  of  $\boldsymbol{\theta}$ . From [\(15\)](#), we get

$$\|P_{T_{\boldsymbol{\theta}}}\|_{2 \rightarrow \infty} = 1 \quad \text{and} \quad J^\circ(e_{\boldsymbol{\theta}}) = \|\text{sign}(\boldsymbol{\theta}_{I_{\boldsymbol{\theta}}^{\text{sat}}})\|_1 / |I_{\boldsymbol{\theta}}^{\text{sat}}| \leq 1,$$

with equality whenever  $\boldsymbol{\theta} \neq 0$ . In addition, the  $\ell_\infty$  norm is the gauge of the hypercube whose vertex set is  $\mathcal{V} = \{\pm 1\}^p$ . Thus

$$|\mathcal{V}| = 2^p.$$

We have the following oracle inequalities.

**Corollary 2** Consider the estimators  $\hat{\boldsymbol{\theta}}_n^{\text{EWA}}$  and  $\hat{\boldsymbol{\theta}}_n^{\text{PEN}}$ , where where  $J$  is anti-sparsity penalty [\(14\)](#), and  $F$  satisfies Assumptions **(H.1)**, **(H.2)** and **(H.4)**. Suppose that  $\text{rank}(\mathbf{X}) = n$  and let  $s(\mathbf{X}) = \max_{i,j} |\mathbf{X}_{i,j}|$ . Choose

$$\lambda_n \geq \tau \sigma s(\mathbf{X}) \sqrt{2\delta \log(2)} \sqrt{\frac{p}{n}} \left( 1 + 2\kappa/\sigma \sqrt{\delta \log(2)} \sqrt{\frac{p}{n}} \right),$$

for some  $\tau > 1$  and  $\delta > 1$ . Then, with probability at least  $1 - 2^{-p(\delta-1)+1}$ , the following holds

$$R_n(\hat{\boldsymbol{\theta}}_n^{\text{EWA}}, \boldsymbol{\theta}_0) \leq \inf_{\substack{I \subset \{1, \dots, p\} \\ \boldsymbol{\theta}: I_{\boldsymbol{\theta}}^{\text{sat}}=I}} \left( R_n(\boldsymbol{\theta}, \boldsymbol{\theta}_0) + \frac{1}{n} \varphi^+ \left( \frac{\lambda_n \sqrt{n}(\tau+1)}{\tau \mathcal{Y}(\{\bar{\boldsymbol{\theta}}: \bar{\boldsymbol{\theta}}_I \in \mathbb{R} \text{sign}(\boldsymbol{\theta}_I)\}, \frac{\tau+1}{\tau-1})} \right) \right) + p\beta, \quad (28)$$

and

$$R_n(\hat{\boldsymbol{\theta}}_n^{\text{PEN}}, \boldsymbol{\theta}_0) \leq \inf_{\substack{I \subset \{1, \dots, p\} \\ \boldsymbol{\theta}: I_{\boldsymbol{\theta}}^{\text{sat}}=I}} \left( R_n(\boldsymbol{\theta}, \boldsymbol{\theta}_0) + \frac{1}{n} \varphi^+ \left( \frac{\lambda_n \sqrt{n}(\tau+1)}{\tau \mathcal{Y}(\{\bar{\boldsymbol{\theta}}: \bar{\boldsymbol{\theta}}_I \in \mathbb{R} \text{sign}(\boldsymbol{\theta}_I)\}, \frac{\tau+1}{\tau-1})} \right) \right). \quad (29)$$

We are not aware of any result of this kind in the literature. The bound imposed on  $\mathbf{X}$  is similar to what is generally assumed in the vector quantization literature [Lyubarskii and Vershynin \(2010\)](#); [Studer et al \(2012\)](#).

### 3.4.2 General penalty

Extending the above reasoning to a general penalty requires a deviation inequality for the supremum of an empirical process in (25) under the Bernstein moment condition (H.4) (iii), but without the need of uniform boundedness. This can be achieved via generic chaining along a tree using entropy with bracketing; see ([van de Geer and Lederer 2013](#), Theorem 8). The resulting deviation bound will thus depend on the entropies with bracketing. These quantities capture the complexity of the set  $\mathbf{X}(\mathcal{C})$  but are intricate to compute in general. This subject deserves further investigation that we leave to a future work.

*Remark 5 (Group Lasso)* Using the union bound, we have

$$\mathbb{P} \left( \max_{i \in \{1, \dots, L\}} \|\mathbf{X}_{b_i}^\top \nabla F(\mathbf{X}\boldsymbol{\theta}_0, \mathbf{y})\|_2 \geq \lambda_n n / \tau \right) \leq L \max_i \mathbb{P} \left( \|\mathbf{X}_{b_i}^\top \nabla F(\mathbf{X}\boldsymbol{\theta}_0, \mathbf{y})\|_2 \geq \lambda_n n / \tau \right).$$

This requires a concentration inequality for quadratic forms of independent random variables satisfying the Bernstein moment assumption above. We are not aware of any such result. But if our moment condition is strengthened to

$$\mathbb{E} \left[ |f'_i((\mathbf{X}\boldsymbol{\theta}_0)_i, \mathbf{y}_i)|^{2m} \right] \leq m! \kappa^{2(m-1)} \sigma_i^2 / 2, \quad \forall 1 \leq i \leq n, \forall m \geq 1,$$

then one can use ([Bellec 2014](#), Theorem 3). Indeed, assuming  $\max_i \|\mathbf{X}_i\|_2 \leq \sqrt{n}$ , which is a natural normalization on the design, we have by independence that

$$\begin{aligned} \mathbb{E} \left[ \|\mathbf{X}_{b_i}^\top \nabla F(\mathbf{X}\boldsymbol{\theta}_0, \mathbf{y})\|_2 \right] &\leq \mathbb{E} \left[ \|\mathbf{X}_{b_i}^\top \nabla F(\mathbf{X}\boldsymbol{\theta}_0, \mathbf{y})\|_2^2 \right]^{1/2} \\ &= \sigma \sqrt{\text{tr}(\mathbf{X}_{b_i}^\top \mathbf{X}_{b_i}) / 2} = \sigma \sqrt{\sum_{j \in b_i} \|\mathbf{X}_j\|_2^2 / 2} \leq \sigma \sqrt{Kn / 2}. \end{aligned}$$

It then follows that taking

$$\lambda_n \geq \tau \frac{\sigma\sqrt{K} + 16\kappa\sqrt{\delta\log(L)}}{\sqrt{n}}, \quad \delta > 1,$$

the oracle inequalities (34) and (35) hold for the group Lasso with probability at least  $1 - L^{1-\delta}$ . A similar result can be proved for the analysis group Lasso just as well (see Section 4.3.3).

#### 4 Oracle inequalities for low-complexity linear regression

In this section, we consider the classical linear regression problem where the  $n$  response-covariate pairs  $(\mathbf{y}_i, \mathbf{X}_i)$  are linked as

$$\mathbf{y} = \mathbf{X}\boldsymbol{\theta}_0 + \boldsymbol{\xi}, \quad (30)$$

where  $\boldsymbol{\xi}$  is a noise vector. The data loss will be set to  $F(\mathbf{u}, \mathbf{y}) = \frac{1}{2}\|\mathbf{y} - \mathbf{u}\|_2^2$ . This in turn entails that  $\varphi = \varphi^+ = \frac{1}{2}(\cdot)^2$  on  $\mathbb{R}_+$  and  $R_n(\boldsymbol{\theta}, \boldsymbol{\theta}_0) = \frac{1}{2n}\|\mathbf{X}\boldsymbol{\theta} - \mathbf{X}\boldsymbol{\theta}_0\|_2^2$ .

In this section, we assume that the noise  $\boldsymbol{\xi}$  is a zero-mean sub-Gaussian vector in  $\mathbb{R}^n$  with parameter  $\sigma$ . That is, its one-dimensional marginals  $\langle \boldsymbol{\xi}, \mathbf{z} \rangle$  are sub-Gaussian random variables  $\forall \mathbf{z} \in \mathbb{R}^n$ , i.e. they satisfy

$$\mathbb{P}(|\langle \boldsymbol{\xi}, \mathbf{z} \rangle| \geq t) \leq 2e^{-t^2/(2\|\mathbf{z}\|_2^2\sigma^2)}, \quad \forall \mathbf{z} \in \mathbb{R}^n. \quad (31)$$

In this case, the bounds of Section 3.4 can be improved.

##### 4.1 General penalty

As we will shortly show, the event  $\mathcal{E}$  will depend on the Gaussian width, a summary geometric quantity which, informally speaking, measures the size of the bulk of a set in  $\mathbb{R}^n$ .

**Definition 2** The Gaussian width of a subset  $\mathcal{S} \subset \mathbb{R}^n$  is defined as

$$w(\mathcal{S}) \stackrel{\text{def}}{=} \mathbb{E}[\sigma_{\mathcal{S}}(\mathbf{g})], \quad \text{where } \mathbf{g} \sim \mathcal{N}(0, \mathbf{Id}_n).$$

The concept of Gaussian width has appeared in the literature in different contexts. In particular, it has been used to establish sample complexity bounds to ensure exact recovery (noiseless case) and mean-square estimation stability (noisy case) for low-complexity penalized estimators from Gaussian measurements; see e.g. Rudelson and Vershynin (2008); Chandrasekaran et al (2012); Tropp (2015a); Vershynin (2015); Vaïter et al (2015b).

The Gaussian width has deep connections to convex geometry and it enjoys many useful properties. It is well-known that it is positively homogeneous, monotonic w.r.t. inclusion, and invariant under orthogonal transformations. Moreover,  $w(\overline{\text{conv}}(\mathcal{S})) = w(\mathcal{S})$ . From Lemma 3(ii)-(iii),  $w(\mathcal{S})$  is a non-negative finite quantity whenever the set  $\mathcal{S}$  is bounded and contains the origin.

We are now ready to state our oracle inequality in probability with sub-Gaussian noise.



**Proposition 2** *Let the data generated by (30) where  $\boldsymbol{\xi}$  is a zero-mean sub-Gaussian random vector with parameter  $\sigma$ . Consider the estimators  $\widehat{\boldsymbol{\theta}}_n^{\text{EWA}}$  and  $\widehat{\boldsymbol{\theta}}_n^{\text{PEN}}$ , where  $F$  and  $J \stackrel{\text{def}}{=} \gamma_{\mathcal{C}}$  satisfy Assumptions (H.1)-(H.2) and (H.3). Suppose that  $\lambda_n \geq \frac{\tau \sigma c_1 \sqrt{2 \log(c_2/\delta) w(\mathbf{X}(\mathcal{C}))}}{n}$ , for some  $\tau > 1$  and  $0 < \delta < \min(c_2, 1)$ , where  $c_1$  and  $c_2$  are positive absolute constants. Then with probability at least  $1 - \delta$ , (19) and (23) hold with the remainder term given by (20) with  $\nu = 1$ .*

The proof requires sophisticated ideas from the theory of generic chaining Talagrand (2005), but we only apply these results. The constants  $c_1$  and  $c_2$  can be traced back to the proof of these results as detailed in Talagrand (2005).

*Proof* First, from (31), we have the bound

$$\mathbb{P}(|\langle \boldsymbol{\xi}, \mathbf{z} - \mathbf{z}' \rangle| \geq t) \leq 2e^{-t^2/(2\|\mathbf{z} - \mathbf{z}'\|_2^2 \sigma^2)}, \quad \forall \mathbf{z}, \mathbf{z}' \in \mathbb{R}^n,$$

i.e. the increment condition (Talagrand 2005, (0.4)) is verified. Thus combining (25) with the probability bound in (Talagrand 2005, page 11), the generic chaining theorem (Talagrand 2005, Theorem 1.2.6) and the majorizing measure theorem (Talagrand 2005, Theorem 2.1.1), we have

$$\begin{aligned} \mathbb{P}\left(J^\circ(\mathbf{X}^\top \boldsymbol{\xi}) \geq \lambda_n n / \tau\right) &\leq \mathbb{P}\left(\sup_{\mathbf{z} \in \mathbf{X}(\mathcal{C})} \langle \boldsymbol{\xi}, \mathbf{z} \rangle \geq \sigma c_1 \sqrt{2 \log(c_2/\delta) w(\mathbf{X}(\mathcal{C}))}\right) \\ &\leq c_2 \exp\left(-\frac{\sigma^2 2 \log(c_2/\delta)}{2\sigma^2}\right) = \delta. \end{aligned}$$

If the noise is Gaussian, an enhanced version can be proved by invoking Gaussian concentration of Lipschitz functions Ledoux (2001).

**Proposition 3** *Let the data generated by (30) with noise  $\boldsymbol{\xi} \sim \mathcal{N}(0, \sigma^2 \mathbf{Id}_n)$ . Consider the estimators  $\widehat{\boldsymbol{\theta}}_n^{\text{EWA}}$  and  $\widehat{\boldsymbol{\theta}}_n^{\text{PEN}}$ , where  $F$  and  $J \stackrel{\text{def}}{=} \gamma_{\mathcal{C}}$  satisfy Assumptions (H.1)-(H.2) and (H.3). Suppose that  $\lambda_n \geq \frac{(1+\delta)\tau \sigma w(\mathbf{X}(\mathcal{C}))}{n}$ , for some  $\tau > 1$  and  $\delta > 0$ . Then with probability at least  $1 - \exp\left(-\frac{\delta^2 w(\mathbf{X}(\mathcal{C}))^2}{2\|\mathbf{X}\|_{J \rightarrow 2}^2}\right)$ , (19) and (23) hold with the remainder term given by (20) with  $\nu = 1$ .*

*Proof* Thanks to sublinearity (see Lemma 4(i) and Lemma 5), the function  $\boldsymbol{\xi} \mapsto J^\circ(\mathbf{X}^\top \boldsymbol{\xi})$  is Lipschitz continuous with Lipschitz constant  $\|\mathbf{X}^\top\|_{2 \rightarrow J^\circ} = \|\mathbf{X}\|_{J \rightarrow 2}$ . From (25), we also have

$$\mathbb{E}\left[J^\circ(\mathbf{X}^\top \boldsymbol{\xi})\right] = \sigma w(\mathbf{X}(\mathcal{C})).$$

Observe that  $\mathbf{X}(\mathcal{C})$  is a convex compact set containing the origin. Setting  $\epsilon = \lambda_n n / \tau - \sigma w(\mathbf{X}(\mathcal{C})) \geq \delta \sigma w(\mathbf{X}(\mathcal{C}))$ , it follows from (25) and the Gaussian

concentration of Lipschitz functions [Ledoux \(2001\)](#) that

$$\begin{aligned} \mathbb{P}\left(J^\circ(\mathbf{X}^\top \boldsymbol{\xi}) \geq \lambda_n n / \tau\right) &\leq \mathbb{P}\left(J^\circ(\mathbf{X}^\top \boldsymbol{\xi}) - \mathbb{E}\left[J^\circ(\mathbf{X}^\top \boldsymbol{\xi})\right] \geq \epsilon\right) \\ &\leq \mathbb{P}\left(J^\circ(\mathbf{X}^\top \boldsymbol{\xi} / \sigma) - w(\mathbf{X}(\mathcal{C})) \geq \delta w(\mathbf{X}(\mathcal{C}))\right) \\ &\leq \exp\left(-\frac{\delta^2 w(\mathbf{X}(\mathcal{C}))^2}{2\|\mathbf{X}\|_{J \rightarrow 2}^2}\right). \end{aligned}$$

Estimating theoretically the Gaussian width of a set<sup>7</sup> is a non-trivial problem that has been extensively studied in the areas of probability in Banach spaces and stochastic processes. There are classical bounds on the Gaussian width (Sudakov's and Dudley's inequalities), but they are difficult to estimate in most cases and neither of these bounds is tight for all sets. When the set is a convex cone (intersected with a sphere), tractable estimates based on polarity arguments were proposed in, e.g., [Chandrasekaran et al \(2012\)](#).

#### 4.2 Polyhedral penalty

When  $\mathcal{C}$  and is polytope, enhanced oracle inequalities can be obtained by invoking a simple union bound argument.

**Proposition 4** *Let the data generated by (30) where  $\boldsymbol{\xi}$  is a zero-mean sub-Gaussian random vector with parameter  $\sigma$ . Consider the estimators  $\hat{\boldsymbol{\theta}}_n^{\text{EWA}}$  and  $\hat{\boldsymbol{\theta}}_n^{\text{PEN}}$ , where  $F$  and  $J \stackrel{\text{def}}{=} \gamma_{\mathcal{C}}$  satisfy Assumptions [\(H.1\)](#)-[\(H.2\)](#) and [\(H.3\)](#), and moreover  $\mathcal{C}$  is a polytope with vertices  $\mathcal{V}$ . Suppose that  $\lambda_n \geq \frac{\tau\sigma(\max_{\mathbf{v} \in \mathcal{V}} \|\mathbf{X}\mathbf{v}\|_2) \sqrt{2\delta \log(|\mathcal{V}|)}}{n}$ , for some  $\tau > 1$  and  $\delta > 1$ . Then with probability at least  $1 - 2|\mathcal{V}|^{1-\delta}$ , [\(19\)](#) and [\(23\)](#) hold with the remainder term given by [\(20\)](#) with  $\nu = 1$ .*

*In particular, if  $\max_{\mathbf{v} \in \mathcal{V}} \|\mathbf{X}\mathbf{v}\|_2 = C\sqrt{n}$ , for a positive constant  $C$ , then one can take  $\lambda_n \geq C\tau\sigma \sqrt{\frac{2\delta \log(|\mathcal{V}|)}{n}}$ .*

*Proof* From [\(25\)](#) we have

$$J^\circ(\mathbf{X}^\top \boldsymbol{\xi}) = \max_{\mathbf{v} \in \mathcal{C}} \langle \mathbf{X}\mathbf{v}, \boldsymbol{\xi} \rangle = \max_{\mathbf{v} \in \mathcal{V}} \langle \mathbf{X}\mathbf{v}, \boldsymbol{\xi} \rangle,$$

where in the last inequality, we used the fact that a convex function attains its maximum on  $\mathcal{C}$  at an extreme point  $\mathcal{V}$ . Let  $\epsilon = \sigma(\max_{\mathbf{v} \in \mathcal{V}} \|\mathbf{X}\mathbf{v}\|_2) \sqrt{2\delta \log(|\mathcal{V}|)}$ .

<sup>7</sup> Not to mention its image with a linear operator as for  $\mathbf{X}(\mathcal{C})$ .

By the union bound, (31) and (25), we have

$$\begin{aligned}
\mathbb{P}\left(J^\circ(\mathbf{X}^\top \boldsymbol{\xi}) \geq \lambda_n n / \tau\right) &\leq \mathbb{P}\left(\max_{\mathbf{v} \in \mathcal{V}} \langle \mathbf{X} \mathbf{v}, \boldsymbol{\xi} \rangle \geq \epsilon\right) \\
&\leq |\mathcal{V}| \max_{\mathbf{v} \in \mathcal{V}} \mathbb{P}(\langle \mathbf{X} \mathbf{v}, \boldsymbol{\xi} \rangle \geq \epsilon) \\
&\leq |\mathcal{V}| \max_{\mathbf{v} \in \mathcal{V}} \mathbb{P}(|\langle \mathbf{X} \mathbf{v}, \boldsymbol{\xi} \rangle| \geq \epsilon) \\
&\leq 2|\mathcal{V}| \exp\left(-\epsilon^2 / (2\sigma^2 \max_{\mathbf{v} \in \mathcal{V}} \|\mathbf{X} \mathbf{v}\|_2^2)\right) \\
&\leq 2|\mathcal{V}|^{1-\delta}.
\end{aligned}$$

### 4.3 Applications

In this section, we exemplify our oracle inequalities for the penalties described in Section 2.5.

#### 4.3.1 Lasso

Recall the derivations for the Lasso in Section 3.4.1. We obtain the following corollary of Proposition 4.

**Corollary 3** *Let the data generated by (30) where  $\boldsymbol{\xi}$  is a zero-mean sub-Gaussian random vector with parameter  $\sigma$ . Assume that  $\mathbf{X}$  is such that  $\max_i \|\mathbf{X}_i\|_2 \leq \sqrt{n}$ . Consider the estimators  $\hat{\boldsymbol{\theta}}_n^{\text{EWA}}$  and  $\hat{\boldsymbol{\theta}}_n^{\text{PEN}}$ , where  $J$  is the Lasso penalty (7) and  $F$  satisfies Assumptions (H.1)-(H.2). Suppose that  $\lambda_n \geq \tau\sigma\sqrt{\frac{2\delta \log(2p)}{n}}$ , for some  $\tau > 1$  and  $\delta > 1$ . Then, with probability at least  $1 - 2(2p)^{1-\delta}$ , the following holds*

$$\begin{aligned}
\frac{1}{n} \|\mathbf{X} \hat{\boldsymbol{\theta}}_n^{\text{EWA}} - \mathbf{X} \boldsymbol{\theta}_0\|_2^2 &\leq \inf_{\substack{I \subset \{1, \dots, p\} \\ \boldsymbol{\theta}: \text{supp}(\boldsymbol{\theta})=I}} \left( \frac{1}{n} \|\mathbf{X} \boldsymbol{\theta} - \mathbf{X} \boldsymbol{\theta}_0\|_2^2 + \frac{\lambda_n^2 (\tau+1)^2 |I|}{\tau^2 \mathcal{Y} \left(\text{Span}\{\mathbf{a}_i\}_{i \in I}, \frac{\tau+1}{\tau-1}\right)^2} \right) \\
&\quad + 2p\beta,
\end{aligned} \tag{32}$$

and

$$\frac{1}{n} \|\mathbf{X} \hat{\boldsymbol{\theta}}_n^{\text{PEN}} - \mathbf{X} \boldsymbol{\theta}_0\|_2^2 \leq \inf_{\substack{I \subset \{1, \dots, p\} \\ \boldsymbol{\theta}: \text{supp}(\boldsymbol{\theta})=I}} \left( \frac{1}{n} \|\mathbf{X} \boldsymbol{\theta} - \mathbf{X} \boldsymbol{\theta}_0\|_2^2 + \frac{\lambda_n^2 (\tau+1)^2 |I|}{\tau^2 \mathcal{Y} \left(\text{Span}\{\mathbf{a}_i\}_{i \in I}, \frac{\tau+1}{\tau-1}\right)^2} \right). \tag{33}$$

The normalization on the design is natural. The remainder term grows as  $\frac{|I| \log(p)}{n}$ . The oracle inequality (33) recovers (Dalalyan et al 2018, Theorem 1) in the exactly sparse case, and (33) the one in (Sun and Zhang 2012, Theorem 4) (see also (Koltchinskii et al 2011, Theorem 11) and (Dalalyan et al 2017, Theorem 2)). It is worth mentioning, however, that (Dalalyan et al 2018,

Theorem 1) handles the inexactly sparse case while we do not. For the choice  $\beta = O(\sigma^2 |I| \log(2p)/(pn))$ , the remainder terms in (32) and (33) are of the same order. Observe that this choice of the temperature parameter is optimal in view of the results of Castillo et al (2015). These authors proved that for the  $\ell_1$  penalty, orthonormal design, noise  $\boldsymbol{\xi} \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_n)$ , and all choices of the form  $\lambda = C\sigma\sqrt{\log(n)/n}$ , then the pseudo-posterior  $\mu_n$  in (2) with temperature  $\beta = \sigma^2/n$  puts asymptotically no mass on the ball centered at  $\boldsymbol{\theta}_0$  of radius  $\sim \sqrt{\log(n)/n}$ .

#### 4.3.2 Group Lasso

Recall the notations in Section 2.5.2, and denote  $I_\boldsymbol{\theta} = \text{supp}_B(\boldsymbol{\theta})$  the set indexing active blocks in  $\boldsymbol{\theta}$ . From (10), we have

$$\|\mathbb{P}_{T_\boldsymbol{\theta}}\|_{2 \rightarrow J} = \sqrt{|I_\boldsymbol{\theta}|} \quad \text{and} \quad J^\circ(e_\boldsymbol{\theta}) = \|e_\boldsymbol{\theta}\|_{\infty, 2} \leq 1,$$

where the last bound holds as an equality whenever  $\boldsymbol{\theta} \neq 0$ .

We have the following oracle inequalities as corollaries of Proposition 2 and Proposition 3.

**Corollary 4** *Let the data generated by (30). Consider the estimators  $\widehat{\boldsymbol{\theta}}_n^{\text{EWA}}$  and  $\widehat{\boldsymbol{\theta}}_n^{\text{PEN}}$ , where  $F$  satisfies Assumptions (H.1)-(H.2), and  $J$  is the group Lasso (9) with  $L$  non-overlapping blocks of equal size  $K$ . Let  $s(\mathbf{X}) = \sqrt{\max_i \|\mathbf{X}_{b_i}^\top \mathbf{X}_{b_i}\|_{2 \rightarrow 2}}/n$ .*

(i)  $\boldsymbol{\xi}$  is a zero-mean sub-Gaussian random vector with parameter  $\sigma$ : suppose that

$\lambda_n \geq 3\tau\sigma s(\mathbf{X})c_1 \frac{\sqrt{2 \log(c_2/\delta)}(\sqrt{K} + \sqrt{2 \log(L)})}{\sqrt{n}}$ , for some  $\tau > 1$  and  $0 < \delta < \min(c_2, 1)$ , where  $c_1$  and  $c_2$  are the positive absolute constants in Proposition 2. Then, with probability at least  $1 - \delta$ , the following holds

$$\begin{aligned} \frac{1}{n} \|\mathbf{X} \widehat{\boldsymbol{\theta}}_n^{\text{EWA}} - \mathbf{X} \boldsymbol{\theta}_0\|_2^2 &\leq \inf_{\substack{I \subset \{1, \dots, L\} \\ \boldsymbol{\theta}: \text{supp}_B(\boldsymbol{\theta})=I}} \left( \frac{1}{n} \|\mathbf{X} \boldsymbol{\theta} - \mathbf{X} \boldsymbol{\theta}_0\|_2^2 \right. \\ &\quad \left. + \frac{\lambda_n^2 (\tau+1)^2 |I|}{\tau^2 \mathcal{Y} \left( \text{Span}\{a_j\}_{j \in b_i, i \in I, \frac{\tau+1}{\tau-1}} \right)^2} \right) + 2p\beta, \end{aligned} \quad (34)$$

and

$$\begin{aligned} \frac{1}{n} \|\mathbf{X} \widehat{\boldsymbol{\theta}}_n^{\text{PEN}} - \mathbf{X} \boldsymbol{\theta}_0\|_2^2 &\leq \inf_{\substack{I \subset \{1, \dots, L\} \\ \boldsymbol{\theta}: \text{supp}_B(\boldsymbol{\theta})=I}} \left( \frac{1}{n} \|\mathbf{X} \boldsymbol{\theta} - \mathbf{X} \boldsymbol{\theta}_0\|_2^2 \right. \\ &\quad \left. + \frac{\lambda_n^2 (\tau+1)^2 |I|}{\tau^2 \mathcal{Y} \left( \text{Span}\{a_j\}_{j \in b_i, i \in I, \frac{\tau+1}{\tau-1}} \right)^2} \right). \end{aligned} \quad (35)$$

(ii)  $\boldsymbol{\xi} \sim \mathcal{N}(0, \sigma^2 \mathbf{Id}_n)$ : suppose that  $\lambda_n \geq \tau \sigma s(\mathbf{X}) \frac{\sqrt{K} + \sqrt{2\delta \log(L)}}{\sqrt{n}}$ , for some  $\tau > 1$  and  $\delta > 1$ . Then, with probability at least  $1 - L^{1-\delta}$ , (34) and (35) hold.

When  $s(\mathbf{X}) = O(1)$ <sup>8</sup>, the first remainder term is on the order  $\frac{|I|(\sqrt{K} + \sqrt{2\log(L)})^2}{n}$ . This is similar to the scaling that has been provided in the literature for EWA with other group sparsity priors and noises [Rigollet and Tsybakov \(2012\)](#); [Duy Luu et al \(2016\)](#). Similar rates were given for  $\hat{\boldsymbol{\theta}}_n^{\text{PEN}}$  with the group Lasso in [Negahban et al \(2012\)](#); [Lounici et al \(2011\)](#); [van de Geer \(2014\)](#).

*Proof*

(i) This is a consequence of Proposition 2, for which we need to bound

$$w(\mathbf{X}(\mathcal{C})) = \mathbb{E} \left[ \max_{i \in \{1, \dots, L\}} \|\mathbf{X}_{b_i}^\top \mathbf{g}\|_2 \right].$$

We first have, for any block  $b_i$

$$\mathbb{E} \left[ \|\mathbf{X}_{b_i}^\top \mathbf{g}\|_2 \right] \leq \mathbb{E} \left[ \|\mathbf{X}_{b_i}^\top \mathbf{g}\|_2^2 \right]^{1/2} \leq s(\mathbf{X}) \sqrt{Kn}.$$

Furthermore,  $\|\mathbf{X}_{b_i}^\top \cdot\|_2$  is Lipschitz continuous with Lipschitz constant  $s(\mathbf{X})\sqrt{n}$ . Thus the union bound and Gaussian concentration of Lipschitz functions [Ledoux \(2001\)](#) yield, for any  $t > 0$ ,

$$\begin{aligned} & \mathbb{P} \left( \max_{i \in \{1, \dots, L\}} \|\mathbf{X}_{b_i}^\top \mathbf{g}\|_2 \geq s(\mathbf{X}) \sqrt{Kn} + t \right) \\ & \leq \sum_{i=1}^L \mathbb{P} \left( \|\mathbf{X}_{b_i}^\top \boldsymbol{\xi}\|_2 - \mathbb{E} \left[ \|\mathbf{X}_{b_i}^\top \boldsymbol{\xi}\|_2 \right] \geq t \right) \leq L \exp \left( -\frac{t^2}{2s(\mathbf{X})^2 n} \right). \end{aligned}$$

Let  $\kappa = s(\mathbf{X})(\sqrt{Kn} + \sqrt{2n \log(L)})$ .  $w(\mathbf{X}(\mathcal{C}))$  can be expressed as

$$\begin{aligned} w(\mathbf{X}(\mathcal{C})) &= \int_0^\infty \mathbb{P} \left( \max_{i \in \{1, \dots, L\}} \|\mathbf{X}_{b_i}^\top \mathbf{g}\|_2 \geq u \right) du \\ &\leq \int_0^\kappa du + \int_\kappa^\infty e^{-\frac{(u - s(\mathbf{X})\sqrt{Kn})^2 - 2s(\mathbf{X})^2 n \log(L)}{2n}} du \\ &= \kappa + s(\mathbf{X})\sqrt{n} \int_{\kappa/(s(\mathbf{X})\sqrt{n})}^\infty e^{-\frac{(s - \sqrt{K})^2 - 2\log(L)}{2}} du \\ &\leq \kappa + s(\mathbf{X})\sqrt{n} \int_{\kappa/(s(\mathbf{X})\sqrt{n})}^\infty e^{-\frac{s - \kappa/(s(\mathbf{X})\sqrt{n})}{2}} du = \kappa + 2s(\mathbf{X})\sqrt{n} \leq 3\kappa. \end{aligned}$$

<sup>8</sup> This is for instance the case if  $\mathbf{X}$  is drawn from the standard Gaussian ensemble and  $K = O(n)$  (the  $O(\cdot)$  is in fact even  $o(\cdot)$  as the remainder term is supposed to go to 0 as  $n \rightarrow +\infty$ ). In this case, classical concentration bounds of the largest eigenvalue of a Wishart matrix allow to conclude that  $s(\mathbf{X}) = O(1 + \sqrt{K/n}) = O(1)$  with high probability.

- (ii) The proof follows the lines of Proposition 3 where we additionally use the union bound. Indeed,

$$\begin{aligned}
& \mathbb{P} \left( \max_{i \in \{1, \dots, L\}} \|\mathbf{X}_{b_i}^\top \boldsymbol{\xi}\|_2 \geq \lambda_n n / \tau \right) \\
& \leq \sum_{i=1}^L \mathbb{P} \left( \|\mathbf{X}_{b_i}^\top \boldsymbol{\xi}\|_2 - \mathbb{E} \left[ \|\mathbf{X}_{b_i}^\top \boldsymbol{\xi}\|_2 \right] \geq \lambda_n n / \tau - \mathbb{E} \left[ \|\mathbf{X}_{b_i}^\top \boldsymbol{\xi}\|_2 \right] \right) \\
& \leq \sum_{i=1}^L \mathbb{P} \left( \|\mathbf{X}_{b_i}^\top \boldsymbol{\xi}\|_2 - \mathbb{E} \left[ \|\mathbf{X}_{b_i}^\top \boldsymbol{\xi}\|_2 \right] \geq \lambda_n n / \tau - \sigma s(\mathbf{X}) \sqrt{Kn} \right) \\
& \leq \sum_{i=1}^L \mathbb{P} \left( \|\mathbf{X}_{b_i}^\top \boldsymbol{\xi}\|_2 - \mathbb{E} \left[ \|\mathbf{X}_{b_i}^\top \boldsymbol{\xi}\|_2 \right] \geq \sigma s(\mathbf{X}) \sqrt{2\delta n \log(L)} \right) \\
& \leq L \exp(-\delta \log(L)) = L^{1-\delta},
\end{aligned}$$

where used the Gaussian concentration of Lipschitz functions Ledoux (2001) in the last inequality.

We observe in passing that another way to prove the oracle inequalities in the sub-Gaussian is to use Dudley's inequality on the sphere in  $\mathbb{R}^K$  after applying a union bound on the  $L$  blocks. In addition, in the Gaussian case, the (similar) bound  $\lambda_n \geq 3\delta\tau\sigma s(\mathbf{X}) \frac{\sqrt{K} + \sqrt{2\log(L)}}{\sqrt{n}}$  can be obtained by combining Proposition 3 and the estimate  $w(\mathbf{X}(\mathcal{C})) \leq 3s(\mathbf{X})(\sqrt{Kn} + \sqrt{2n\log(L)})$  in the proof of (i). The corresponding probability of success would be at least  $1 - L^{-9(\delta-1)^2}$ .

#### 4.3.3 Analysis group Lasso

We now turn to the prior penalty (11). Recall the notations in Section 2.5.3, and remind  $\Lambda_\theta = \bigcup_{i \in \text{supp}_B(\mathbf{D}^\top \theta)} b_i$ . We assume that  $\mathbf{D}$  is a frame of  $\mathbb{R}^p$ , hence surjective, meaning that there exist  $c, d > 0$  such that for any  $\boldsymbol{\omega} \in \mathbb{R}^p$

$$d \|\boldsymbol{\omega}\|_2^2 \leq \|\mathbf{D}^\top \boldsymbol{\omega}\|_2^2 \leq c \|\boldsymbol{\omega}\|_2^2.$$

This together with (12)-(13) and Cauchy-Schwarz inequality entail

$$\begin{aligned}
\|P_{T_\theta}\|_{2 \rightarrow J} &= \sup_{\|\boldsymbol{\omega}_{T_\theta}\|_2 \leq 1} \|\mathbf{D}^\top \boldsymbol{\omega}_{T_\theta}\|_{1,2} \leq \sqrt{c} \sup_{\|\mathbf{D}^\top \boldsymbol{\omega}_{T_\theta}\|_2 \leq 1} \|\mathbf{D}^\top \boldsymbol{\omega}_{T_\theta}\|_{1,2} \\
&= \sqrt{c} \sup_{\|\mathbf{D}_{\Lambda_\theta}^\top \boldsymbol{\omega}_{T_\theta}\|_2 \leq 1} \|\mathbf{D}_{\Lambda_\theta}^\top \boldsymbol{\omega}_{T_\theta}\|_{1,2} \\
&= \sqrt{c} \sqrt{|\text{supp}_B(\mathbf{D}^\top \theta)|}.
\end{aligned}$$

Note, however, that from (12), we do not have in general  $C(\mathbf{D}, \theta) \stackrel{\text{def}}{=} \left\| \mathbf{D}^+ P_{\text{Ker}(\mathbf{D}_{\Lambda_\theta}^\top)} \mathbf{D} e_{\mathbf{D}^\top \theta} \right\|_{\infty, 2} \leq 1$ .

With exactly the same arguments to those for proving Corollary 4, replacing  $\mathbf{X}$  by  $\mathbf{X}\mathbf{D}$ , we arrive at the following oracle inequalities.

**Corollary 5** *Let the data generated by (30). Consider the estimators  $\hat{\boldsymbol{\theta}}_n^{\text{EWA}}$  and  $\hat{\boldsymbol{\theta}}_n^{\text{PEN}}$ , where  $F$  satisfies Assumptions (H.1)-(H.2), and  $J$  is the analysis group Lasso (11) with  $L$  blocks of equal size  $K$ . Assume that  $\mathbf{D}$  is a frame, and let  $s(\mathbf{X}\mathbf{D}) = \sqrt{\max_i \left\| \mathbf{D}_{b_i}^\top \mathbf{X}^\top \mathbf{X} \mathbf{D}_{b_i} \right\|_{2 \rightarrow 2}}/n$ .*

(i)  $\boldsymbol{\xi}$  is a zero-mean sub-Gaussian random vector with parameter  $\sigma$ : suppose that

$\lambda_n \geq 3\tau\sigma s(\mathbf{X}\mathbf{D})c_1 \frac{\sqrt{\log(c_2/\delta)}(\sqrt{K} + \sqrt{2\log(L)})}{\sqrt{n}}$ , for some  $\tau > 1$  and  $0 < \delta < \min(c_2, 1)$ , where  $c_1$  and  $c_2$  are the positive absolute constants in Proposition 2. Then, with probability at least  $1 - \delta$ , the following holds

$$\begin{aligned} \frac{1}{n} \left\| \mathbf{X} \hat{\boldsymbol{\theta}}_n^{\text{EWA}} - \mathbf{X} \boldsymbol{\theta}_0 \right\|_2^2 \leq & \inf_{\substack{I \subset \{1, \dots, L\} \\ \boldsymbol{\theta}: \text{supp}_{\mathcal{B}}(\mathbf{D}^\top \boldsymbol{\theta}) = I}} \left( \frac{1}{n} \left\| \mathbf{X} \boldsymbol{\theta} - \mathbf{X} \boldsymbol{\theta}_0 \right\|_2^2 \right. \\ & \left. + \frac{c\lambda_n^2 (\tau C(\mathbf{D}, \boldsymbol{\theta}) + 1)^2 |I|}{\tau^2 \Upsilon \left( \text{Ker}(\mathbf{D}_{A_\theta}^\top), \frac{\tau C(\mathbf{D}, \boldsymbol{\theta}) + 1}{\tau - 1} \right)^2} \right) + 2p\beta, \end{aligned} \quad (36)$$

and

$$\begin{aligned} \frac{1}{n} \left\| \mathbf{X} \hat{\boldsymbol{\theta}}_n^{\text{PEN}} - \mathbf{X} \boldsymbol{\theta}_0 \right\|_2^2 \leq & \inf_{\substack{I \subset \{1, \dots, L\} \\ \boldsymbol{\theta}: \text{supp}_{\mathcal{B}}(\mathbf{D}^\top \boldsymbol{\theta}) = I}} \left( \frac{1}{n} \left\| \mathbf{X} \boldsymbol{\theta} - \mathbf{X} \boldsymbol{\theta}_0 \right\|_2^2 \right. \\ & \left. + \frac{c\lambda_n^2 (\tau C(\mathbf{D}, \boldsymbol{\theta}) + 1)^2 |I|}{\tau^2 \Upsilon \left( \text{Ker}(\mathbf{D}_{A_\theta}^\top), \frac{\tau C(\mathbf{D}, \boldsymbol{\theta}) + 1}{\tau - 1} \right)^2} \right). \end{aligned} \quad (37)$$

(ii)  $\boldsymbol{\xi} \sim \mathcal{N}(0, \sigma^2 \mathbf{Id}_n)$ : suppose that  $\lambda_n \geq \tau\sigma s(\mathbf{X}\mathbf{D}) \frac{\sqrt{K} + \sqrt{2\delta \log(L)}}{\sqrt{n}}$ , for some  $\tau > 1$  and  $\delta > 1$ . Then, with probability at least  $1 - L^{1-\delta}$ , (36) and (37) hold.

To the best of our knowledge, this result is new to the literature. The scaling of the remainder term is the same as in Duy Luu et al (2016) and Rigollet and Tsybakov (2012) with analysis sparsity priors different from ours (the authors in the latter also assume that  $\mathbf{D}$  is invertible).

#### 4.3.4 Anti-sparsity

Recall the derivations for the  $\ell_\infty$  norm example in Section 3.4.1. We have the following oracle inequalities from Proposition 4.

**Corollary 6** *Let the data generated by (30) where  $\boldsymbol{\xi}$  is a zero-mean sub-Gaussian random vector with parameter  $\sigma$ . Assume that  $\mathbf{X}$  is such that*

$\max_{i,j} |\mathbf{X}_{i,j}| \leq 1/p$ . Consider the estimators  $\hat{\boldsymbol{\theta}}_n^{\text{EWA}}$  and  $\hat{\boldsymbol{\theta}}_n^{\text{PEN}}$ , where  $F$  satisfies Assumptions **(H.1)**-**(H.2)**, and  $J$  is the anti-sparsity penalty (14). Suppose that  $\lambda_n \geq \tau\sigma\sqrt{2\delta\log(2)}\sqrt{\frac{p}{n}}$ , for some  $\tau > 1$  and  $\delta > 1$ . Then, with probability at least  $1 - 2^{-p(\delta-1)+1}$ , the following holds

$$\begin{aligned} \frac{1}{n} \|\mathbf{X}\hat{\boldsymbol{\theta}}_n^{\text{EWA}} - \mathbf{X}\boldsymbol{\theta}_0\|_2^2 &\leq \inf_{\substack{I \subset \{1, \dots, p\} \\ \boldsymbol{\theta}: I_{\boldsymbol{\theta}}^{\text{sat}}=I}} \left( \frac{1}{n} \|\mathbf{X}\boldsymbol{\theta} - \mathbf{X}\boldsymbol{\theta}_0\|_2^2 \right. \\ &\quad \left. + \frac{\lambda_n^2(\tau+1)^2}{\tau^2 \mathcal{Y}\left(\{\bar{\boldsymbol{\theta}}: \bar{\boldsymbol{\theta}}_I \in \mathbb{R} \text{sign}(\boldsymbol{\theta}_I)\}, \frac{\tau+1}{\tau-1}\right)^2} \right) + 2p\beta, \end{aligned} \quad (38)$$

and

$$\begin{aligned} \frac{1}{n} \|\mathbf{X}\hat{\boldsymbol{\theta}}_n^{\text{PEN}} - \mathbf{X}\boldsymbol{\theta}_0\|_2^2 &\leq \inf_{\substack{I \subset \{1, \dots, p\} \\ \boldsymbol{\theta}: I_{\boldsymbol{\theta}}^{\text{sat}}=I}} \left( \frac{1}{n} \|\mathbf{X}\boldsymbol{\theta} - \mathbf{X}\boldsymbol{\theta}_0\|_2^2 \right. \\ &\quad \left. + \frac{\lambda_n^2(\tau+1)^2}{\tau^2 \mathcal{Y}\left(\{\bar{\boldsymbol{\theta}}: \bar{\boldsymbol{\theta}}_I \in \mathbb{R} \text{sign}(\boldsymbol{\theta}_I)\}, \frac{\tau+1}{\tau-1}\right)^2} \right). \end{aligned} \quad (39)$$

The first remainder term scales as  $\frac{p}{n}$  which reflects that anti-sparsity regularization requires an overdetermined regime to ensure good stability performance. This is in agreement with (Vaiter et al 2015a, Theorem 7). This phenomenon was also observed by Donoho and Tanner (2010) who studied sample complexity thresholds for noiseless recovery from random projections of the hypercube.

#### 4.3.5 Nuclear norm

We now turn to the nuclear norm case. Recall the notations of Section 2.5.5. For matrices  $\boldsymbol{\theta} \in \mathbb{R}^{p_1 \times p_2}$ , a measurement map  $\mathbf{X}$  takes the form of a linear operator whose  $i$ th component is given by the Frobenius scalar product

$$\mathbf{X}(\boldsymbol{\theta})_i = \text{tr}((\mathbf{X}^i)^\top \boldsymbol{\theta}) = \langle \mathbf{X}^i, \boldsymbol{\theta} \rangle_{\text{F}},$$

where  $\mathbf{X}^i$  is a matrix in  $\mathbb{R}^{p_1 \times p_2}$ . We denote  $\|\cdot\|_{\text{F}}$  the associated norm. From (17), it is immediate to see that whenever  $\boldsymbol{\theta} \neq 0$ ,

$$\mathcal{J}^\circ(e_{\boldsymbol{\theta}}) = \|\|\mathbf{U}\mathbf{V}^\top\|\|_{2 \rightarrow 2} = 1.$$

Moreover, from (17), we have

$$\begin{aligned} \|\mathbf{P}_{T_{\boldsymbol{\theta}}}\|_{\text{F} \rightarrow *} &= \sup_{\boldsymbol{\theta}' \in T_{\boldsymbol{\theta}}} \frac{\|\boldsymbol{\theta}'\|_*}{\|\boldsymbol{\theta}'\|_{\text{F}}} = \sup_{\boldsymbol{\theta}' \in T_{\boldsymbol{\theta}}} \frac{\|\lambda(\boldsymbol{\theta}')\|_1}{\|\lambda(\boldsymbol{\theta}')\|_2} \leq \sup_{\boldsymbol{\theta}' \in T_{\boldsymbol{\theta}}} \sqrt{\text{rank}(\boldsymbol{\theta}')} \\ &\leq \sqrt{\min(r, p_1) + \min(r, p_2)} \leq \sqrt{2r}. \end{aligned}$$



To apply Proposition 2 and Proposition 3, we need to bound  $w(\mathbf{X}(\mathcal{C}))$  ( $\mathcal{C}$  is the nuclear ball), or equivalently, to bound

$$\mathbb{E} [\|\mathbf{X}^*(\mathbf{g})\|_{2 \rightarrow 2}] = \mathbb{E} \left[ \left\| \sum_{i=1}^n \mathbf{X}^i \mathbf{g}_i \right\|_{2 \rightarrow 2} \right], \quad \mathbf{g} \sim \mathcal{N}(0, \sigma^2 \mathbf{Id}_n),$$

which is the expectation of the operator norm of a random series with matrix coefficients. Thus using (Tropp 2015b, Theorem 4.1.1(4.1.5)) to get this bound, and inserting it into Proposition 2 and Proposition 3, we get the following oracle inequalities for the nuclear norm. Define

$$s(\mathbf{X}) = \sqrt{\max \left( \left\| \sum_{i=1}^n \mathbf{X}^i (\mathbf{X}^i)^\top \right\|_{2 \rightarrow 2}, \left\| \sum_{i=1}^n (\mathbf{X}^i)^\top \mathbf{X}^i \right\|_{2 \rightarrow 2} \right) / n}.$$

**Corollary 7** *Let the data generated by (30) with a linear operator  $\mathbf{X} : \mathbb{R}^{p_1 \times p_2} \rightarrow \mathbb{R}^n$ . Consider the estimators  $\hat{\boldsymbol{\theta}}_n^{\text{EWA}}$  and  $\hat{\boldsymbol{\theta}}_n^{\text{PEN}}$ , where  $F$  satisfies Assumptions (H.1)-(H.2), and  $J$  is the nuclear norm (16).*

- (i)  $\boldsymbol{\xi}$  is a zero-mean sub-Gaussian random vector with parameter  $\sigma$ : suppose that  $\lambda_n \geq 2\tau\sigma s(\mathbf{X})c_1 \sqrt{\frac{\log(c_2/\delta) \log(p_1+p_2)}{n}}$ , for some  $\tau > 1$  and  $0 < \delta < \min(c_2, 1)$ , where  $c_1$  and  $c_2$  are the positive absolute constants in Proposition 2. Then, with probability at least  $1 - \delta$ , the following holds

$$\begin{aligned} \frac{1}{n} \|\mathbf{X} \hat{\boldsymbol{\theta}}_n^{\text{EWA}} - \mathbf{X} \boldsymbol{\theta}_0\|_2^2 &\leq \inf_{\substack{r \in \{1, \dots, \min(p_1, p_2)\} \\ \boldsymbol{\theta}: \text{rank}(\boldsymbol{\theta})=r}} \left( \frac{1}{n} \|\mathbf{X} \boldsymbol{\theta} - \mathbf{X} \boldsymbol{\theta}_0\|_2^2 \right. \\ &\quad \left. + \frac{2\lambda_n^2 (\tau+1)^2 r}{\tau^2 \mathcal{T} \left( T_{\boldsymbol{\theta}}, \frac{\tau+1}{\tau-1} \right)^2} \right) + 2p_1 p_2 \beta, \end{aligned} \quad (40)$$

and

$$\begin{aligned} \frac{1}{n} \|\mathbf{X} \hat{\boldsymbol{\theta}}_n^{\text{PEN}} - \mathbf{X} \boldsymbol{\theta}_0\|_2^2 &\leq \inf_{\substack{r \in \{1, \dots, \min(p_1, p_2)\} \\ \boldsymbol{\theta}: \text{rank}(\boldsymbol{\theta})=r}} \left( \frac{1}{n} \|\mathbf{X} \boldsymbol{\theta} - \mathbf{X} \boldsymbol{\theta}_0\|_2^2 \right. \\ &\quad \left. + \frac{2\lambda_n^2 (\tau+1)^2 r}{\tau^2 \mathcal{T} \left( T_{\boldsymbol{\theta}}, \frac{\tau+1}{\tau-1} \right)^2} \right). \end{aligned} \quad (41)$$

- (ii)  $\boldsymbol{\xi} \sim \mathcal{N}(0, \sigma^2 \mathbf{Id}_n)$ : suppose that  $\lambda_n \geq (1 + \delta)\tau\sigma s(\mathbf{X}) \sqrt{\frac{2 \log(p_1+p_2)}{n}}$ , for some  $\tau > 1$  and  $\delta > 0$ . Then, with probability at least  $1 - (p_1 + p_2)^{-\delta^2}$ , (40) and (41) hold.

The set over which the infimum is taken just reminds us that the nuclear norm is partly smooth (see above) relative to the constant rank manifold

(which is a Riemannian submanifold of  $\mathbb{R}^{p_1 \times p_2}$ ) (Daniilidis et al 2014, Theorem 3.19). In the iid Gaussian noise case, we recover the same rate as in (Dalalyan et al 2018, Theorem 3) for  $\hat{\boldsymbol{\theta}}_n^{\text{EWA}}$  and in (Koltchinskii et al 2011, Theorem 2) for  $\hat{\boldsymbol{\theta}}_n^{\text{PEN}}$ . If  $s(\mathbf{X}) = O(\sqrt{p_1 + p_2})$ , then the first remainder term scales as  $\frac{r(p_1 + p_2) \log(p_1 + p_2)}{n}$ . For low-rank matrix recovery, the same rate was also independently proved in Mai and Alquier (2015); Suzuki (2015)<sup>9</sup> for EWA and the posterior conditional mean respectively, in the temperature regime  $\beta = C/n$ , though with completely different priors, but without requiring the compatibility factor assumption.

The assumption  $s(\mathbf{X}) = O(\sqrt{p_1 + p_2})$  on the design is mild and verified in many situations. Indeed, by Jensen's inequality we have

$$s(\mathbf{X})^2 \leq n^{-1} \sum_{i=1}^n \max \left( \left\| \mathbf{X}^i (\mathbf{X}^i)^\top \right\|_{2 \rightarrow 2}, \left\| (\mathbf{X}^i)^\top \mathbf{X}^i \right\|_{2 \rightarrow 2} \right) \leq \left\| \mathbf{X}^i \right\|_{2 \rightarrow 2}^2.$$

If, for example,  $(\mathbf{X}^i)_i$  are independent copies of a standard random Gaussian matrix, then classical concentration bounds of the largest eigenvalue of a Wishart matrix entail that  $\left\| \mathbf{X}^i \right\|_{2 \rightarrow 2}$  concentrates around its mean  $\mathbb{E} \left[ \left\| \mathbf{X}^i \right\|_{2 \rightarrow 2} \right] \leq \sqrt{p_1} + \sqrt{p_2} \leq \sqrt{2(p_1 + p_2)}$ .

#### 4.4 Discussion of minimax optimality

In this section, we discuss the optimality of the estimators  $\hat{\boldsymbol{\theta}}_n^{\text{EWA}}$  and  $\hat{\boldsymbol{\theta}}_n^{\text{PEN}}$  (we remind the reader that the design  $\mathbf{X}$  is fixed). Recall the discussion on stratification at the end of Section 3.1. Let  $\mathcal{M}_0 \in \mathcal{M}$  be the stratum active at  $\boldsymbol{\theta}_0 \in \mathcal{M}_0$ . In this setting, choosing  $\beta = C(1 + \delta)^2 \sigma^2 w(\mathbf{X}(C))^2 \left\| \mathbb{P}_{T_{\boldsymbol{\theta}_0}} \right\|_{2 \rightarrow J}^2 / (pn^2)$  for some constant  $C > 0$ , (22) and Proposition 3 ensure that

$$\begin{aligned} & \frac{1}{n} \left\| \mathbf{X} \hat{\boldsymbol{\theta}}_n^{\text{EWA}} - \mathbf{X} \boldsymbol{\theta}_0 \right\|_2^2 \\ & \leq \frac{(1 + \delta)^2 \sigma^2 w(\mathbf{X}(C))^2 \left\| \mathbb{P}_{T_{\boldsymbol{\theta}_0}} \right\|_{2 \rightarrow J}^2}{n^2} \left( \frac{(\tau J^\circ(e_{\boldsymbol{\theta}_0}) + 1)^2}{\Upsilon\left(T_{\boldsymbol{\theta}_0}, \frac{\tau J^\circ(e_{\boldsymbol{\theta}_0}) + 1}{\tau - 1}\right)^2} + C \right) \\ & \frac{1}{n} \left\| \mathbf{X} \hat{\boldsymbol{\theta}}_n^{\text{PEN}} - \mathbf{X} \boldsymbol{\theta}_0 \right\|_2^2 \\ & \leq \frac{(1 + \delta)^2 \sigma^2 w(\mathbf{X}(C))^2 \left\| \mathbb{P}_{T_{\boldsymbol{\theta}_0}} \right\|_{2 \rightarrow J}^2}{n^2} \frac{(\tau J^\circ(e_{\boldsymbol{\theta}_0}) + 1)^2}{\Upsilon\left(T_{\boldsymbol{\theta}_0}, \frac{\tau J^\circ(e_{\boldsymbol{\theta}_0}) + 1}{\tau - 1}\right)^2}, \end{aligned}$$

with high probability. In particular, for a polyhedral gauge penalty, in which case  $\mathcal{M}_0 = T_{\boldsymbol{\theta}_0}$  (see Vaiteer et al (2015a)), and under the normalization

<sup>9</sup> The noise is iid Gaussian in Suzuki (2015) and subexponential in Mai and Alquier (2015). The assumptions on the design are also different.

$\max_{\mathbf{v} \in \mathcal{V}} \|\mathbf{X}\mathbf{v}\|_2 \leq \sqrt{n}$  and with the choice  $\beta = 2C\delta\sigma^2 \|\| \mathbb{P}_{\mathcal{M}_0} \|\|_{2 \rightarrow J}^2 \log(|\mathcal{V}|)/(pn)$ , Proposition 4 entails

$$\begin{aligned} \frac{1}{n} \|\mathbf{X}\hat{\boldsymbol{\theta}}_n^{\text{EWA}} - \mathbf{X}\boldsymbol{\theta}_0\|_2^2 &\leq \frac{2\delta\sigma^2 \|\| \mathbb{P}_{\mathcal{M}_0} \|\|_{2 \rightarrow J}^2 \log(|\mathcal{V}|)}{n} \left( \frac{(\tau J^\circ(e_{\boldsymbol{\theta}_0}) + 1)^2}{\mathcal{Y}\left(\mathcal{M}_0, \frac{\tau J^\circ(e_{\boldsymbol{\theta}_0}) + 1}{\tau - 1}\right)^2} + C \right) \\ \frac{1}{n} \|\mathbf{X}\hat{\boldsymbol{\theta}}_n^{\text{PEN}} - \mathbf{X}\boldsymbol{\theta}_0\|_2^2 &\leq \frac{2\delta\sigma^2 \|\| \mathbb{P}_{\mathcal{M}_0} \|\|_{2 \rightarrow J}^2 \log(|\mathcal{V}|)}{n} \frac{(\tau J^\circ(e_{\boldsymbol{\theta}_0}) + 1)^2}{\mathcal{Y}\left(\mathcal{M}_0, \frac{\tau J^\circ(e_{\boldsymbol{\theta}_0}) + 1}{\tau - 1}\right)^2}, \end{aligned}$$

with high probability. Thus the risk bounds only depend on  $\mathcal{M}_0$ . A natural question that arises is whether the above bounds are optimal, i.e. whether an estimator can achieve a significantly better prediction risk than  $\hat{\boldsymbol{\theta}}_n^{\text{EWA}}$  and  $\hat{\boldsymbol{\theta}}_n^{\text{PEN}}$  uniformly on  $\mathcal{M}_0$ . A classical way to answer this question is the minimax point of view. This amounts to finding a lower bound on the minimax probabilities of the form

$$\inf_{\hat{\boldsymbol{\theta}}} \sup_{\boldsymbol{\theta} \in \mathcal{M}_0} \Pr\left(\frac{1}{n} \|\mathbf{X}\hat{\boldsymbol{\theta}} - \mathbf{X}\boldsymbol{\theta}\|_2^2 \geq \psi_n\right),$$

where  $\psi_n$  is the rate, which ideally, should be comparable to the risk bounds above. A standard path to derive such a lower bound is to exhibit a subset of  $\mathcal{M}_0$  of well-separated points while controlling its diameter, see (Tsybakov 2008, Chapter 2) or (Massart 2007, Section 4.3). This however must be worked out on a case-by-case basis.

*Example 2 (Lasso)* In this case,  $\mathcal{M}_0 = T_{\boldsymbol{\theta}_0}$  is the subspace of vectors whose support is contained in that of  $\boldsymbol{\theta}_0$ . Let  $I = \text{supp}(\boldsymbol{\theta}_0)$  and  $s = \|\boldsymbol{\theta}_0\|_0$ . Define the set

$$\mathcal{B}_0 = \{\boldsymbol{\theta} \in \mathbb{R}^p : \boldsymbol{\theta}_I \in \{0, 1\}^s \text{ and } \boldsymbol{\theta}_{I^c} = 0\}.$$

We have  $\mathcal{B}_0 \subset \mathcal{M}_0$  and  $\|\boldsymbol{\theta} - \boldsymbol{\theta}'\|_0 \leq 2s$  for all  $(\boldsymbol{\theta}, \boldsymbol{\theta}') \in \mathcal{B}_0$ . Define  $\mathcal{F}_0 \stackrel{\text{def}}{=} \{r\mathbf{X}\boldsymbol{\theta} : \boldsymbol{\theta} \in \mathcal{B}_0\}$ , for  $r > 0$  to be specified later. Due to the Varshamov-Gilbert lemma (Massart 2007, Lemma 4.7), given  $a \in ]0, 1[$ , there exists a subset  $\mathcal{B} \subset \mathcal{B}_0$  with cardinality  $|\mathcal{B}| \geq 2^{\rho s/2}$  such that for two distinct elements  $\mathbf{X}\boldsymbol{\theta}$  and  $\mathbf{X}\boldsymbol{\theta}'$  in  $\mathcal{F}_0$

$$\begin{aligned} \|\mathbf{X}(\boldsymbol{\theta} - \boldsymbol{\theta}')\|_2^2 &\geq \underline{\kappa} r^2 \|\boldsymbol{\theta} - \boldsymbol{\theta}'\|_2^2 \geq 2(1-a)\underline{\kappa} r^2 s, \\ \|\mathbf{X}(\boldsymbol{\theta} - \boldsymbol{\theta}')\|_2^2 &\leq \bar{\kappa} r^2 \|\boldsymbol{\theta} - \boldsymbol{\theta}'\|_2^2 \leq 4\bar{\kappa} r^2 s, \end{aligned}$$

where

$$\underline{\kappa} = \inf_{\boldsymbol{\theta} \in \mathcal{M}_0} \frac{\|\mathbf{X}\boldsymbol{\theta}\|_2^2}{\|\boldsymbol{\theta}\|_2^2} \leq \bar{\kappa} = \sup_{\boldsymbol{\theta} \in \mathcal{M}_0} \frac{\|\mathbf{X}\boldsymbol{\theta}\|_2^2}{\|\boldsymbol{\theta}\|_2^2}.$$

Standard results from random matrix theory, see Tropp (2015a), ensure that  $\underline{\kappa} > 0$  for a Gaussian design with high probability as long as  $n \geq s + C\sqrt{s}$  for some positive absolute constant  $C$ .

Then choosing  $r^2 = \frac{c\rho\sigma^2}{4\bar{\kappa}}$ , where  $c \in ]0, 1/8[$  and  $\rho = (1+a)\log(1+a) + (1-a)\log(1-a)$ , we get the bounds

$$\begin{aligned}\|\mathbf{X}(\boldsymbol{\theta} - \boldsymbol{\theta}')\|_2^2 &\geq \frac{\sigma^2 c(1-a)\rho\bar{\kappa}}{2\bar{\kappa}}s, \\ \|\mathbf{X}(\boldsymbol{\theta} - \boldsymbol{\theta}')\|_2^2 &\leq 2\sigma^2 c \log(|\mathcal{B}|).\end{aligned}$$

We are now in position to apply (Tsybakov 2008, Theorem 2.5) to conclude that there exists  $\eta \in ]0, 1[$  (that depends on  $a$ ) such that

$$\inf_{\hat{\boldsymbol{\theta}}} \sup_{\boldsymbol{\theta} \in \mathcal{M}_0} \Pr\left(\frac{1}{n}\|\mathbf{X}\hat{\boldsymbol{\theta}} - \mathbf{X}\boldsymbol{\theta}\|_2^2 \geq \frac{\sigma^2 c(1-a)\rho\bar{\kappa}}{4\bar{\kappa}} \frac{s}{n}\right) \geq \eta.$$

This lower bound together with Corollary 3 show that  $\hat{\boldsymbol{\theta}}_n^{\text{EWA}}$  (with  $\beta = O(\sigma^2 s \log(2p)/(pn))$ ) and  $\hat{\boldsymbol{\theta}}_n^{\text{PEN}}$  are nearly minimax (up to a logarithmic factor) over  $\mathcal{M}_0$ .

One can generalize this reasoning to get a minimax lower bound over the larger class of  $s$ -sparse vectors, i.e.  $\bigcup\{V = \text{Span}\{(\mathbf{a}_j)_{1 \leq j \leq p}\} : \dim(V) = s\}$ , which is a finite union of subspaces that contains  $\mathcal{M}_0$ . Let  $(a, b) \in ]0, 1]^2$  such that  $1 \leq s \leq abp$  and  $a(-1 + b - \log(b)) \geq \log(2)$ <sup>10</sup>,  $c \in ]0, 1/8[$ . Then combining (Tsybakov 2008, Theorem 2.5) and (Massart 2007, Lemma 4.6 and Lemma 4.10), we have for  $\eta \stackrel{\text{def}}{=} \frac{1}{1+(ab)^{\rho s/2}} \left(1 - 2c - \sqrt{\frac{2c}{-\rho \log(ab)}}\right) \in ]0, 1[$

$$\inf_{\hat{\boldsymbol{\theta}}} \sup_{\boldsymbol{\theta} \in \mathcal{M}_0} \Pr\left(\frac{1}{n}\|\mathbf{X}\hat{\boldsymbol{\theta}} - \mathbf{X}\boldsymbol{\theta}\|_2^2 \geq \frac{\sigma^2 c\rho(1-\alpha)\bar{\kappa}}{2\bar{\kappa}} \frac{s \log(p/s)}{n}\right) \geq \eta,$$

where  $\rho = -a(-1 + b - \log(b))/\log(ab)$ , and  $\underline{\kappa}$  and  $\bar{\kappa}$  are now the restricted isometry constants of  $\mathbf{X}$  of degree  $2s$ , i.e.

$$\underline{\kappa} = \inf_{\|\boldsymbol{\theta}\|_0 \leq 2s} \frac{\|\mathbf{X}\boldsymbol{\theta}\|_2^2}{\|\boldsymbol{\theta}\|_2^2} \leq \bar{\kappa} = \sup_{\|\boldsymbol{\theta}\|_0 \leq 2s} \frac{\|\mathbf{X}\boldsymbol{\theta}\|_2^2}{\|\boldsymbol{\theta}\|_2^2}.$$

For this lower bound to be meaningful,  $\underline{\kappa}$  should be positive. From the compressed sensing literature, many random designs are known to verify this condition for  $n$  large enough compared to  $s$ , e.g. sub-Gaussian designs with  $n \gtrsim s \log(p)$ .

One can see that the difference between this lower bound and the one on  $\mathcal{M}_0$  lies in the  $\log(p/s)$  factor, which basically derives from the control over the union of subspaces. The minimax prediction risk (in expectation) over the  $\ell_0$ -ball were studied in Rigollet and Tsybakov (2011); Raskutti et al (2011); Verzelen (2012); Ye and Zhang (2010); Wang et al (2014), where similar lower bounds were obtained.

<sup>10</sup> E.g. take  $b = 1/(1 + e^{\sqrt[3]{2}})$ .

*Example 3 (Group Lasso)* For the group Lasso with  $L$  groups of equal size  $K$ ,  $\mathcal{M}_0$  is the subspace group sparse vectors whose group support is included in that of  $\boldsymbol{\theta}_0$ . Let  $s$  be the number of non-zero (active) groups in  $\boldsymbol{\theta}_0$ . Following exactly the same reasoning as for the Lasso, one can show that the risk lower bound in probability scales as  $C\sigma^2sK/n$ , which together with Corollary 4, shows that  $\hat{\boldsymbol{\theta}}_n^{\text{EWA}}$  (with  $\beta = O(\sigma^2s(\sqrt{K} + \sqrt{2\log(L)})^2/(pm))$ ) and  $\hat{\boldsymbol{\theta}}_n^{\text{PEN}}$  are nearly minimax (up again to a logarithmic factor) over  $\mathcal{M}_0$ . One can also derive the lower bound  $C\sigma^2s(K + \log(L/s))/n$  over the set of  $s$ -block sparse vectors. Such minimax lower bound is comparable to the one in Lounici et al (2011).

*Example 4 (Anti-sparsity)* Denote the saturation support of  $\boldsymbol{\theta}_0$  as  $I^{\text{sat}}$  and recall the subspace  $T_{\boldsymbol{\theta}_0}$  from (15). Thus,  $\mathcal{M}_0 = T_{\boldsymbol{\theta}_0}$  is the subspace of vectors which are collinear to  $\text{sign}(\boldsymbol{\theta}_0)$  on  $I^{\text{sat}}$  and free on its complement. Observe that  $\dim(\mathcal{M}_0) = p - s + 1$ , where  $s = |I^{\text{sat}}|$ . Define the set

$$\mathcal{B}_0 = \{\boldsymbol{\theta} \in \mathbb{R}^p : \boldsymbol{\theta}_{I^{\text{sat}}} = \text{sign}(\boldsymbol{\theta}_{I^{\text{sat}}}) \text{ and } \boldsymbol{\theta}_{(I^{\text{sat}})^c} \in \{0, 1\}^{p-s}\}.$$

By construction,  $\mathcal{B}_0 \subset \mathcal{M}_0$ , and  $\|\boldsymbol{\theta} - \boldsymbol{\theta}'\|_0 \leq 2(p - s)$  for all  $(\boldsymbol{\theta}, \boldsymbol{\theta}') \in \mathcal{B}_0$ . Thus following the same arguments as for the Lasso example (using again Varshamov-Gilbert lemma and (Tsybakov 2008, Theorem 2.5)), we conclude that there exists  $\eta \in ]0, 1[$  (that depends on  $a$ ) such that

$$\inf_{\hat{\boldsymbol{\theta}}} \sup_{\boldsymbol{\theta} \in \mathcal{M}_0} \Pr\left(\frac{1}{n}\|\mathbf{X}\hat{\boldsymbol{\theta}} - \mathbf{X}\boldsymbol{\theta}\|_2^2 \geq \frac{\sigma^2c(1-a)\rho\kappa p - s}{4\bar{\kappa}n}\right) \geq \eta,$$

where the restricted isometry constants are defined similarly to the Lasso but with respect to the model subspace  $\mathcal{M}_0$  of the  $\ell_\infty$  norm. Again, for a Gaussian design,  $\kappa > 0$  with high probability as long as  $n \geq (p - s + 1) + C\sqrt{p - s + 1}$  Tropp (2015a).

The obtained minimax lower bound is consistent with the sample complexity thresholds derived in Donoho and Tanner (2010) for noiseless recovery from random projections of the hypercube. For a saturation support size  $s$  small compared to  $p$ , the bound of Corollary 6 (with  $\beta = O(\sigma^2/n^2)$ ) comes close to the minimax lower bound.

*Example 5 (Nuclear norm)* Let  $r = \text{rank}(\boldsymbol{\theta}_0)$ , where  $\boldsymbol{\theta}_0 \in \mathbb{R}^{p_1 \times p_2}$ , and  $p = \max(p_1, p_2)$ . For the nuclear norm,  $\mathcal{M}_0$  is the manifold of rank- $r$  matrices. Thus arguing as in (Koltchinskii et al 2011, Theorem 5) (who use the Varshamov-Gilbert lemma Massart (2007) to find the covering set), one can show that the minimax risk lower bound over  $\mathcal{M}_0$  is  $C\sigma^2r/n$ . In view of Corollary 7, we deduce that  $\hat{\boldsymbol{\theta}}_n^{\text{EWA}}$  (with  $\beta = O(\sigma^2r \log(p_1 + p_2)/(p_1p_2n))$ ) and  $\hat{\boldsymbol{\theta}}_n^{\text{PEN}}$  are nearly minimax over the constant rank manifolds.

## A Pre-requisites from convex analysis

We here collect some ingredients from convex analysis that are essential to our exposition.

*Monotone conjugate*

**Lemma 2** *Let  $g$  be a non-decreasing function on  $\mathbb{R}_+$  that vanishes at 0. Then the following hold:*

- (i)  $g^+$  is a proper closed convex and non-decreasing function on  $\mathbb{R}_+$  that vanishes at 0.
- (ii) If  $g$  is also closed and convex, then  $g^{++} = g$ .
- (iii) Let  $f : t \in \mathbb{R} \mapsto g(|t|)$  such that  $f$  is differentiable on  $\mathbb{R}$ , where  $g$  is finite-valued, strictly convex and strongly coercive. Then  $g^+$  is likewise finite-valued, strictly convex, strongly coercive, and  $f^* = g^+ \circ |\cdot|$  is differentiable on  $\mathbb{R}$ . In particular, both  $g$  and  $g^+$  are strictly increasing on  $\mathbb{R}_+$ .

*Proof* (i) By (Bauschke and Combettes 2011, Proposition 13.11),  $g^+$  is a closed convex function. We have  $\inf_{t \geq 0} g(t) = -\sup_{t \geq 0} t \cdot 0 - g(t) = -g^+(0)$ . Since  $g$  is non-decreasing and  $g(0) = 0$ , then  $g^+(0) = -\inf_{t \geq 0} g(t) = -g(0) = 0$ . In addition, by (5), we have  $g^+(a) \geq a \cdot 0 - g(0) = 0, \forall a \in \mathbb{R}_+$ . This shows that  $g^+$  is non-negative and  $\text{dom}(g^+) \neq \emptyset$ , and in turn, it is also proper.

Let  $a, b$  in  $\mathbb{R}_+$  such that  $a < b$ . Then

$$g^+(a) - g^+(b) = \left( \sup_{t \geq 0} ta - g(t) \right) - \left( \sup_{t' \geq 0} t'b - g(t') \right) \leq \sup_{t \geq 0} (ta - g(t) - tb + g(t)) = \sup_{t \geq 0} t(a - b) = 0.$$

That is,  $g^+$  is non-decreasing on  $\mathbb{R}_+$ .

- (ii) This follows from (Rockafellar 1996, Theorem 12.4).
- (iii) By definition of  $f$ ,  $f$  is a finite-valued function on  $\mathbb{R}$ , strictly convex, differentiable and strongly coercive. It then follows from (Hiriart-Urruty and Lemaréchal 2001, Corollary X.4.1.4) that  $f^*$  enjoys the same properties. In turn, using the fact that both  $f$  and  $f^*$  are even, we have  $g^+$  is strongly coercive, and strict convexity of  $f$  (resp.  $f^*$ ) is equivalent to that of  $g$  (resp.  $g^+$ ). Altogether, this shows the first claim. We now prove that  $g$  vanishes only at 0 (and similarly for  $g^+$ ). As  $g$  is non-decreasing and strictly convex, we have, for any  $\rho \in ]0, 1[$  and  $a, b$  in  $\mathbb{R}_+$  such that  $a < b$ ,

$$g(a) \leq g(\rho a + (1 - \rho)b) < \rho g(a) + (1 - \rho)g(b) \leq \rho g(b) + (1 - \rho)g(b) = g(b).$$

*Support function* The support function of  $\mathcal{C} \subset \mathbb{R}^p$  is

$$\sigma_{\mathcal{C}}(\boldsymbol{\omega}) = \sup_{\boldsymbol{\theta} \in \mathcal{C}} \langle \boldsymbol{\omega}, \boldsymbol{\theta} \rangle.$$

We recall the following properties whose proofs can be found in e.g. Rockafellar (1996); Hiriart-Urruty and Lemaréchal (2001).

**Lemma 3** *Let  $\mathcal{C}$  be a non-empty set.*

- (i)  $\sigma_{\mathcal{C}}$  is proper lower semicontinuous (lsc) and sublinear.
- (ii)  $\sigma_{\mathcal{C}}$  is finite-valued if and only if  $\mathcal{C}$  is bounded.
- (iii) If  $0 \in \mathcal{C}$ , then  $\sigma_{\mathcal{C}}$  is non-negative.
- (iv) If  $\mathcal{C}$  is convex and compact with  $0 \in \text{int}(\mathcal{C})$ , then  $\sigma_{\mathcal{C}}$  is finite-valued and coercive.

*Gauges and polars*

**Definition 3 (Polar set)** Let  $\mathcal{C}$  be a nonempty convex set. The set  $\mathcal{C}^\circ$  given by

$$\mathcal{C}^\circ = \left\{ \boldsymbol{\eta} \in \mathbb{R}^p : \langle \boldsymbol{\eta}, \boldsymbol{\theta} \rangle \leq 1 \quad \forall \boldsymbol{\theta} \in \mathcal{C} \right\}$$

is called the polar of  $\mathcal{C}$ .

The set  $\mathcal{C}^\circ$  is closed convex and contains the origin. When  $\mathcal{C}$  is also closed and contains the origin, then it coincides with its bipolar, i.e.  $\mathcal{C}^{\circ\circ} = \mathcal{C}$ .

Let  $\mathcal{C} \subseteq \mathbb{R}^p$  be a non-empty closed convex set containing the origin. The *gauge* of  $\mathcal{C}$  is the function  $\gamma_{\mathcal{C}}$  defined on  $\mathbb{R}^p$  by

$$\gamma_{\mathcal{C}}(\boldsymbol{\theta}) = \inf \{ \lambda > 0 : \boldsymbol{\theta} \in \lambda \mathcal{C} \}.$$

As usual,  $\gamma_{\mathcal{C}}(\boldsymbol{\theta}) = +\infty$  if the infimum is not attained.

Lemma 4 hereafter recaps the main properties of a gauge that we need. In particular, (ii) is a fundamental result of convex analysis that states that there is a one-to-one correspondence between gauge functions and closed convex sets containing the origin. This allows to identify sets from their gauges, and vice versa.

**Lemma 4**

- (i)  $\gamma_{\mathcal{C}}$  is a non-negative, lsc and sublinear function.
- (ii)  $\mathcal{C}$  is the unique closed convex set containing the origin such that

$$\mathcal{C} = \{ \boldsymbol{\theta} \in \mathbb{R}^p : \gamma_{\mathcal{C}}(\boldsymbol{\theta}) \leq 1 \}.$$

- (iii)  $\gamma_{\mathcal{C}}$  is finite-valued if, and only if,  $0 \in \text{int}(\mathcal{C})$ , in which case  $\gamma_{\mathcal{C}}$  is Lipschitz continuous.
- (iv)  $\gamma_{\mathcal{C}}$  is finite-valued and coercive if, and only if,  $\mathcal{C}$  is compact and  $0 \in \text{int}(\mathcal{C})$ .

See [Vaïter et al \(2015a\)](#) for the proof.

Observe that thanks to sublinearity, local Lipschitz continuity valid for any finite-valued convex function is strengthened to global Lipschitz continuity. Moreover,  $\gamma_{\mathcal{C}}$  is a norm, having  $\mathcal{C}$  as its unit ball, if and only if  $\mathcal{C}$  is bounded with nonempty interior and symmetric.

We now define the polar gauge.

**Definition 4 (Polar Gauge)** The polar of a gauge  $\gamma_{\mathcal{C}}$  is the function  $\gamma_{\mathcal{C}}^\circ$  defined by

$$\gamma_{\mathcal{C}}^\circ(\boldsymbol{\omega}) = \inf \{ \mu \geq 0 : \langle \boldsymbol{\theta}, \boldsymbol{\omega} \rangle \leq \mu \gamma_{\mathcal{C}}(\boldsymbol{\theta}), \forall \boldsymbol{\theta} \}.$$

An immediate consequence is that gauges polar to each other have the property

$$\langle \boldsymbol{\theta}, \mathbf{u} \rangle \leq \gamma_{\mathcal{C}}(\boldsymbol{\theta}) \gamma_{\mathcal{C}}^\circ(\mathbf{u}) \quad \forall (\boldsymbol{\theta}, \mathbf{u}) \in \text{dom}(\gamma_{\mathcal{C}}) \times \text{dom}(\gamma_{\mathcal{C}}^\circ), \quad (42)$$

just as dual norms satisfy a duality inequality. In fact, polar pairs of gauges correspond to the best inequalities of this type.

**Lemma 5** Let  $\mathcal{C} \subseteq \mathbb{R}^p$  be a closed convex set containing the origin. Then,

- (ii)  $\gamma_{\mathcal{C}}^\circ$  is a gauge function and  $\gamma_{\mathcal{C}^\circ}^\circ = \gamma_{\mathcal{C}}$ .
- (iii)  $\gamma_{\mathcal{C}}^\circ = \gamma_{\mathcal{C}^\circ}$ , or equivalently

$$\mathcal{C}^\circ = \{ \boldsymbol{\theta} \in \mathbb{R}^p : \gamma_{\mathcal{C}}^\circ(\boldsymbol{\theta}) \leq 1 \}.$$

- (iv) The gauge of  $\mathcal{C}$  and the support function of  $\mathcal{C}$  are mutually polar, i.e.

$$\gamma_{\mathcal{C}} = \sigma_{\mathcal{C}^\circ} \quad \text{and} \quad \gamma_{\mathcal{C}^\circ} = \sigma_{\mathcal{C}}.$$

See [Rockafellar \(1996\)](#); [Hiriart-Urruty and Lemaréchal \(2001\)](#); [Vaïter et al \(2015a\)](#) for the proof.

## B Expectation of the inner product

We start with some definitions and notations that will be used in the proof. For a non-empty closed convex set  $\mathcal{C} \in \mathbb{R}^p$ , we denote  $(\mathcal{C})^0$  its minimal selection, i.e. the element of minimal norm in  $\mathcal{C}$ . This element is of course unique. For a proper lsc and convex function  $f$  and  $\gamma > 0$ , its Moreau envelope (or Moreau-Yosida regularization) is defined by

$$\gamma f(\boldsymbol{\theta}) \stackrel{\text{def}}{=} \min_{\bar{\boldsymbol{\theta}} \in \mathbb{R}^p} \frac{1}{2\gamma} \|\bar{\boldsymbol{\theta}} - \boldsymbol{\theta}\|_2^2 + f(\bar{\boldsymbol{\theta}}).$$

The Moreau envelope enjoys several important properties that we collect in the following lemma.

**Lemma 6** *Let  $f$  be a finite-valued and convex function. Then*

- (i)  $(\gamma f(\boldsymbol{\theta}))_{\gamma > 0}$  is a decreasing net, and  $\forall \boldsymbol{\theta} \in \mathbb{R}^p$ ,  $\gamma f(\boldsymbol{\theta}) \nearrow f(\boldsymbol{\theta})$  as  $\gamma \searrow 0$ .
- (ii)  $\gamma f \in C^1(\mathbb{R}^p)$  with  $\gamma^{-1}$ -Lipschitz continuous gradient.
- (iii)  $\forall \boldsymbol{\theta} \in \mathbb{R}^p$ ,  $\nabla \gamma f(\boldsymbol{\theta}) \rightarrow (\partial f(\boldsymbol{\theta}))^0$  and  $\|\nabla \gamma f(\boldsymbol{\theta})\|_2 \nearrow \|(\partial f(\boldsymbol{\theta}))^0\|_2$  as  $\gamma \searrow 0$ .

*Proof* (i) (Bauschke and Combettes 2011, Proposition 12.32). (ii) (Bauschke and Combettes 2011, Proposition 12.29). (iii) Since  $f$  is finite-valued and convex, it is subdifferentiable everywhere and its subdifferential is a maximal monotone operator with full domain  $\mathbb{R}^p$ , and the result follows from (Bauschke and Combettes 2011, Corollary 23.46(i)).

We are now equipped to prove the following important result<sup>11</sup>. To study EWA, Leung and Barron (2006) used Stein's identity and Catoni (2003, 2007); Dalalyan and Tsybakov (2008); Mai and Alquier (2015) used PAC-Bayesian bounds. Our result hereafter turns out to be instrumental to study EWA in the low-temperature regime for general penalties.

**Proposition 5** *Let the density  $\mu_n$  in (2), where*

- (a)  $F$  satisfies Assumptions (H.1)-(H.2);
- (b)  $J$  is a finite-valued lower-bounded convex function, and  $\exists R > 0$  and  $\rho \geq 0$ , such that  $\forall \boldsymbol{\theta} \in \mathbb{R}^p$ ,  $\|(\partial J(\boldsymbol{\theta}))^0\|_2 \leq R \|\boldsymbol{\theta}\|_2^{\rho}$ ;
- (c) and  $V_n$  is coercive.

Then,  $\forall \bar{\boldsymbol{\theta}} \in \mathbb{R}^p$ ,

$$\mathbb{E}_{\mu_n} \left[ \langle (\partial V_n(\boldsymbol{\theta}))^0, \bar{\boldsymbol{\theta}} - \boldsymbol{\theta} \rangle \right] = -p\beta.$$

This result covers of course the situation where  $J$  fulfills (H.3). In this case, since  $\partial J(\boldsymbol{\theta}) \subset \mathcal{C}^\circ$  by Theorem 1(i), we have  $\rho = 0$  and  $R = \text{diam}(\mathcal{C}^\circ)$ , the diameter of the convex compact set  $\mathcal{C}^\circ$  containing the origin. It can be shown that, when  $F(\cdot, \mathbf{y})$  is strongly coercive, the coercivity assumption (c) can be equivalently stated as  $J_\infty(\boldsymbol{\theta}) > 0$ ,  $\forall \boldsymbol{\theta} \in \ker(\mathbf{X}) \setminus \{0\}$ , where  $J_\infty$  is the recession/asymptotic function of  $J$ ; see e.g. Rockafellar and Wets (1998).

*Proof* Let  $V_n^\gamma(\boldsymbol{\theta}) \stackrel{\text{def}}{=} \frac{1}{n} F(\mathbf{X}\boldsymbol{\theta}, \mathbf{y}) + \lambda_n \gamma J(\boldsymbol{\theta})$  and define  $\mu_n^\gamma(\boldsymbol{\theta}) \stackrel{\text{def}}{=} \exp(-V_n^\gamma(\boldsymbol{\theta})/\beta)/Z$ , where  $0 < Z < +\infty$  is the normalizing constant of the density  $\mu_n$ . Assumption (H.1) and Lemma 6(ii)-(iii) tell us that  $V_n^\gamma \in C^1(\mathbb{R}^p)$  and  $\nabla V_n^\gamma(\boldsymbol{\theta}) \rightarrow (\partial V_n(\boldsymbol{\theta}))^0$  as  $\gamma \rightarrow 0$ . Thus

$$\mathbb{E}_{\mu_n} \left[ \langle (\partial V_n(\boldsymbol{\theta}))^0, \bar{\boldsymbol{\theta}} - \boldsymbol{\theta} \rangle \right] = \int_{\mathbb{R}^p} \lim_{\gamma \rightarrow 0} \langle \mu_n^\gamma(\boldsymbol{\theta}) \nabla V_n^\gamma(\boldsymbol{\theta}), \bar{\boldsymbol{\theta}} - \boldsymbol{\theta} \rangle d\boldsymbol{\theta}.$$

We now check that  $\langle \mu_n^\gamma(\boldsymbol{\theta}) \nabla V_n^\gamma(\boldsymbol{\theta}), \bar{\boldsymbol{\theta}} - \boldsymbol{\theta} \rangle$  is dominated by an integrable function. From the definition of the Moreau envelope, we have

$$V_n^\gamma(\boldsymbol{\theta}) = \min_{\bar{\boldsymbol{\theta}} \in \mathbb{R}^p} \frac{1}{n} F(\mathbf{X}\boldsymbol{\theta}, \mathbf{y}) + \lambda_n (J(\boldsymbol{\theta} - \bar{\boldsymbol{\theta}}) + \frac{1}{2\gamma} \|\bar{\boldsymbol{\theta}}\|_2^2).$$

<sup>11</sup> It will be proved here using Moreau-Yosida regularization. Yet another alternative proof could be based on mollifiers for approximating subdifferentials.



From coercivity of  $V_n$ , the objective in the min is also coercive in  $(\boldsymbol{\theta}, \bar{\boldsymbol{\theta}})$  by (Rockafellar and Wets 1998, Exercise 3.29(b)). It then follows from (Rockafellar and Wets 1998, Theorem 3.31) that  $V_n^\gamma$  is also coercive. In turn, (Rockafellar and Wets 1998, Theorem 11.8(c) and 3.26(a)) allow to assert that for some  $a \in ]0, +\infty[$ ,  $\exists b \in ]-\infty, +\infty[$  such that for all  $\gamma > 0$  and  $\boldsymbol{\theta} \in \mathbb{R}^p$

$$\mu_n^\gamma(\boldsymbol{\theta}) \leq \exp(-a \|\boldsymbol{\theta}\|_2 - b)/Z. \quad (43)$$

Lemma 6-(iii) and assumption (b) on  $J$  entail that for any  $\boldsymbol{\theta} \in \mathbb{R}^p$ ,

$$\|\nabla^\gamma J(\boldsymbol{\theta})\|_2 \leq \|(\partial J(\boldsymbol{\theta}))^0\|_2 \leq R \|\boldsymbol{\theta}\|_2^\rho.$$

Altogether, we have

$$\begin{aligned} |\langle \mu_n^\gamma(\boldsymbol{\theta}) \nabla V_n^\gamma(\boldsymbol{\theta}), \bar{\boldsymbol{\theta}} - \boldsymbol{\theta} \rangle| &\leq \mu_n^\gamma(\boldsymbol{\theta}) \left( |\langle \mathbf{X}^\top \frac{1}{n} \nabla F(\mathbf{X}\boldsymbol{\theta}, \mathbf{y}), \bar{\boldsymbol{\theta}} - \boldsymbol{\theta} \rangle| + \lambda_n \|\nabla^\gamma J(\boldsymbol{\theta})\|_2 \|\bar{\boldsymbol{\theta}} - \boldsymbol{\theta}\|_2 \right) \\ &\leq CZ^{-1} \exp(-F(\mathbf{X}\boldsymbol{\theta}, \mathbf{y})/(n\beta)) \left( |\langle \frac{1}{n} \nabla F(\mathbf{X}\boldsymbol{\theta}, \mathbf{y}), \mathbf{X}(\bar{\boldsymbol{\theta}} - \boldsymbol{\theta}) \rangle| \right. \\ &\quad \left. + (Z \exp b)^{-1} \lambda_n R \exp(-a \|\boldsymbol{\theta}\|_2) \|\boldsymbol{\theta}\|_2^\rho \|\bar{\boldsymbol{\theta}} - \boldsymbol{\theta}\|_2 \right), \end{aligned}$$

where the constant  $C > 0$  reflects the lower-boundedness of  $J$ . It is easy to see that the function in this upper-bound is integrable, where we also use (H.2). Hence, we can apply the dominated convergence theorem to get

$$\mathbb{E}_{\mu_n} \left[ \langle (\partial V_n(\boldsymbol{\theta}))^0, \bar{\boldsymbol{\theta}} - \boldsymbol{\theta} \rangle \right] = \lim_{\gamma \rightarrow 0} \int_{\mathbb{R}^p} \langle \mu_n^\gamma(\boldsymbol{\theta}) \nabla V_n^\gamma(\boldsymbol{\theta}), \bar{\boldsymbol{\theta}} - \boldsymbol{\theta} \rangle d\boldsymbol{\theta}.$$

Now, by simple differential calculus (chain and product rules), we have

$$\begin{aligned} \langle \mu_n^\gamma(\boldsymbol{\theta}) \nabla V_n^\gamma(\boldsymbol{\theta}), \bar{\boldsymbol{\theta}} - \boldsymbol{\theta} \rangle &= -\beta \langle \nabla \mu_n^\gamma(\boldsymbol{\theta}), \bar{\boldsymbol{\theta}} - \boldsymbol{\theta} \rangle \\ &= -\beta \sum_{i=1}^p \frac{\partial}{\partial \theta_i} (\mu_n^\gamma(\boldsymbol{\theta}) (\bar{\theta}_i - \theta_i)) - p\beta \mu_n^\gamma(\boldsymbol{\theta}). \end{aligned}$$

Integrating the first term, we get by Fubini theorem and the Newton-Leibniz formula

$$\begin{aligned} &\int_{\mathbb{R}^{p-1}} \left( \int_{\mathbb{R}} \frac{\partial}{\partial \theta_i} (\mu_n^\gamma(\boldsymbol{\theta}) (\bar{\theta}_i - \theta_i)) d\theta_i \right) d\boldsymbol{\theta}_1 \cdots d\boldsymbol{\theta}_{i-1} d\boldsymbol{\theta}_{i+1} \cdots d\boldsymbol{\theta}_p \\ &= \int_{\mathbb{R}^{p-1}} [\mu_n^\gamma(\boldsymbol{\theta}) (\bar{\theta}_i - \theta_i)]_{\mathbb{R}} d\boldsymbol{\theta}_1 \cdots d\boldsymbol{\theta}_{i-1} d\boldsymbol{\theta}_{i+1} \cdots d\boldsymbol{\theta}_p = 0, \end{aligned}$$

where we used coercivity of  $V_n^\gamma$  (see (43)) to conclude that  $\lim_{|\theta_i| \rightarrow +\infty} \mu_n^\gamma(\boldsymbol{\theta}) (\bar{\theta}_i - \theta_i) = 0$ . For the second term, we have from Lemma 6(i) that  $\mu_n^\gamma \rightarrow \mu_n$  as  $\gamma \rightarrow 0$ . Thus, arguing again as in (43), we can apply the dominated convergence theorem to conclude that

$$\lim_{\gamma \rightarrow 0} \int_{\mathbb{R}^p} \mu_n^\gamma(\boldsymbol{\theta}) d\boldsymbol{\theta} = \int_{\mathbb{R}^p} \mu_n(\boldsymbol{\theta}) d\boldsymbol{\theta} = 1.$$

This concludes the proof.

**Acknowledgements** This work was supported by Conseil Régional de Basse-Normandie and partly by Institut Universitaire de France.

## References

- Bach F (2008) Consistency of the group lasso and multiple kernel learning. *Journal of Machine Learning Research* 9:1179–1225
- Bakin S (1999) Adaptive regression and model selection in data mining problems. Thesis (Ph.D.)—Australian National University, 1999
- Bauschke HH, Combettes PL (2011) *Convex analysis and monotone operator theory in Hilbert spaces*. Springer
- Bellec P (2014) Concentration of quadratic forms under a bernstein moment assumption. Technical report, Ecole Polytechnique
- Bickel PJ, Ritov Y, Tsybakov A (2009) Simultaneous analysis of lasso and Dantzig selector. *The Annals of Statistics* 37(4):1705–1732
- Bogdan M, van den Berg E, Sabatti C, Su W, Candès EJ (2014) Slope – adaptive variable selection via convex optimization. *Annals of Applied Statistics* 9(3):1103–1140
- Bühlmann P, van de Geer S (2011) *Statistics for High-Dimensional Data: Methods, Theory and Applications*. Springer Series in Statistics, Springer-Verlag Berlin Heidelberg
- Candès E, Plan Y (2009) Near-ideal model selection by  $\ell_1$  minimization. *The Annals of Statistics* 37(5A):2145–2177
- Candès EJ, Recht B (2009) Exact matrix completion via convex optimization. *Foundations of Computational mathematics* 9(6):717–772
- Candès EJ, Li X, Ma Y, Wright J (2011) Robust principal component analysis? *Journal of the ACM* 58(3):11:1–11:37
- Candès EJ, Strohmer T, Vershynin V (2013) Phaselift: Exact and stable signal recovery from magnitude measurements via convex programming. *Communications on Pure and Applied Mathematics* 66(8):1241–1274
- Castillo I, Schmidt-Hieber J, van der Vaart A (2015) Bayesian linear regression with sparse priors. *The Annals of Statistics* 43(5):1986–2018
- Catoni O (2003) A PAC-bayesian approach to adaptive classification
- Catoni O (2007) *PAC-Bayesian Supervised Classification (The Thermodynamics of Statistical Learning)*, Lecture Notes-Monograph Series, vol 56. IMS
- Chandrasekaran V, Recht B, Parrilo PA, Willsky A (2012) The convex geometry of linear inverse problems. *Foundations of Computational Mathematics* 12(6):805–849
- Chen S, Donoho D, Saunders M (1999) Atomic decomposition by basis pursuit. *SIAM journal on scientific computing* 20(1):33–61
- Chen X, Lin Q, Kim S, Carbonell JG, Xing EP (2010) An efficient proximal-gradient method for general structured sparse learning. *Preprint arXiv:1005.4717*
- Coste M (2002) An introduction to semialgebraic geometry. Technical report, Institut de Recherche Mathématiques de Rennes
- Dalalyan A, Tsybakov A (2009) pac-bayesian bounds for the expected error of aggregation by exponential weights. Technical report, Université Paris 6, CREST and CERTIS, Ecole des Ponts ParisTech, personal communication
- Dalalyan A, Tsybakov AB (2008) Aggregation by exponential weighting, sharp PAC-bayesian bounds and sparsity. *Machine Learning* 72(1-2):39–61, DOI 10.1007/s10994-008-5051-0, URL <http://dx.doi.org/10.1007/s10994-008-5051-0>
- Dalalyan AS, Tsybakov AB (2007) Aggregation by exponential weighting and sharp oracle inequalities. In: Proceedings of the 20th Annual Conference on Learning Theory, Springer-Verlag, Berlin, Heidelberg, COLT’07, pp 97–111, URL <http://dl.acm.org/citation.cfm?id=1768841.1768854>
- Dalalyan AS, Tsybakov AB (2012) Sparse regression learning by aggregation and langevin monte-carlo. *Journal of Computer and System Sciences* 78(5):1423–1443, DOI 10.1016/j.jcss.2011.12.023, URL <http://dx.doi.org/10.1016/j.jcss.2011.12.023>
- Dalalyan AS, Hebiri M, Lederer J (2017) On the prediction performance of the lasso. *Bernoulli* 23(1):552–581, DOI 10.3150/15-BEJ756, URL <http://dx.doi.org/10.3150/15-BEJ756>
- Dalalyan AS, Grappin E, Paris Q (2018) On the Exponentially Weighted Aggregate with the Laplace Prior. *The Annals of Statistics* 46(5):2452–2478

- Daniilidis A, Drusvyatskiy D, Lewis AS (2014) Orthogonal invariance and identifiability. *SIAM Journal on Matrix Analysis and Applications* 35(2):580–598, DOI 10.1137/130916710
- Donoho D (2006) For most large underdetermined systems of linear equations the minimal  $\ell^1$ -norm solution is also the sparsest solution. *Communications on Pure and Applied Mathematics* 59(6):797–829
- Donoho D, Tanner J (2010) Counting the faces of randomly-projected hypercubes and orthants. *Discrete and Computational Geometry* 43(3):522–541
- Durmus A, Moulines E, Pereyra M (2016) Sampling from convex non continuously differentiable functions, when Moreau meets Langevin, URL <https://hal.archives-ouvertes.fr/hal-01267115>, preprint hal-01267115
- Duy Luu T, Fadili JM, Chesneau C (2016) PAC-Bayesian risk bounds for group-analysis sparse regression by exponential weighting. Technical report, hal-01367742, URL <https://hal.archives-ouvertes.fr/hal-01367742>
- Fadili MJ, Peyré G, Vaïter S, Deledalle C, Salmon J (2013) Stable recovery with analysis decomposable priors. In: Proc. Sampta'13, pp 113–116
- Fazel M, Hindi H, Boyd SP (2001) A rank minimization heuristic with application to minimum order system approximation. In: Proceedings of the American Control Conference, IEEE, vol 6, pp 4734–4739
- van de Geer S (2008) High-dimensional generalized linear models and the Lasso. *The Annals of Statistics* 36:614–645
- van de Geer S (2014) Weakly decomposable regularization penalties and structured sparsity. *Scandinavian Journal of Statistics* 41(1):72–86, DOI 10.1111/sjos.12032, URL <http://dx.doi.org/10.1111/sjos.12032>
- van de Geer S, Bühlmann P (2009) On the conditions used to prove oracle results for the lasso. *Electronic Journal of Statistics* 3:1360–1392, DOI 10.1214/09-EJS506, URL <http://dx.doi.org/10.1214/09-EJS506>
- van de Geer S, Lederer J (2013) The Bernstein–Orlicz norm and deviation inequalities. *Probab Theory Relat Fields* 157:225–250
- Guedj B, Alquier P (2013) PAC-bayesian estimation and prediction in sparse additive models. *Electronic Journal of Statistics* 7:264–291, DOI 10.1214/13-EJS771, URL <http://dx.doi.org/10.1214/13-EJS771>
- Hiriart-Urruty JB, Lemaréchal C (2001) *Convex Analysis And Minimization Algorithms*, vol I and II. Springer
- Jacob L, Obozinski G, Vert JP (2009) Group lasso with overlap and graph lasso. In: Danyluk AP, Bottou L, Littman ML (eds) ICML'09, vol 382, p 55
- Jégou H, Furon T, Fuchs JJ (2012) Anti-sparse coding for approximate nearest neighbor search. In: IEEE ICASSP, pp 2029–2032
- Koltchinskii V (2008) Oracle inequalities in empirical risk minimization and sparse recovery problems. In: Lectures from the 38th Probability Summer School held in Saint-Flour, Lecture Notes in Mathematics, vol 2033, Springer
- Koltchinskii V, Lounici K, Tsybakov AB (2011) Nuclear-norm penalization and optimal rates for noisy low-rank matrix completion. *The Annals of Statistics* 39(5):2302–2329, DOI 10.1214/11-AOS894, URL <http://dx.doi.org/10.1214/11-AOS894>
- Lecué G (2007) Simultaneous adaptation to the margin and to complexity in classification. *The Annals of Statistics* 35(4):1698–1721, DOI 10.1214/009053607000000055, URL <http://dx.doi.org/10.1214/009053607000000055>
- Ledoux M (2001) *The concentration of measure phenomenon*. Mathematical surveys and monographs, American Mathematical Society, Providence (R.I.)
- Leung G, Barron AR (2006) Information theory and mixing least-squares regressions. *IEEE Transactions on Information Theory* 52(8):3396–3410
- Lounici K, Pontil M, van de Geer S, Tsybakov AB (2011) Oracle inequalities and optimal inference under group sparsity. *The Annals of Statistics* 39(4):2164–2204, DOI 10.1214/11-AOS896, URL <http://dx.doi.org/10.1214/11-AOS896>
- Lyubarskii Y, Vershynin R (2010) Uncertainty principles and vector quantization. *IEEE Transactions on Information Theory* 56(7):3491–3501
- Mai TT, Alquier P (2015) A bayesian approach for noisy matrix completion: Optimal rate under general sampling distribution. *Electronic Journal of Statistics* 9(1):823–841, DOI

- 10.1214/15-EJS1020, URL <https://doi.org/10.1214/15-EJS1020>
- Massart P (2007) *Concentration inequalities and model selection*. Ecole d'Été de Probabilités de Saint-Flour XXXIII - 2003, Springer Verlag
- Negahban S, Ravikumar P, Wainwright MJ, Yu B (2012) A unified framework for high-dimensional analysis of M-estimators with decomposable regularizers. *Statistical Science* 27(4):538–557
- Nemirovski A (2000) Topics in non-parametric statistics. In: Emery M, Nemirovski A, Voiculescu D (eds) *Ecole d'Été de Probabilités de Saint-Flour XXVIII-1998*, Springer, Lecture Notes in Mathematics, vol 1738, pp 87–285
- Osborne M, Presnell B, Turlach B (2000) A new approach to variable selection in least squares problems. *IMA journal of numerical analysis* 20(3):389–403
- Peyré G, Fadili J, Chesneau C (2011) Adaptive Structured Block Sparsity Via Dyadic Partitioning. In: EUSIPCO, Barcelona, Spain
- Raskutti G, Wainwright MJ, Yu B (2011) Minimax rates of estimation for high-dimensional linear regression over  $\ell_q$ -balls. *IEEE Transactions on Information Theory* 57(10):6976–6994, DOI 10.1109/TIT.2011.2165799
- Recht B, Fazel M, Parrilo PA (2010) Guaranteed minimum-rank solutions of linear matrix equations via nuclear norm minimization. *SIAM review* 52(3):471–501
- Rigollet P, Tsybakov A (2007) Linear and convex aggregation of density estimators. *Mathematical Methods of Statistics* 16(3):260–280, DOI 10.3103/S1066530707030052, URL <http://dx.doi.org/10.3103/S1066530707030052>
- Rigollet P, Tsybakov A (2011) Exponential screening and optimal rates of sparse estimation. *The Annals of Statistics* 39(2):731–771, DOI 10.1214/10-AOS854, URL <http://dx.doi.org/10.1214/10-AOS854>
- Rigollet P, Tsybakov AB (2012) Sparse estimation by exponential weighting. *Statistical Science* 27(4):558–575, DOI 10.1214/12-STS393, URL <http://dx.doi.org/10.1214/12-STS393>
- Rockafellar R (1996) *Convex analysis*, vol 28. Princeton University Press
- Rockafellar R, Wets R (1998) *Variational analysis*, vol 317. Springer Verlag
- Rudelson M, Vershynin R (2008) On sparse reconstruction from Fourier and Gaussian measurements. *Communications on Pure and Applied Mathematics* 61(8):1025–1045
- Rudin L, Osher S, Fatemi E (1992) Nonlinear total variation based noise removal algorithms. *Physica D: Nonlinear Phenomena* 60(1-4):259–268
- Studer C, Yin W, Baraniuk RG (2012) Signal representations with minimum  $\ell_\infty$ -norm. In: 50th Annual Allerton Conference on Communication, Control, and Computing,
- Su W, Candès EJ (2015) Slope is adaptive to unknown sparsity and asymptotically minimax. *The Annals of Statistics* 44(3):1038–1068
- Sun T, Zhang CH (2012) Scaled sparse linear regression. *Biometrika* 99(4):879, DOI 10.1093/biomet/ass043
- Suzuki T (2015) Convergence rate of bayesian tensor estimator and its minimax optimality. In: Proceedings of the 32nd International Conference on Machine Learning (ICML), vol 37, pp 1273–1282
- Talagrand M (2005) *The generic chaining. Upper and lower bounds of stochastic processes*. Springer-Verlag, Berlin
- Tibshirani R (1996) Regression shrinkage and selection via the Lasso. *Journal of the Royal Statistical Society Series B Methodological* 58(1):267–288
- Tibshirani R, Saunders M, Rosset S, Zhu J, Knight K (2005) Sparsity and smoothness via the fused Lasso. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 67(1):91–108
- Tropp J (2015a) Convex recovery of a structured signal from independent random linear measurements. In: Pfander G (ed) *Sampling Theory, a Renaissance, Applied and Numerical Harmonic Analysis (ANHA)*, Birkhäuser/Springer
- Tropp JA (2015b) An introduction to matrix concentration inequalities. *Foundations and Trends in Machine Learning* 8(1-2):1–230, DOI 10.1561/22000000048, URL <http://dx.doi.org/10.1561/22000000048>
- Tsybakov AB (2008) *Introduction to Nonparametric Estimation*, 1st edn. Springer
- Vaiter S, Golbabaee M, Fadili J, Peyré G (2015a) Model selection with low complexity priors. *Information and Inference: A Journal of the IMA* 4(3):230

- Vaiter S, Peyré G, Fadili MJ (2015b) Low complexity regularization of linear inverse problems. In: Pfander G (ed) Sampling Theory, a Renaissance, Applied and Numerical Harmonic Analysis (ANHA), Birkhäuser/Springer
- Vaiter S, Deledalle C, Fadili MJ, Peyré G, Dossal C (2017) The degrees of freedom of partly smooth regularizers. *Annals of the Institute of Statistical Mathematics* 69(4):791–832
- Vaiter S, Peyré G, Fadili MJ (2018) Model consistency of partly smooth regularizers. *IEEE Transactions on Information Theory* 64(3):1725 – 1737
- Vershynin R (2015) Estimation in high dimension : A geometric perspective. In: Pfander G (ed) Sampling Theory, a Renaissance, Applied and Numerical Harmonic Analysis (ANHA), Birkhäuser/Springer
- Verzelen N (2012) Minimax risks for sparse regressions: Ultra-high dimensional phenomena. *Electronic Journal of Statistics* 6(0):38–90, DOI 10.1214/12-ejs666, URL <http://dx.doi.org/10.1214/12-ejs666>
- Wang Z, Paterlini S, Gao F, Yang Y (2014) Adaptive minimax regression estimation over sparse lq-hulls. *Journal of Machine Learning Research* 15(1):1675–1711, URL <http://dl.acm.org/citation.cfm?id=2627435.2638589>
- Wei F, Huang J (2010) Consistent group selection in high-dimensional linear regression. *Bernoulli* 16(4):1369–1384
- Yang Y (2004) Aggregating regression procedures to improve performance. *Bernoulli* 10(1):25–47
- Ye F, Zhang CH (2010) Rate minimaxity of the lasso and dantzig selector for the lq loss in lr balls. *Journal of Machine Learning Research* 11:3519–3540, URL <http://dl.acm.org/citation.cfm?id=1756006.1953043>
- Yuan M, Lin Y (2006) Model selection and estimation in regression with grouped variables. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 68(1):49–67