



The degrees of freedom of partly smooth regularizers

Samuel Vaïter, Charles-Alban Deledalle, Jalal M. Fadili, Gabriel Peyré,
Charles H Dossal

► To cite this version:

Samuel Vaïter, Charles-Alban Deledalle, Jalal M. Fadili, Gabriel Peyré, Charles H Dossal. The degrees of freedom of partly smooth regularizers . *Annals of the Institute of Statistical Mathematics*, 2017, 69 (4), pp.791 - 832. 10.1007/s10463-016-0563-z . hal-00981634v4

HAL Id: hal-00981634

<https://hal.science/hal-00981634v4>

Submitted on 10 Feb 2016

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.

The Degrees of Freedom of Partly Smooth Regularizers

Samuel Vaïter · Charles Deledalle ·
Jalal Fadili · Gabriel Peyré ·
Charles Dossal

Abstract In this paper, we are concerned with regularized regression problems where the prior regularizer is a proper lower semicontinuous and convex function which is also partly smooth relative to a Riemannian submanifold. This encompasses as special cases several known penalties such as the Lasso (ℓ^1 -norm), the group Lasso (ℓ^1 - ℓ^2 -norm), the ℓ^∞ -norm, and the nuclear norm. This also includes so-called analysis-type priors, i.e. composition of the previously mentioned penalties with linear operators, typical examples being the total variation or fused Lasso penalties. We study the sensitivity of *any* regularized minimizer to perturbations of the observations and provide its precise local parameterization. Our main sensitivity analysis result shows that the predictor moves locally stably along the same active submanifold as the observations undergo small perturbations. This local stability is a consequence of the smoothness of the regularizer when restricted to the active submanifold, which in turn plays a pivotal role to get a closed form expression for the variations of the predictor w.r.t. observations. We also show that, for a variety of regularizers, including polyhedral ones or the group Lasso and its analysis counterpart, this divergence formula holds Lebesgue almost everywhere. When the perturbation is random (with an appropriate continuous distribution), this allows us to derive an unbiased estimator of the degrees of freedom and of the risk of the estimator prediction. Our results hold true without requiring the design matrix to be full column rank. They generalize those already known in the literature such as the Lasso problem, the general Lasso problem (analysis ℓ^1 -penalty), or the group Lasso where existing results for the latter assume that the design is full column rank.

Samuel Vaïter, Gabriel Peyré
CEREMADE, CNRS, Université Paris-Dauphine, Place du Maréchal De Lattre De Tassigny,
75775 Paris Cedex 16, France
E-mail: {samuel.vaïter,gabriel.peyre}@ceremade.dauphine.fr

Jalal Fadili
GREYC, CNRS-ENSICAEN-Université de Caen, 6, Bd du Maréchal Juin, 14050 Caen
Cedex, France
E-mail: Jalal.Fadili@greyc.ensicaen.fr

Charles Deledalle, Charles Dossal
IMB, CNRS, Université Bordeaux 1, 351, Cours de la libération, 33405 Talence Cedex,
France
E-mail: {charles.deledalle,charles.dossal}@math.u-bordeaux1.fr

Keywords Degrees of freedom · Partial smoothness · Manifold · Sparsity · Model selection · o-minimal structures · Semi-algebraic sets · Group Lasso · Total variation

1 Introduction

1.1 Regression and Regularization

We consider a model

$$\mathbb{E}(Y|X) = h(X\beta_0), \quad (1)$$

where $Y = (Y_1, \dots, Y_n)$ is the response vector, $\beta_0 \in \mathbb{R}^p$ is the unknown vector of linear regression coefficients, $X \in \mathbb{R}^{n \times p}$ is the fixed design matrix whose columns are the p covariate vectors, and the expectation is taken with respect to some σ -finite measure. h is a known real-valued and smooth function $\mathbb{R}^n \rightarrow \mathbb{R}^n$. The goal is to design an estimator of β_0 and to study its properties. In the sequel, we do not make any specific assumption on the number of observations n with respect to the number of predictors p . Recall that when $n < p$, (1) is underdetermined, whereas when $n \geq p$ and all the columns of X are linearly independent, it is overdetermined.

Many examples fall within the scope of model (1). We here review two of them.

Example 1 (GLM) One naturally thinks of generalized linear models (GLMs) (McCullagh and Nelder 1989) which assume that conditionally on X , Y_i are independent with distribution that belongs to a given (one-parameter) standard exponential family. Recall that the random variable $Z \in \mathbb{R}$ has a distribution in this family if its distribution admits a density with respect to some reference σ -finite measure on \mathbb{R} of the form

$$p(z; \theta) = B(z) \exp(z\theta - \varphi(\theta)), \quad \theta \in \Theta \subseteq \mathbb{R},$$

where Θ is the natural parameter space and θ is the canonical parameter. For model (1), the distribution of Y belongs to the n -parameter exponential family and its density reads

$$f(y|X; \beta_0) = \left(\prod_{i=1}^n B_i(y_i) \right) \exp \left(\langle y, X\beta_0 \rangle - \sum_{i=1}^n \varphi_i((X\beta_0)_i) \right), \quad X\beta_0 \in \Theta^n, \quad (2)$$

where $\langle \cdot, \cdot \rangle$ is the inner product, and the canonical parameter vector is the linear predictor $X\beta_0$. In this case, $h(\mu) = (h_i(\mu_i))_{1 \leq i \leq n}$, where h_i is the *inverse* of the link function in the language of GLM. Each h_i is a monotonic differentiable function, and a typical choice is the canonical link $h_i = \varphi'_i$, where φ'_i is one-to-one if the family is regular (Brown 1986).

Example 2 (Transformations) The second example is where h plays the role of a transformation such as variance-stabilizing transformations (VSTs), symmetrizing transformations, or bias-corrected transformations. There is an enormous body of literature on transformations, going back to the early 1940s. A typical example is when Y_i are independent Poisson random variables $\sim \mathcal{P}((X\beta_0)_i)$, in which case h_i takes the form of the Anscombe bias-corrected VST. See (DasGupta 2008, Chapter 4) for a comprehensive treatment and more examples.

1.2 Variational 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 (Tikhonov and Arsenin 1997), to non-Hilbertian norms that have sparked considerable interest in the recent years.

Given observations (y_1, \dots, y_n) , we consider the class of estimators obtained by solving the convex optimization problem

$$\hat{\beta}(y) \in \underset{\beta \in \mathbb{R}^p}{\operatorname{Argmin}} F(\beta, y) + J(\beta) . \quad (\mathcal{P}(y))$$

The fidelity term F is of the following form

$$F(\beta, y) = F_0(X\beta, y) \quad (3)$$

where $F_0(\cdot, y)$ is a general loss function assumed to be a proper, convex and sufficiently smooth function of its first argument; see Section 3 for a detailed exposition of the smoothness assumptions. The regularizing penalty J is proper lower semicontinuous and convex, and promotes some specific notion of simplicity/low-complexity on $\hat{\beta}(y)$; see Section 3 for a precise description of the class of regularizing penalties J that we consider in this paper. The type of convex optimization problem in $(\mathcal{P}(y))$ is referred to as a regularized M -estimator in Negahban et al (2012), where J is moreover assumed to have a special decomposability property.

We now provide some illustrative examples of loss functions F and regularizing penalty J routinely used in signal processing, imaging sciences and statistical machine learning.

Example 3 (Generalized linear models) Generalized linear models in the exponential family falls into the class of losses we consider. Indeed, taking the negative log-likelihood corresponding to (2) gives¹

$$F_0(\mu, y) = \sum_{i=1}^n \varphi_i(\mu_i) - \langle y, \mu \rangle . \quad (4)$$

¹ Strictly speaking, the minimization may have to be over a convex subset of \mathbb{R}^p .

It is well-known that if the exponential family is regular, then φ_i is proper, infinitely differentiable, its Hessian is definite positive, and thus it is strictly convex (Brown 1986). Therefore, $F_0(\cdot, y)$ shares exactly the same properties. We recover the squared loss $F_0(\mu, y) = \frac{1}{2}\|y - \mu\|^2$ for the standard linear models (Gaussian case), and the logistic loss $F_0(\mu, y) = \sum_{i=1}^n \log(1 + \exp(\mu_i)) - \langle y, \mu \rangle$ for logistic regression (Bernoulli case).

GLM estimators with losses (4) and ℓ^1 or $\ell^1 - \ell^2$ (group) penalties have been previously considered and some of their properties studied including in (Bunea 2008; van de Geer 2008; de Geer 2008; Meier et al 2008; Bach 2010; Kakade et al 2010); see also (Bühlmann and van de Geer 2011, Chapter 3, 4 and 6).

Example 4 (Lasso) The Lasso regularization is used to promote the sparsity of the minimizers, see (Chen et al 1999; Tibshirani 1996; Osborne et al 2000; Donoho 2006; Candès and Plan 2009; Bickel et al 2009), and (Bühlmann and van de Geer 2011) for a comprehensive review. It corresponds to choosing J as the ℓ^1 -norm

$$J(\beta) = \|\beta\|_1 = \sum_{i=1}^p |\beta_i|. \quad (5)$$

It is also referred to as ℓ^1 -synthesis in the signal processing community, in contrast to the more general ℓ^1 -analysis sparsity penalty detailed below.

Example 5 (General Lasso) To allow for more general sparsity penalties, it may be desirable to promote sparsity through a linear operator $D = (d_1, \dots, d_q) \in \mathbb{R}^{p \times q}$. This leads to the so-called analysis-type sparsity penalty (a.k.a. general Lasso after Tibshirani and Taylor (2012)) where the ℓ^1 -norm is pre-composed by D^* , hence giving

$$J(\beta) = \|D^* \beta\|_1 = \sum_{j=1}^q |\langle d_j, \beta \rangle|. \quad (6)$$

This of course reduces to the usual lasso penalty (5) when $D = \text{Id}_p$. The penalty (6) encapsulates several important penalties including that of the 1-D total variation (Rudin et al 1992), and the fused Lasso (Tibshirani et al 2005). In the former, D^* is a finite difference approximation of the derivative, and in the latter, D^* is the concatenation of the identity matrix Id_p and the finite difference matrix to promote both the sparsity of the vector and that of its variations.

Example 6 (ℓ^∞ Anti-sparsity) In some cases, the vector to be reconstructed is expected to be flat. Such a prior can be captured using the ℓ^∞ norm (a.k.a. Tchebycheff norm)

$$J(\beta) = \|\beta\|_\infty = \max_{i \in \{1, \dots, p\}} |\beta_i|. \quad (7)$$

More generally, it is worth mentioning that a finite-valued function J is polyhedral convex (including Lasso, general Lasso, ℓ^∞) if and only if it can be expressed

as $\max_{i \in \{1, \dots, q\}} \langle d_i, \beta \rangle - b_i$, where the vectors d_i define the facets of the sublevel set at 1 of the penalty (Rockafellar 1996). The ℓ^∞ regularization has found applications in computer vision (Jégou et al 2012), vector quantization (Lyubarskii and Vershynin 2010), or wireless network optimization (Studer et al 2012).

Example 7 (Group Lasso) When the covariates are assumed to be clustered in a few active groups/blocks, the group Lasso has been advocated since it promotes sparsity of the 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 reads

$$J(\beta) = \|\beta\|_{1,2} = \sum_{b \in \mathcal{B}} \|\beta_b\|_2. \quad (8)$$

where $\beta_b = (\beta_i)_{i \in b}$ is the sub-vector of β whose entries are indexed by the block $b \in \mathcal{B}$ where \mathcal{B} is a disjoint union of the set of indices i.e. $\bigcup_{b \in \mathcal{B}} = \{1, \dots, p\}$ such that $b, b' \in \mathcal{B}, b \cap b' = \emptyset$. The mixed ℓ^1 - ℓ^2 norm defined in (8) has the attractive property to be invariant under (groupwise) orthogonal transformations.

Example 8 (General 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 group sparsity. Given a collection of linear operators $\{D_b\}_{b \in \mathcal{B}}$, that are not all orthogonal, the analysis group sparsity penalty is

$$J(\beta) = \|D^* \beta\|_{1,2} = \sum_{b \in \mathcal{B}} \|D_b^* \beta\|_2. \quad (9)$$

This encompasses the 2-D isotropic total variation (Rudin et al 1992), where β is a 2-D discretized image, and each $D_b^* \beta \in \mathbb{R}^2$ is a finite difference approximation of the gradient of β at a pixel indexed by b . The overlapping group Lasso (Jacob et al 2009) is also a special case of (9) by taking $D_b^* : \beta \mapsto \beta_b$ to be a block extractor operator (Peyré et al 2011; Chen et al 2010).

Example 9 (Nuclear norm) The natural extension of low-complexity priors to matrix-valued objects $\beta \in \mathbb{R}^{p_1 \times p_2}$ (where $p = p_1 p_2$) is to penalize the singular values of the matrix. Let $U_\beta \in \mathbb{R}^{p_1 \times p_1}$ and $V_\beta \in \mathbb{R}^{p_2 \times p_2}$ be the orthonormal matrices of left and right singular vectors of β , and $\lambda : \mathbb{R}^{p_1 \times p_2} \rightarrow \mathbb{R}^{p_2}$ is the mapping that returns the singular values of β in non-increasing order. If $j \in \Gamma_0(\mathbb{R}^{p_2})$, i.e. convex, lower semi-continuous and proper, is an absolutely permutation-invariant function, then one can consider the penalty $J(\beta) = j(\lambda(\beta))$. This is a so-called spectral function, and moreover, it can be also shown that $J \in \Gamma_0(\mathbb{R}^{p_1 \times p_2})$ (Lewis 2003b). The most popular spectral penalty is the nuclear norm obtained for $j = \|\cdot\|_1$,

$$J(\beta) = \|\beta\|_* = \|\lambda(\beta)\|_1. \quad (10)$$

This penalty is the best convex candidate to enforce a low-rank prior. It has been widely used for various applications, including low rank matrix completion (Recht et al 2010; Candès and Recht 2009), robust PCA (Candès et al 2011), model reduction (Fazel et al 2001), and phase retrieval (Candès et al 2013).

1.3 Sensitivity Analysis

A chief goal of this paper is to investigate the sensitivity of any solution $\hat{\beta}(y)$ to the parameterized problem $(\mathcal{P}(y))$ to (small) perturbations of y . Sensitivity analysis² is a major branch of optimization and optimal control theory. Comprehensive monographs on the subject are (Bonnans and Shapiro 2000; Mordukhovich 1992). The focus of sensitivity analysis is the dependence and the regularity properties of the optimal solution set and the optimal values when the auxiliary parameters (e.g. y here) undergo a perturbation. In its simplest form, sensitivity analysis of first-order optimality conditions, in the parametric form of the Fermat rule, relies on the celebrated implicit function theorem.

The set of regularizers J we consider is that of partly smooth functions relative to a Riemannian submanifold as detailed in Section 3. The notion of partial smoothness was introduced in (Lewis 2003a). This concept, as well as that of identifiable surfaces (Wright 1993), captures essential features of the geometry of non-smoothness which are along the so-called “active/identifiable manifold”. For convex functions, a closely related idea was developed in (Lemaréchal et al 2000). Loosely speaking, a partly smooth function behaves smoothly as we move on the identifiable manifold, and sharply if we move normal to the manifold. In fact, the behaviour of the function and of its minimizers (or critical points) depend essentially on its restriction to this manifold, hence offering a powerful framework for sensitivity analysis theory. In particular, critical points of partly smooth functions move stably on the manifold as the function undergoes small perturbations (Lewis 2003a; Lewis and Zhang 2013).

Getting back to our class of regularizers, the core of our proof strategy relies on the identification of the active manifold associated to a particular minimizer $\hat{\beta}(y)$ of $(\mathcal{P}(y))$. We exhibit explicitly a certain set of observations, denoted \mathcal{H} (see Definition 3), outside which the initial non-smooth optimization $(\mathcal{P}(y))$ boils down locally to a smooth optimization along the active manifold. This part of the proof strategy is in close agreement with the one developed in (Lewis 2003a) for the sensitivity analysis of partly smooth functions. See also (Bolte et al 2011, Theorem 13) for the case of linear optimization over a convex semialgebraic partly smooth feasible set, where the authors proves a sensitivity result with a zero-measure transition space. However, it is important to stress that neither the results of (Lewis 2003a) nor those of (Bolte et al 2011; Drusvyatskiy and Lewis 2011) can be applied straightforwardly in

² The meaning of sensitivity is different here from what is usually intended in statistical sensitivity and uncertainty analysis.

our context for two main reasons (see also Remark 1 for a detailed discussion). In all these works, a non-degeneracy assumption is crucial while it does not necessarily hold in our case, and this is precisely the reason we consider the boundary of the sets $\mathcal{H}_{\mathcal{M}}$ in the definition of the transition set \mathcal{H} . Moreover, in the latter papers, the authors are concerned with a particular type of perturbations (see Remark 1) which does not allow to cover our class of regularized problems except for restrictive cases such as X injective. For our class of problems $(\mathcal{P}(y))$, we were able to go beyond these works by solving additional key challenges that are important in a statistical context, namely: (i) we provide an analytical description of the set \mathcal{H} involving the boundary of $\mathcal{H}_{\mathcal{M}}$, which entails that \mathcal{H} is potentially of dimension strictly less than n , hence of zero Lebesgue measure, as we will show under a mild o-minimality assumption. (ii) we prove a general sensitivity analysis result valid for any proper lower semicontinuous convex partly smooth regularizer J ; (iii) we compute the first-order expansion of $\hat{\beta}(y)$ and provide an analytical form of the weak derivative of $y \mapsto X\hat{\beta}(y)$ valid outside a set involving \mathcal{H} . If this set is of zero-Lebesgue measure, this allows us to get an unbiased estimator of the risk on the prediction $X\hat{\beta}(Y)$.

1.4 Degrees of Freedom and Unbiased Risk Estimation

The degrees of freedom (DOF) of an estimator quantifies the complexity of a statistical modeling procedure (Efron 1986). It is at the heart of several risk estimation procedures and thus allows one to perform parameter selection through risk minimization.

In this section, we will assume that F_0 in (3) is strictly convex, so that the response (or the prediction) $\hat{\mu}(y) = X\hat{\beta}(y)$ is uniquely defined as a single-valued mapping of y (see Lemma 2). That is, it does not depend on a particular choice of solution $\hat{\beta}(y)$ of $(\mathcal{P}(y))$.

Let $\mu_0 = X\beta_0$. Suppose that h in (1) is the identity and that the observations $Y \sim \mathcal{N}(\mu_0, \sigma^2 \text{Id}_n)$. Following (Efron 1986), the DOF is defined as

$$df = \sum_{i=1}^n \frac{\text{cov}(Y_i, \hat{\mu}_i(Y))}{\sigma^2} .$$

The well-known Stein's lemma (Stein 1981) asserts that, if $y \mapsto \hat{\mu}(y)$ is weakly differentiable function (i.e. typically in a Sobolev space over an open subset of \mathbb{R}^n), such that each coordinate $y \mapsto \hat{\mu}_i(y) \in \mathbb{R}$ has an essentially bounded weak derivative³

$$\mathbb{E} \left(\left| \frac{\partial \hat{\mu}_i}{\partial y_i}(Y) \right| \right) < \infty, \quad \forall i ,$$

³ We write the same symbol as for the derivative, and rigorously speaking, this has to be understood to hold Lebesgue-a.e.

then its divergence is an unbiased estimator of its DOF, i.e.

$$\widehat{df} = \text{div}(\widehat{\mu})(Y) \stackrel{\text{def.}}{=} \text{tr}(D\widehat{\mu}(Y)) \quad \text{and} \quad \mathbb{E}(\widehat{df}) = df ,$$

where $D\widehat{\mu}$ is the Jacobian of $y \mapsto \widehat{\mu}(y)$. In turn, this allows to get an unbiased estimator of the prediction risk $\mathbb{E}(\|\widehat{\mu}(Y) - \mu_0\|^2)$ through the SURE (Stein 1981).

Extensions of the SURE to independent variables from an exponential family are considered in (Hudson 1978) for the continuous case, and (Hwang 1982) in the discrete case. Eldar (2009) generalizes the SURE principle to continuous multivariate exponential families.

1.5 Contributions

We consider a large class of losses F_0 , and of regularizing penalties J which are proper, lower semicontinuous, convex and partly smooth functions relative to a Riemannian submanifold, see Section 3. For this class of regularizers and losses, we first establish in Theorem 1 a general sensitivity analysis result, which provides the local parametrization of any solution to $(\mathcal{P}(y))$ as a function of the observation vector y . This is achieved without placing any specific assumption on X , should it be full column rank or not. We then derive an expression of the divergence of the prediction with respect to the observations (Theorem 2) which is valid outside a set of the form $\mathcal{G} \cap \mathcal{H}$, where \mathcal{G} is defined in Section 5. Using tools from o-minimal geometry, we prove that the transition set \mathcal{H} is of Lebesgue measure zero. If \mathcal{G} is also negligible, then the divergence formula is valid Lebesgue-a.e.. In turn, this allows us to get an unbiased estimate of the DOF and of the prediction risk (Theorem 3 and Theorem 4) for model (1) under two scenarios: (i) Lipschitz continuous non-linearity h and an additive i.i.d. Gaussian noise; (ii) GLMs with a continuous exponential family. Our results encompass many previous ones in the literature as special cases (see discussion in the next section). It is important however to mention that though our sensitivity analysis covers the case of the nuclear norm (also known as the trace norm), unbiasedness of the DOF and risk estimates is not guaranteed in general for this regularizer as the restricted positive definiteness assumption (see Section 4) may not hold at any minimizer (see Example 27), and thus \mathcal{G} may not be always negligible.

1.6 Relation to prior works

In the case of standard Lasso (i.e. ℓ^1 penalty (5)) with $Y \sim \mathcal{N}(X\beta_0, \sigma^2 \text{Id}_n)$ and X of full column rank, (Zou et al 2007) showed that the number of nonzero coefficients is an unbiased estimate for the DOF. Their work was generalized in (Dossal et al 2013) to any arbitrary design matrix. Under the same Gaussian linear regression model, unbiased estimators of the DOF for the general Lasso

penalty (6), were given independently in (Tibshirani and Taylor 2012; Vaiter et al 2013).

A formula of an estimate of the DOF for the group Lasso when the design is orthogonal within each group was conjectured in (Yuan and Lin 2006). Kato (2009) studied the DOF of a general shrinkage estimator where the regression coefficients are constrained to a closed convex set \mathcal{C} . His work extends that of (Meyer and Woodroffe 2000) which treats the case where \mathcal{C} is a convex polyhedral cone. When X is full column rank, (Kato 2009) derived a divergence formula under a smoothness condition on the boundary of \mathcal{C} , from which an unbiased estimator of the degrees of freedom was obtained. When specializing to the constrained version of the group Lasso, the author provided an unbiased estimate of the corresponding DOF under the same group-wise orthogonality assumption on X as (Yuan and Lin 2006). Hansen and Sokol (2014) studied the DOF of the metric projection onto a closed set (non-necessarily convex), and gave a precise representation of the bias when the projector is not sufficiently differentiable. An estimate of the DOF for the group Lasso was also given by (Solo and Ulfarsson 2010) using heuristic derivations that are valid only when X is full column rank, though its unbiasedness is not proved.

Vaiter et al (2012) also derived an estimator of the DOF of the group Lasso and proved its unbiasedness when X is full column rank, but without the orthogonality assumption required in (Yuan and Lin 2006; Kato 2009). When specialized to the group Lasso penalty, our results establish that the DOF estimator formula in (Vaiter et al 2012) is still valid while removing the full column rank assumption. This of course allows one to tackle the more challenging rank-deficient or underdetermined case $p > n$.

2 Notations and preliminaries

Vectors and matrices Given a non-empty closed set $\mathcal{C} \subset \mathbb{R}^p$, we denote $P_{\mathcal{C}}$ the orthogonal projection on \mathcal{C} . For a subspace $T \subset \mathbb{R}^p$, we denote

$$\beta_T = P_T \beta \quad \text{and} \quad X_T = X P_T.$$

For a set of indices $I \subset \mathbb{N}^*$, we will denote β_I (resp. X_I) the subvector (resp. submatrix) whose entries (resp. columns) are those of β (resp. of X) indexed by I . For a linear operator A , A^* is its adjoint. For a matrix M , M^\top is its transpose and M^+ its Moore-Penrose pseudo-inverse.

Sets In the following, for a non-empty set $\mathcal{C} \subset \mathbb{R}^p$, we denote $\text{conv } \mathcal{C}$ and $\text{cone } \mathcal{C}$ respectively its convex and conical hulls. $\iota_{\mathcal{C}}$ is the indicator function of \mathcal{C} (takes 0 in \mathcal{C} and $+\infty$ otherwise), and $N_{\mathcal{C}}(\beta)$ is the cone normal to \mathcal{C} at β . For a non-empty convex set \mathcal{C} , its affine hull $\text{aff } \mathcal{C}$ is the smallest affine manifold containing it. It is a translate of $\text{par } \mathcal{C}$, the subspace parallel to \mathcal{C} , i.e. $\text{par } \mathcal{C} = \text{aff } \mathcal{C} - \beta = \mathbb{R}(\mathcal{C} - \mathcal{C})$ for any $\beta \in \mathcal{C}$. The relative interior $\text{ri } \mathcal{C}$ (resp. relative boundary $\text{rbd } \mathcal{C}$) of \mathcal{C} is its interior (resp. boundary) for the topology relative to its affine hull.

Functions For a C^1 vector field $v : y \in \mathbb{R}^n \mapsto v(y)$, $Dv(y)$ denotes its Jacobian at y . For a C^2 smooth function \tilde{f} , $d\tilde{f}(\beta)[\xi] = \langle \nabla \tilde{f}(\beta), \xi \rangle$ is its directional derivative, $\nabla \tilde{f}(\beta)$ is its (Euclidean) gradient and $\nabla^2 \tilde{f}(\beta)$ is its (Euclidean) Hessian at β . For a bivariate function $g : (\beta, y) \in \mathbb{R}^p \times \mathbb{R}^n \rightarrow \mathbb{R}$ that is C^2 with respect to the first variable β , for any y , we will denote $\nabla g(\beta, y)$ and $\nabla^2 g(\beta, y)$ the gradient and Hessian of g at β with respect to the first variable.

A function $f : \beta \in \mathbb{R}^p \mapsto \mathbb{R} \cup \{+\infty\}$ is lower semicontinuous (lsc) if its epigraph is closed. $\Gamma_0(\mathbb{R}^p)$ is the class of convex and lsc functions which are proper (i.e. not everywhere $+\infty$). ∂f is the (set-valued) subdifferential operator of $f \in \Gamma_0(\mathbb{R}^p)$. If f is differentiable at β then $\nabla f(\beta)$ is its unique subgradient, i.e. $\partial f(\beta) = \{\nabla f(\beta)\}$.

Consider a function $J \in \Gamma_0(\mathbb{R}^p)$ such that $\partial J(\beta) \neq \emptyset$. We denote S_β the subspace parallel to $\partial J(\beta)$ and its orthogonal complement T_β , i.e.

$$S_\beta = \text{par}(\partial J(\beta)) \quad \text{and} \quad T_\beta = S_\beta^\perp. \quad (11)$$

We also use the notation

$$e(\beta) = P_{\text{aff}(\partial J(\beta))}(0),$$

i.e. the projection of 0 onto the affine hull of $\partial J(\beta)$.

Differential and Riemannian geometry Let \mathcal{M} be a C^2 -smooth embedded submanifold of \mathbb{R}^p around $\beta^* \in \mathcal{M}$. To lighten notation, henceforth we shall state C^2 -manifold instead of C^2 -smooth embedded submanifold of \mathbb{R}^p . $\mathcal{T}_\beta(\mathcal{M})$ denotes the tangent space to \mathcal{M} at any point $\beta \in \mathcal{M}$ near β^* . The natural embedding of a submanifold \mathcal{M} into \mathbb{R}^p permits to define a Riemannian structure on \mathcal{M} , and we simply say \mathcal{M} is a Riemannian manifold. For a vector $v \in \mathcal{T}_\beta(\mathcal{M})^\perp$, the Weingarten map of \mathcal{M} at β is the operator $\mathfrak{A}_\beta(\cdot, v) : \mathcal{T}_\beta(\mathcal{M}) \rightarrow \mathcal{T}_\beta(\mathcal{M})$ defined as

$$\mathfrak{A}_\beta(\xi, v) = -P_{\mathcal{T}_\beta(\mathcal{M})} dV[\xi]$$

where V is any local extension of v to a normal vector field on \mathcal{M} . The definition is independent of the choice of the extension V , and $\mathfrak{A}_\beta(\cdot, v)$ is a symmetric linear operator which is closely tied to the second fundamental form of \mathcal{M} ; see (Chavel 2006, Proposition II.2.1).

Let f be a real-valued function which is C^2 on \mathcal{M} around β^* . The covariant gradient of f at β is the vector $\nabla_{\mathcal{M}} f(\beta) \in \mathcal{T}_\beta(\mathcal{M})$ such that

$$\langle \nabla_{\mathcal{M}} f(\beta), \xi \rangle = \frac{d}{dt} f(P_{\mathcal{M}}(\beta + t\xi)) \big|_{t=0}, \forall \xi \in \mathcal{T}_\beta(\mathcal{M}).$$

The covariant Hessian of f at β is the symmetric linear mapping $\nabla_{\mathcal{M}}^2 f(\beta)$ from $\mathcal{T}_\beta(\mathcal{M})$ into itself defined as

$$\langle \nabla_{\mathcal{M}}^2 f(\beta) \xi, \xi \rangle = \frac{d^2}{dt^2} f(P_{\mathcal{M}}(\beta + t\xi)) \big|_{t=0}, \forall \xi \in \mathcal{T}_\beta(\mathcal{M}).$$

This definition agrees with the usual definition using geodesics or connections (Miller and Malick 2005). Assume now that \mathcal{M} is a Riemannian embedded

submanifold of \mathbb{R}^p , and that a function f has a smooth restriction on \mathcal{M} . This can be characterized by the existence of a smooth extension (representative) of f , i.e. a smooth function \tilde{f} on \mathbb{R}^p such that \tilde{f} and f agree on \mathcal{M} . Thus, the Riemannian gradient $\nabla_{\mathcal{M}}f(\beta)$ is also given by

$$\nabla_{\mathcal{M}}f(\beta) = P_{\mathcal{T}_{\beta}(\mathcal{M})} \nabla \tilde{f}(\beta) \quad (12)$$

and, $\forall \xi \in \mathcal{T}_{\beta}(\mathcal{M})$, the Riemannian Hessian reads

$$\begin{aligned} \nabla_{\mathcal{M}}^2 f(\beta) \xi &= P_{\mathcal{T}_{\beta}(\mathcal{M})} d(\nabla_{\mathcal{M}}f)(\beta)[\xi] = P_{\mathcal{T}_{\beta}(\mathcal{M})} d\left(\beta \mapsto P_{\mathcal{T}_{\beta}(\mathcal{M})} \nabla \tilde{f}(\beta)\right)[\xi] \\ &= P_{\mathcal{T}_{\beta}(\mathcal{M})} \nabla^2 \tilde{f}(\beta) P_{\mathcal{T}_{\beta}(\mathcal{M})} \xi + \mathfrak{A}_{\beta}(\xi, P_{\mathcal{T}_{\beta}(\mathcal{M})}^{\perp} \nabla \tilde{f}(\beta)) , \end{aligned} \quad (13)$$

where the last equality comes from (Absil et al 2013, Theorem 1). When \mathcal{M} is an affine or linear subspace of \mathbb{R}^p , then obviously $\mathcal{M} = \beta + \mathcal{T}_{\beta}(\mathcal{M})$, and $\mathfrak{A}_{\beta}(\xi, P_{\mathcal{T}_{\beta}(\mathcal{M})}^{\perp} \nabla \tilde{f}(\beta)) = 0$, hence (13) becomes

$$\nabla_{\mathcal{M}}^2 f(\beta) = P_{\mathcal{T}_{\beta}(\mathcal{M})} \nabla^2 \tilde{f}(\beta) P_{\mathcal{T}_{\beta}(\mathcal{M})} . \quad (14)$$

Similarly to the Euclidean case, for a real-valued bivariate function g that is C^2 on \mathcal{M} around the first variable β , for any y , we will denote $\nabla_{\mathcal{M}}g(\beta, y)$ and $\nabla_{\mathcal{M}}^2g(\beta, y)$ the Riemannian gradient and Hessian of g at β with respect to the first variable. See e.g. (Lee 2003; Chavel 2006) for more material on differential and Riemannian manifolds.

3 Partly Smooth Functions

3.1 Partial Smoothness

Toward the goal of studying the sensitivity behaviour of $\hat{\beta}(y)$ and $\hat{\mu}(y)$ with regularizers $J \in \Gamma_0(\mathbb{R}^p)$, we restrict our attention to a subclass of these functions that fulfill some regularity assumptions according to the following definition.

Definition 1 *Let $J \in \Gamma_0(\mathbb{R}^p)$ and a point β such that $\partial J(\beta) \neq \emptyset$. J is said to be partly smooth at β relative to a set $\mathcal{M} \subseteq \mathbb{R}^p$ if*

1. *Smoothness: \mathcal{M} is a C^2 -manifold and J restricted to \mathcal{M} is C^2 around β .*
2. *Sharpness: $\mathcal{T}_{\beta}(\mathcal{M}) = T_{\beta} \stackrel{\text{def.}}{=} \text{par}(\partial J(\beta))^{\perp}$.*
3. *Continuity: The set-valued mapping ∂J is continuous at β relative to \mathcal{M} .*

J is said to be partly smooth relative to the manifold \mathcal{M} if J is partly smooth at each point $\beta \in \mathcal{M}$ relative to \mathcal{M} .

Observe that \mathcal{M} being affine or linear is equivalent to $\mathcal{M} = \beta + T_{\beta}$. A closed convex set \mathcal{C} is partly smooth at a point $\beta \in \mathcal{C}$ relative to a C^2 -manifold \mathcal{M} locally contained in \mathcal{C} if its indicator function $\iota_{\mathcal{C}}$ maintains this property.

Lewis (2003a, Proposition 2.10) allows to prove the following fact (known as local normal sharpness).

Fact 1 *If J is partly smooth at β relative to \mathcal{M} , then all $\beta' \in \mathcal{M}$ near β satisfy*

$$\mathcal{T}_{\beta'}(\mathcal{M}) = T_{\beta'}.$$

In particular, when \mathcal{M} is affine or linear, then

$$\forall \beta' \in \mathcal{M} \text{ near } \beta, \quad T_{\beta'} = T_{\beta}.$$

It can also be shown that the class of partly smooth functions enjoys a powerful calculus. For instance, under mild conditions, it is closed under positive combination, pre-composition by a linear operator and spectral lifting, with closed-form expressions of the resulting partial smoothness manifolds and their tangent spaces, see (Lewis 2003a; Vaïter et al 2014).

It turns out that except the nuclear norm, the regularizing penalties that we exemplified in Section 1 are partly smooth relative to a linear subspace. The nuclear norm is partly smooth relative to the fixed-rank manifold.

Example 10 (Lasso) We denote $(a_i)_{1 \leq i \leq p}$ the canonical basis of \mathbb{R}^p . Then, $J = \|\cdot\|_1$ is partly smooth at β relative to

$$\mathcal{M} = T_{\beta} = \text{Span}\{(a_i)_{i \in \text{supp}(\beta)}\} \quad \text{where} \quad \text{supp}(\beta) \stackrel{\text{def.}}{=} \{i \in \{1, \dots, p\} : \beta_i \neq 0\}.$$

Example 11 (General Lasso) Vaïter et al (2015, Proposition 9) relates the partial smoothness subspace associated to a convex partly smooth regularizer $J \circ D^*$, where D is a linear operator, to that of J . In particular, for $J = \|\cdot\|_1$, $J \circ D^*$ is partly smooth at β relative to

$$\mathcal{M} = T_{\beta} = \text{Ker}(D_{\Lambda^c}^*) \quad \text{where} \quad \Lambda = \text{supp}(D^*\beta).$$

Example 12 (ℓ^∞ Anti-sparsity) It can be readily checked that $J = \|\cdot\|_\infty$ is partly smooth at β relative to

$$\mathcal{M} = T_{\beta} = \{\beta' : \beta'_I \in \mathbb{R} \text{sign}(\beta_I)\} \quad \text{where} \quad I = \{i : \beta_i = \|\beta\|_\infty\}.$$

Example 13 (Group Lasso) The partial smoothness subspace associated to β when the blocks are of size greater than 1 can be defined similarly, but using the notion of block support. Using the block structure \mathcal{B} , one has that the group Lasso regularizer is partly smooth at β relative to

$$\mathcal{M} = T_{\beta} = \text{Span}\{(a_i)_{i \in \text{supp}_{\mathcal{B}}(\beta)}\},$$

where

$$\text{supp}_{\mathcal{B}}(\beta) = \{i \in \{1, \dots, p\} : \exists b \in \mathcal{B}, \beta_b \neq 0 \text{ and } i \in b\}.$$

Example 14 (General Group Lasso) Using again (Vaïter et al 2015, Proposition 9), we can describe the partial smoothness subspace for $J = \|D^* \cdot\|_{\mathcal{B}}$, which reads

$$\mathcal{M} = T_{\beta} = \text{Ker}(D_{\Lambda^c}^*) \quad \text{where} \quad \Lambda = \text{supp}_{\mathcal{B}}(D^*\beta).$$

Example 15 (Nuclear norm) Piecing together (Daniilidis et al 2013, Theorem 3.19) and Example 10, the nuclear norm can be shown to be partly smooth at $\beta \in \mathbb{R}^{p_1 \times p_2}$ relative to the set

$$\mathcal{M} = \{\beta' : \text{rank}(\beta') = r\}, \quad r = \text{rank}(\beta),$$

which is a C^2 -manifold around β of dimension $(p_1 + p_2 - r)r$; see (Lee 2003, Example 8.14).

Example 16 (Indicator function of a partly smooth set \mathcal{C}) Let \mathcal{C} be a closed convex and partly smooth set at $\beta \in \mathcal{C}$ relative to \mathcal{M} . Observe that when $\beta \in \text{ri}\mathcal{C}$, $\mathcal{M} = \mathbb{R}^p$. For $\beta \in \text{rbd}\mathcal{C}$, \mathcal{M} is locally contained in $\text{rbd}\mathcal{C}$.

We now consider an instructive example of a partly smooth function relative to a non-flat active submanifold that will serve as a useful illustration in the rest of the paper.

Example 17 ($J = \max(\|\cdot\| - 1, 0)$) We have $J \in \Gamma_0(\mathbb{R}^p)$ and continuous. It is then differentiable Lebesgue-a.e., except on the unit sphere \mathbb{S}^{p-1} . For β outside \mathbb{S}^{p-1} , J is partly smooth at β relative to \mathbb{R}^p . For $\beta \in \mathbb{S}^{p-1}$, J is partly smooth at β relative to \mathbb{S}^{p-1} . Obviously, \mathbb{S}^{p-1} is a C^2 -smooth manifold.

3.2 Riemannian Gradient and Hessian

We now give expressions of the Riemannian gradient and Hessian for the case of partly smooth functions relative to a C^2 -manifold. This is summarized in the following fact which follows by combining (12), (13), Definition 1 and Daniilidis et al (2009, Proposition 17).

Fact 2 *If J is partly smooth relative at β relative to \mathcal{M} , then for any $\beta' \in \mathcal{M}$ near β*

$$\nabla_{\mathcal{M}} J(\beta') = P_{T_{\beta'}} (\partial J(\beta')) = e(\beta'),$$

and this does not depend on the smooth representation \tilde{J} of J on \mathcal{M} . In turn,

$$\nabla_{\mathcal{M}}^2 J(\beta) = P_{T_{\beta}} \nabla^2 \tilde{J}(\beta) P_{T_{\beta}} + \mathfrak{A}(\cdot, P_{S_{\beta}} \nabla \tilde{J}(\beta)).$$

Let's now exemplify this fact by providing the expressions of the Riemannian Hessian for the examples discussed above.

Example 18 (Polyhedral penalty) Polyhedrality of J implies that it is affine nearby β along the partial smoothness subspace $\mathcal{M} = \beta + T_{\beta}$, and its subdifferential is locally constant nearby β along \mathcal{M} . In turn, the Riemannian Hessian of J vanishes locally, i.e. $\nabla_{\mathcal{M}}^2 J(\beta') = 0$ for all $\beta' \in \mathcal{M}$ near β . Of course, this holds for the Lasso, general Lasso and ℓ^∞ anti-sparsity penalties since they are all polyhedral.

Example 19 (Group Lasso) Using the expression of $\mathcal{M} = T_\beta$ in Example 13, it is straightforward to show that

$$\nabla_{\mathcal{M}}^2 J(\beta) = \delta_\beta \circ Q_{\beta^\perp},$$

where, for $\Lambda = \text{supp}_B(\beta)$,

$$\begin{aligned} \delta_\beta : T_\beta \rightarrow T_\beta, v \mapsto & \begin{cases} v_b / \|\beta_b\| & \text{if } \beta_b \neq 0 \\ 0 & \text{otherwise} \end{cases} \\ \text{and} \\ Q_{\beta^\perp} : T_\beta \rightarrow T_\beta, v \mapsto & \begin{cases} v_b - \frac{\langle \beta_b, v_b \rangle}{\|\beta_b\|^2} \beta_b & \text{if } \beta_b \neq 0 \\ 0 & \text{otherwise} \end{cases}. \end{aligned}$$

Example 20 (General Group Lasso) Applying the chain rule to Example 19, we get

$$\nabla_{\mathcal{M}}^2 J(\beta) = \text{P}_{\text{Ker}(D_{\Lambda^c}^*)} D(\delta_{D^*\beta} \circ Q_{(D^*\beta)^\perp}) D^* \text{P}_{\text{Ker}(D_{\Lambda^c}^*)},$$

where $\Lambda = \text{supp}_B(D^*\beta)$ and the operator $\delta_{D^*\beta} \circ Q_{(D^*\beta)^\perp}$ is defined similarly to Example 19.

Example 21 (Nuclear norm) For $\beta \in \mathbb{R}^{p_1 \times p_2}$ with $\text{rank}(\beta) = r$, let $\beta = U \text{diag}(\lambda(\beta)) V^*$ be a reduced rank- r SVD decomposition, where $U \in \mathbb{R}^{p_1 \times r}$ and $V \in \mathbb{R}^{p_2 \times r}$ have orthonormal columns, and $\lambda(\beta) \in (\mathbb{R}_+ \setminus \{0\})^r$ is the vector of singular values $(\lambda_1(\beta), \dots, \lambda_r(\beta))$ in non-increasing order. From the partial smoothness of the nuclear norm at β (Example 15) and its subdifferential, one can deduce that

$$\begin{aligned} \mathcal{T}_\beta(\mathcal{M}) = T_\beta &= \{UA^* + BV^* : A \in \mathbb{R}^{p_2 \times r}, B \in \mathbb{R}^{p_1 \times r}\} \text{ and} \\ \nabla_{\mathcal{M}} \|\cdot\|_*(\beta) &= e(\beta) = UV^*. \end{aligned} \quad (15)$$

It can be checked that the orthogonal projector on T_β is given by

$$\text{P}_{T_\beta} W = UU^*W + WV V^* - UU^*WV V^*$$

Let $\xi \in T_\beta$ and $W \in S_\beta$. Then, from (Absil et al 2013, Section 4.5), the Weingarten map reads

$$\mathfrak{A}_\beta(\xi, W) = W\xi^* \beta^{+*} + \beta^{+*} \xi^* W \quad \text{where} \quad \beta^{+*} \stackrel{\text{def}}{=} U \text{diag}(\lambda(\beta))^{-1} V^*. \quad (16)$$

In turn, from Fact 2, the Riemannian Hessian of the nuclear norm reads

$$\begin{aligned} \nabla_{\mathcal{M}}^2 \|\cdot\|_*(\beta)(\xi) &= \text{P}_{T_\beta} \widetilde{\nabla^2 \|\cdot\|_*(\beta)} (\text{P}_{T_\beta} \xi) \\ &\quad + \text{P}_{S_\beta} \widetilde{\nabla \|\cdot\|_*(\beta)} \xi^* \beta^{+*} + \beta^{+*} \xi^* \text{P}_{S_\beta} \widetilde{\nabla \|\cdot\|_*(\beta)}, \end{aligned}$$

where $\widetilde{\|\cdot\|_*}$ is any smooth representative of the nuclear norm at β on \mathcal{M} . Owing to the smooth transfer principle (Daniilidis et al 2013, Corollary 2.3),

the nuclear norm has a C^2 -smooth (and even convex) representation on \mathcal{M} near β which is

$$\widetilde{\|\beta'\|_*} = \widetilde{\|\lambda(\beta')\|_1} = \sum_{i=1}^r \lambda_i(\beta').$$

Combining this with (Lewis 1995, Corollary 2.5), we then have $\nabla \widetilde{\|\cdot\|_*}(\beta) = UV^*$, and thus $\mathfrak{A}_\beta(\xi, P_{S_\beta} \nabla \widetilde{\|\cdot\|_*}(\beta)) = 0$, or equivalently,

$$\nabla_{\mathcal{M}}^2 \|\cdot\|_*(\beta)(\xi) = P_{T_\beta} \nabla^2 \widetilde{\|\cdot\|_*}(\beta)(P_{T_\beta} \xi). \quad (17)$$

The expression of the Hessian $\nabla^2 \widetilde{\|\cdot\|_*}(\beta)$ can be obtained from the derivative of UV^* using either (Candès et al 2012, Theorem 4.3) or (Deledalle et al 2012, Theorem 1) when β is full-rank with distinct singular values, or from (Lewis and Sendov 2001, Theorem 3.3) in the case where β is symmetric with possibly repeated eigenvalues.

Example 22 (Indicator function of a partly smooth set \mathcal{C}) Let \mathcal{C} be a closed convex and partly smooth set at $\beta \in \mathcal{C}$ relative to \mathcal{M} . From Example 16, it is then clear that the zero-function is a smooth representative of $\iota_{\mathcal{C}}$ on \mathcal{M} around β . In turn, the Riemannian gradient and Hessian of $\iota_{\mathcal{C}}$ vanish around β on \mathcal{M} .

Example 23 ($J = \max(\|\cdot\| - 1, 0)$) Let $\beta \in \mathbb{S}^{p-1}$. We have $T_\beta = (\mathbb{R}\beta)^\perp$, and the orthogonal projector onto T_β is

$$P_{T_\beta} = \text{Id} - \beta\beta^\top.$$

The Weingarten map then reduces to

$$\mathfrak{A}_\beta(\xi, v) = -\xi \langle \beta, v \rangle, \quad \xi \in T_\beta \text{ and } v \in S_\beta.$$

Moreover, the zero-function is a smooth representative of J on \mathbb{S}^{p-1} . It then follows that $\nabla_{\mathcal{M}}^2 J(\beta) = 0$.

4 Sensitivity Analysis of $\hat{\beta}(y)$

In all the following, we consider the variational regularized problem $(\mathcal{P}(y))$. We recall that $J \in \Gamma_0(\mathbb{R}^p)$ and is partly smooth. We also suppose that the fidelity term fulfills the following conditions:

$$\forall y \in \mathbb{R}^n, \quad F(\cdot, y) \in C^2(\mathbb{R}^p) \quad \text{and} \quad \forall \beta \in \mathbb{R}^p, \quad F(\beta, \cdot) \in C^2(\mathbb{R}^n). \quad (C_F)$$

Combining (13) and the first part of assumption (C_F) , we have for all $y \in \mathbb{R}^n$

$$\nabla_{\mathcal{M}}^2 F(\beta, y)(\beta, y)\xi = P_{T_\beta} \nabla^2 F(\beta, y) P_{T_\beta} + \mathfrak{A}_\beta(\xi, P_{S_\beta} \nabla F(\beta, y)) P_{T_\beta}. \quad (18)$$

When \mathcal{M} is affine or linear, equation (18) becomes

$$\nabla_{\mathcal{M}}^2 F(\beta, y)(\beta, y)\xi = P_{T_\beta} \nabla^2 F(\beta, y) P_{T_\beta}. \quad (19)$$

4.1 Restricted positive definiteness

In this section, we aim at computing the derivative of the (set-valued) map $y \mapsto \widehat{\beta}(y)$ whenever this is possible. The following condition plays a pivotal role in this analysis.

Definition 2 (Restricted Positive Definiteness) *A vector $\beta \in \mathbb{R}^p$ is said to satisfy the restricted positive definiteness condition if, and only if,*

$$\langle (\nabla_{\mathcal{M}}^2 F(\beta, y) + \nabla_{\mathcal{M}}^2 J(\beta))\xi, \xi \rangle > 0 \quad \forall 0 \neq \xi \in T_{\beta}. \quad (\mathcal{C}_{\beta, y})$$

Condition $(\mathcal{C}_{\beta, y})$ has a convenient re-writing in the following case.

Lemma 1 *Let $J \in \Gamma_0(\mathbb{R}^p)$ be partly smooth at $\beta \in \mathbb{R}^p$ relative to \mathcal{M} , and set $T = T_{\beta}$. Assume that $\nabla_{\mathcal{M}}^2 F(\beta, y)$ and $\nabla_{\mathcal{M}}^2 J(\beta)$ are positive semidefinite on T . Then*

$$(\mathcal{C}_{\beta, y}) \text{ holds if and only if } \text{Ker}(\nabla_{\mathcal{M}}^2 F(\beta, y)) \cap \text{Ker}(\nabla_{\mathcal{M}}^2 J(\beta)) \cap T = \{0\}.$$

For instance, the positive semidefiniteness assumption is satisfied when \mathcal{M} is an affine or linear subspace.

When F takes the form (3) with F_0 the squared loss, condition $(\mathcal{C}_{\beta, y})$ can be interpreted as follows in the examples we discussed so far.

Example 24 (Polyhedral penalty) Recall that a polyhedral penalty is partly smooth at β relative to $\mathcal{M} = \beta + T_{\beta}$. Combining this with Example 18, condition $(\mathcal{C}_{\beta, y})$ specializes to

$$\text{Ker}(X_{T_{\beta}}) = \{0\}.$$

Lasso Applying this to the Lasso (see Example 10), $(\mathcal{C}_{\beta, y})$ reads $\text{Ker}(X_{\Lambda}) = \{0\}$, with $\Lambda = \text{supp}(\beta)$. This condition is already known in the literature, see for instance (Dossal et al 2013).

General Lasso In this case, Example 11 entails that $(\mathcal{C}_{\beta, y})$ becomes

$$\text{Ker}(X) \cap \text{Ker}(D_{\Lambda^c}^* \beta) = \{0\}, \quad \text{where } \Lambda = \text{supp}(D^* \beta).$$

This condition was proposed in (Vaïter et al 2013).

Example 25 (Group Lasso) For the case of the group Lasso, by virtue of Lemma 2(ii) and Example 19, one can see that condition $(\mathcal{C}_{\beta, y})$ amounts to assuming that the system $\{X_b \beta_b : b \in \mathcal{B}, \beta_b \neq 0\}$ is linearly independent. This condition appears in (Liu and Zhang 2009) to establish ℓ^2 -consistency of the group Lasso. It goes without saying that condition $(\mathcal{C}_{\beta, y})$ is much weaker than imposing that X_{Λ} is full column rank, which is standard when analyzing the Lasso.

Example 26 (General group Lasso) For the general group Lasso, let $I_\beta = \{i : b_i \in \mathcal{B} \text{ and } D_{b_i}^* \beta \neq 0\}$, i.e. the set indexing the active blocks of $D^* \beta$. Combining Example 14 and Example 20, one has

$$\begin{aligned} \text{Ker}(\nabla_{\mathcal{M}}^2 J(\beta)) \cap \text{Ker}(D_{\Lambda^c}^*) = \\ \{h \in \mathbb{R}^p : D_{b_i}^* h = 0 \ \forall i \notin I_\beta \text{ and } D_{b_i}^* h \in \mathbb{R} D_{b_i}^* \beta \ \forall i \in I_\beta\}, \end{aligned}$$

where $\Lambda = \text{supp}_{\mathcal{B}}(D^* \beta)$. Indeed, $\delta_{D^* \beta}$ is a diagonal strictly positive linear operator, and $Q_{(D^* \beta)^\perp}$ is a block-wise linear orthogonal projector, and we get for $h \in \text{Ker}(D_{\Lambda^c}^*)$,

$$\begin{aligned} h \in \text{Ker}(\nabla_{\mathcal{M}}^2 J(\beta)) &\iff \langle h, \nabla_{\mathcal{M}}^2 J(\beta) h \rangle = 0 \\ &\iff \langle D^* h, (\delta_{D^* \beta} \circ Q_{(D^* \beta)^\perp}) D^* h \rangle = 0 \\ &\iff \sum_{i \in I_\beta} \frac{\|P_{(D_{b_i}^* \beta)^\perp}(D_{b_i}^* h)\|^2}{\|D_{b_i}^* \beta\|} = 0 \\ &\iff D_{b_i}^* \beta \in \mathbb{R} D_{b_i}^* \beta \quad \forall i \in I_\beta. \end{aligned}$$

In turn, by Lemma 2(ii), condition $(\mathcal{C}_{\beta, y})$ is equivalent to saying that 0 is the only vector in the set

$$\{h \in \mathbb{R}^p : Xh = 0 \text{ and } D_{b_i}^* h = 0 \ \forall i \notin I_\beta \text{ and } D_{b_i}^* h \in \mathbb{R} D_{b_i}^* \beta \ \forall i \in I_\beta\}.$$

Observe that when D is a Parseval tight frame, i.e. $DD^* = \text{Id}$, the above condition is also equivalent to saying that the system $\{(XD)_{b_i} D_{b_i}^* \beta : i \in I_\beta\}$ is linearly independent.

Example 27 (Nuclear norm) We have seen in Example 21 that the nuclear norm has a C^2 -smooth representative which is also convex. It then follows from (17) that the Riemannian Hessian of the nuclear norm at β is positive semidefinite on T_β , where T_β is given in (15).

As far as F is concerned, one cannot conclude in general on positive semidefiniteness of its Riemannian Hessian. Let's consider the case where $\beta \in \mathbf{S}^p$, the vector space of real $p_1 \times p_1$ symmetric matrices endowed with the trace (Frobenius) inner product $\langle \beta, \beta' \rangle = \text{tr}(\beta \beta')$. From (16) and (18), we have for any $\xi \in T_\beta \cap \mathbf{S}^{p_1}$

$$\begin{aligned} \langle \xi, \nabla_{\mathcal{M}}^2 F(\beta, y)(\xi) \rangle &= \langle \xi, P_{T_\beta} \nabla^2 F(\beta, y)(P_{T_\beta} \xi) \rangle \\ &\quad + 2 \langle \xi U \text{diag}(\lambda(\beta))^{-1} U^\top \xi, P_{S_\beta} \nabla F(\beta, y) \rangle. \end{aligned}$$

Assume that β is a global minimizer of $(\mathcal{P}(y))$, which by Lemma 3, implies that

$$P_{S_\beta} \nabla F(\beta, y) = U_\perp \text{diag}(\gamma) U_\perp^\top$$

where $U_\perp \in \mathbb{R}^{n \times (p_1 - r)}$ is a matrix whose columns are orthonormal to U , and $\gamma \in [-1, 1]^{p_1 - r}$. We then get

$$\begin{aligned} \langle \xi, \nabla_{\mathcal{M}}^2 F(\beta, y)(\xi) \rangle &= \langle \xi, P_{T_\beta} \nabla^2 F(\beta, y)(P_{T_\beta} \xi) \rangle \\ &\quad + 2 \langle U_\perp^\top \xi U \text{diag}(\lambda(\beta))^{-1} U^\top \xi U_\perp, \text{diag}(\gamma) \rangle. \end{aligned}$$

It is then sufficient that β is such that the entries of γ are positive for $\nabla_{\mathcal{M}}^2 F(\beta, y)$ to be indeed positive semidefinite on T . In this case, Lemma 1 applies.

In a nutshell, Lemma 1 does not always apply to the nuclear norm as $\nabla_{\mathcal{M}}^2 F(\beta, y)$ is not always guaranteed to be positive semidefinite in this case. One may then wonder whether there exist partly smooth functions J , with a non-flat active submanifold, for which Lemma 1 applies, at least at some minimizer of $(\mathcal{P}(y))$. The answer is affirmative for instance for the regularizer of Example 17.

Example 28 ($J = \max(\|\cdot\| - 1, 0)$) Let $\beta \in \mathbb{S}^{p-1}$. From Example 23, we have for $\xi \in T_{\beta}$

$$\langle \xi, \nabla_{\mathcal{M}}^2 F(\beta, y) \xi \rangle = \langle \xi, \nabla^2 F(\beta, y) \xi \rangle - \|\xi\|^2 \langle \beta, \nabla F(\beta, y) \rangle.$$

Assume that β is a global minimizer of $(\mathcal{P}(y))$, which by Lemma 3, implies that

$$-\nabla F(\beta, y) \in \beta[0, 1] \Rightarrow -\langle \beta, \nabla F(\beta, y) \rangle \in [0, 1].$$

Thus, $\langle \xi, \nabla_{\mathcal{M}}^2 F(\beta, y) \xi \rangle \geq 0$, for all $\xi \in T_{\beta}$. Since from Example 23, $\nabla_{\mathcal{M}}^2 J(\beta) = 0$, Lemma 1 applies at β . Condition $(\mathcal{C}_{\beta, y})$ then holds if, and only if, $\nabla_{\mathcal{M}}^2 F(\beta, y)$ is positive definite on T_{β} . For the case of a quadratic loss, this is equivalent to

$$\ker(X) \cap T_{\beta} = \{0\} \quad \text{or} \quad \beta \text{ is not a minimizer of } F(\cdot, y).$$

4.2 Sensitivity analysis: Main result

Let us now turn to the sensitivity of any minimizer $\widehat{\beta}(y)$ of $(\mathcal{P}(y))$ to perturbations of y . Because of non-smoothness of the regularizer J , it is a well-known fact in sensitivity analysis that one cannot hope for a global claim, i.e. an everywhere smooth mapping⁴ $y \mapsto \widehat{\beta}(y)$. Rather, the sensitivity behaviour is local. This is why the reason we need to introduce the following transition space \mathcal{H} , which basically captures points of non-smoothness of $\widehat{\beta}(y)$.

Let's denote the set of all possible partial smoothness active manifolds \mathcal{M}_{β} associated to J as

$$\mathcal{M} = \{\mathcal{M}_{\beta}\}_{\beta \in \mathbb{R}^p}. \quad (20)$$

For any $\mathcal{M} \in \mathcal{M}$, we denote $\widehat{\mathcal{M}}$ the set of vectors sharing the same partial smoothness manifold \mathcal{M} ,

$$\widehat{\mathcal{M}} = \{\beta' \in \mathbb{R}^p : \mathcal{M}_{\beta'} = \mathcal{M}\}.$$

For instance, when $J = \|\cdot\|_1$, $\widehat{\mathcal{M}}_{\beta}$ is the cone of all vectors sharing the same support as β .

⁴ To be understood here as a set-valued mapping.

Definition 3 The transition space \mathcal{H} is defined as

$$\mathcal{H} = \bigcup_{\mathcal{M} \in \mathcal{M}} \mathcal{H}_{\mathcal{M}}, \quad \text{where } \mathcal{H}_{\mathcal{M}} = \text{bd}(\Pi_{n+p,n}(\mathcal{A}_{\mathcal{M}})),$$

where \mathcal{M} is given by (20), and we denote

$$\Pi_{n+p,n} : \begin{cases} \mathbb{R}^n \times \widehat{\mathcal{M}} \longrightarrow \mathbb{R}^n \\ (y, \beta) \longmapsto y \end{cases}$$

the canonical projection on the first n coordinates, $\text{bd}\mathcal{C}$ is the boundary of the set \mathcal{C} , and

$$\mathcal{A}_{\mathcal{M}} = \{(y, \beta) \in \mathbb{R}^n \times \widehat{\mathcal{M}} : -\nabla F(\beta, y) \in \text{rbd } \partial J(\beta)\}.$$

Remark 1 Before stating our result, some comments about this definition are in order. When bd is removed in the definition of $\mathcal{H}_{\mathcal{M}}$, we recover the classical setting of sensitivity analysis under partial smoothness, where $\mathcal{H}_{\mathcal{M}}$ contains the set of degenerate minimizers (those such that 0 is in the relative boundary of the subdifferential of $F(\cdot, y) + J$). This is considered for instance in (Bolte et al 2011; Drusvyatskiy and Lewis 2011) who studied sensitivity of the minimizers of $\beta \mapsto f_{\mathbf{v}}(\beta) \stackrel{\text{def.}}{=} f(\beta) - \langle \mathbf{v}, \beta \rangle$ to perturbations of \mathbf{v} when $f \in \Gamma_0(\mathbb{R}^p)$ and partly smooth; see also (Drusvyatskiy et al 2015) for the semialgebraic non-necessarily non-convex case. These authors showed that for \mathbf{v} outside a set of Lebesgue measure zero, $f_{\mathbf{v}}$ has a non-degenerate minimizer with quadratic growth of $f_{\mathbf{v}}$, and for each $\bar{\mathbf{v}}$ near \mathbf{v} , the perturbed function $f_{\bar{\mathbf{v}}}$ has a unique minimizer that lies on the active manifold of $f_{\mathbf{v}}$ with quadratic growth of $f_{\bar{\mathbf{v}}}$. These results however do not apply to our setting in general. To see this, consider the case of $(\mathcal{P}(y))$ where F takes the form (3) with F_0 the quadratic (the same applies to other losses in the exponential family just as well). Then, $(\mathcal{P}(y))$ is equivalent to minimizing $f_{\mathbf{v}}$, with $f = J + \|X \cdot\|^2$ and $\mathbf{v} = 2X^\top y$. It goes without saying that, in general (i.e. for any X), a property valid for \mathbf{v} outside a zero Lebesgue measure set does **not** imply it holds for y outside a zero Lebesgue measure set. To circumvent such a difficulty, our key contribution is to consider the boundary of $\mathcal{H}_{\mathcal{M}}$. This turns out to be crucial to get a set of dimension potentially strictly less than n , hence negligible, as we will show under a mild o -minimality assumption (see Section 6). However, doing so, uniqueness of the minimizer is not longer guaranteed.

In the particular case of the Lasso (resp. general Lasso), i.e. F_0 is the squared loss, $J = \|\cdot\|_1$ (resp. $J = \|D^* \cdot\|_1$), the transition space \mathcal{H} specializes to the one introduced in (Dossal et al 2013) (resp. (Vaiter et al 2013)). In these specific cases, since J is a polyhedral gauge, \mathcal{H} is in fact a union of affine hyperplanes. The geometry of this set can be significantly more complex for other regularizers. For instance, for $J = \|\cdot\|_{1,2}$, it can be shown to be a semi-algebraic set (union of algebraic hyper-surfaces). Section 6 is devoted to a detailed analysis of this set \mathcal{H} .

We are now equipped to state our main sensitivity analysis result, whose proof is deferred to Section 8.3.

Theorem 1 Assume that (C_F) holds. Let $y \notin \mathcal{H}$, and $\hat{\beta}(y)$ a solution of $(\mathcal{P}(y))$ where $J \in \Gamma_0(\mathbb{R}^p)$ is partly smooth at $\hat{\beta}(y)$ relative to $\mathcal{M} \stackrel{\text{def.}}{=} \mathcal{M}_{\hat{\beta}(y)}$ and such that $(\mathcal{C}_{\hat{\beta}(y),y})$ holds. Then, there exists an open neighborhood $\mathcal{V} \subset \mathbb{R}^n$ of y , and a mapping $\tilde{\beta} : \mathcal{V} \rightarrow \mathcal{M}$ such that

1. For all $\bar{y} \in \mathcal{V}$, $\tilde{\beta}(\bar{y})$ is a solution of $(\mathcal{P}(\bar{y}))$, and $\tilde{\beta}(y) = \hat{\beta}(y)$.
2. the mapping $\tilde{\beta}$ is $C^1(\mathcal{V})$ and

$$\forall \bar{y} \in \mathcal{V}, \quad D\tilde{\beta}(\bar{y}) = -(\nabla_{\mathcal{M}}^2 F(\tilde{\beta}(\bar{y}), \bar{y}) + \nabla_{\mathcal{M}}^2 J(\tilde{\beta}(\bar{y})))^+ P_{T_{\tilde{\beta}(\bar{y})}} D(\nabla F)(\tilde{\beta}(\bar{y}), \bar{y}), \quad (21)$$

where $D(\nabla F)(\beta, y)$ is the Jacobian of $\nabla F(\beta, \cdot)$ with respect to the second variable evaluated at y .

Theorem 1 can be extended to the case where the data fidelity is of the form $F(\beta, \theta)$ for some parameter θ , with no particular role of y here.

5 Sensitivity Analysis of $\hat{\mu}(y)$

We assume in this section that F takes the form (3) with

$$\forall (\mu, y) \in \mathbb{R}^n \times \mathbb{R}^n, \quad \nabla^2 F_0(\mu, y) \text{ is positive definite.} \quad (C_{dp})$$

This in turn implies that $F_0(\cdot, y)$ is strictly convex for any y (the converse is obviously not true). Recall that this condition is mild and holds in many situations, in particular for some losses (4) in the exponential family, see Section 1.2 for details.

We have the following simple lemma.

Lemma 2 Suppose the condition (C_{dp}) is satisfied. The following holds,

- (i) All minimizers of $(\mathcal{P}(y))$ share the same image under X and J .
- (ii) If the partial smoothness submanifold \mathcal{M} at β is affine or linear, then $(\mathcal{C}_{\beta,y})$ holds if, and only if, $\text{Ker}(X) \cap \text{Ker}(\nabla_{\mathcal{M}}^2 J(\beta)) \cap T = \{0\}$, where $T = T_{\beta}$ and $\nabla_{\mathcal{M}}^2 J(\beta)$ is given in Fact 2.

Owing to this lemma, we can now define the prediction

$$\hat{\mu}(y) = X\hat{\beta}(y) \quad (22)$$

without ambiguity given any solution $\hat{\beta}(y)$, which in turn defines a single-valued mapping $\hat{\mu}$. The following theorem provides a closed-form expression of the local variations of $\hat{\mu}$ as a function of perturbations of y . For this, we define the following set that rules out the points y where $(\mathcal{C}_{\hat{\beta}(y),y})$ does not hold for any any minimizer.

Definition 4 (Non-injectivity set) The Non-injectivity set \mathcal{G} is

$$\mathcal{G} = \left\{ y \notin \mathcal{H} : (\mathcal{C}_{\hat{\beta}(y),y}) \text{ does not hold for any minimizer } \hat{\beta}(y) \text{ of } (\mathcal{P}(y)) \right\}.$$

Theorem 2 *Under assumptions (C_F) and (C_{dp}) , the mapping $y \mapsto \hat{\mu}(y)$ is $C^1(\mathbb{R}^n \setminus (\mathcal{H} \cup \mathcal{G}))$. Moreover, for all $y \notin \mathcal{H} \cup \mathcal{G}$,*

$$\operatorname{div}(\hat{\mu})(y) \stackrel{\text{def.}}{=} \operatorname{tr}(D\hat{\mu}(y)) = \operatorname{tr}(\Delta(y)) \quad (23)$$

where

$$\begin{aligned} \Delta(y) &= -X_T (\nabla_{\mathcal{M}}^2 F(\hat{\mu}(y), y) + \nabla_{\mathcal{M}}^2 J(\hat{\beta}(y)))^+ X_T^\top D(\nabla F_0)(\hat{\mu}(y), y), \\ \nabla_{\mathcal{M}}^2 F(\hat{\mu}(y), y) &= X_T^\top \nabla^2 F_0(\hat{\mu}(y), y) X_T + \mathfrak{A}_\beta(\cdot, X_S^\top \nabla F_0(\hat{\mu}(y), y)) \end{aligned}$$

and $\hat{\beta}(y)$ is any solution of $(\mathcal{P}(y))$ such that $(\mathcal{C}_{\hat{\beta}(y), y})$ holds and $J \in \Gamma_0(\mathbb{R}^p)$ is partly smooth at $\hat{\beta}(y)$ relative to \mathcal{M} , with $T = S^\perp = T_{\hat{\beta}(y)}$.

This result is proved in Section 8.5.

A natural question that arises is whether the set \mathcal{G} is of full (Hausdorff) dimension or not, and in particular, whether there always exists a solution $\hat{\beta}(y)$ such that $(\mathcal{C}_{\hat{\beta}(y), y})$ holds, i.e. \mathcal{G} is empty. Though we cannot provide an affirmative answer to this for any partly smooth regularizer, and this has to be checked on a case-by-case basis, it turns out that \mathcal{G} is indeed empty for many regularizers of interest as established in the next result.

Proposition 1 *The set \mathcal{G} is empty when:*

- (i) $J \in \Gamma_0(\mathbb{R}^p)$ is polyhedral, and in particular, when J is the Lasso, the general Lasso or the ℓ^∞ penalties.
- (ii) J is the general group Lasso penalty, and a fortiori the group Lasso.

The proof of these results is constructive.

We now exemplify the divergence formula (23) when F_0 is the squared loss.

Example 29 (Polyhedral penalty) Thanks to Example 18, it is immediate to see that (23) boils down to

$$\operatorname{div}(\hat{\mu})(y) = \operatorname{rank} X_{T_{\hat{\beta}(y)}} = \dim T_{\hat{\beta}(y)}$$

where we used the rank-nullity theorem and that Lemma 2(ii) holds at $\hat{\beta}(y)$, which always exists by Proposition 1.

Example 30 (Lasso and General Lasso) Combining together Example 11 and Example 29 yields

$$\operatorname{div}(\hat{\mu})(y) = \dim \operatorname{Ker}(D_{\Lambda^c}^*), \quad \Lambda = \operatorname{supp}(D^* \hat{\beta}(y)),$$

where $\hat{\beta}(y)$ is such that Lemma 2(ii) holds. For the Lasso, Example 10 allows to specialize the formula to

$$\operatorname{div}(\hat{\mu})(y) = |\operatorname{supp}(\hat{\beta}(y))|.$$

The general Lasso case was investigated in (Vaiter et al 2013) and (Tibshirani and Taylor 2012), and the Lasso in (Dossal et al 2013) and (Tibshirani and Taylor 2012).

Example 31 (ℓ^∞ Anti-sparsity) By virtue of Example 29 and Example 12, we obtain in this case

$$\operatorname{div}(\hat{\mu})(y) = p - |I| + 1, \quad \text{where } I = \{i : \hat{\beta}_i(y) = \|\hat{\beta}(y)\|_\infty\}$$

and $\hat{\beta}(y)$ is such that Lemma 2(ii) holds, and such a vector always exists by Proposition 1.

Example 32 (Group Lasso and General Group Lasso) For the general group Lasso, piecing together Example 14 and Example 20, it follows that

$$\operatorname{div}(\hat{\mu})(y) = \operatorname{tr} \left(X_T \left(X_T^\top X_T + P_T D \left(\delta_{D^* \hat{\beta}(y)} \circ Q_{(D^* \hat{\beta}(y))^\perp} \right) D^* P_T \right)^+ X_T^\top \right)$$

where $T = \operatorname{Ker}(D_{\Lambda^c}^*)$, $\Lambda = \operatorname{supp}_{\mathcal{B}}(D^* \hat{\beta}(y))$, and $\hat{\beta}(y)$ is such that Lemma 2(ii) holds; such a vector always exists by Proposition 1. For the group Lasso, we get using Example 13 that

$$\operatorname{div}(\hat{\mu})(y) = \operatorname{tr} \left(X_\Lambda \left(X_\Lambda^\top X_\Lambda + (\delta_{D^* \hat{\beta}(y)} \circ Q_{(D^* \hat{\beta}(y))^\perp})_{\Lambda, \Lambda} \right)^{-1} X_\Lambda^\top \right)$$

where $(\delta_{D^* \hat{\beta}(y)} \circ Q_{(D^* \hat{\beta}(y))^\perp})_{\Lambda, \Lambda}$ is the submatrix whose rows and columns are those of $\delta_{D^* \hat{\beta}(y)} \circ Q_{(D^* \hat{\beta}(y))^\perp}$ indexed by $\Lambda = \operatorname{supp}_{\mathcal{B}}(\hat{\beta}(y))$. This result was proved in (Vaiter et al 2012) in the overdetermined case. An immediate consequence of this formula is obtained when X is orthonormal⁵, in which case one recovers the expression of Yuan and Lin (2006),

$$\operatorname{div}(\hat{\mu})(y) = |\Lambda| - \sum_{b \in \mathcal{B}, D_b^* \hat{\beta}(y) \neq 0} \frac{|b| - 1}{\|y_b\|}.$$

The general group Lasso formula is new to the best of our knowledge, and will be illustrated in the numerical experiments on the isotropic 2-D total variation regularization widely used in image processing.

We could also provide a divergence formula for the nuclear norm, but as we discussed in Example 27, we cannot always guarantee the existence of a solution that satisfies $(\mathcal{C}_{\hat{\beta}(y), y})$. However, one can still find other partly smooth functions J with a non-flat submanifold for which this existence can be certified. The function of Example 17 is again a prototypical example.

⁵ Obviously, Lemma 2(ii) holds in such a case at the unique minimizer $\hat{\beta}(y)$.

Example 33 ($J = \max(\|\cdot\| - 1, 0)$) For $\beta \in \mathbb{S}^{p-1}$. If β is a minimizer of $(\mathcal{P}(y))$ is not a minimizer of $F(\cdot, y)$, from Example 28, we have that $\nabla_{\mathcal{M}}^2 F(\beta, y)$ is positive definite on $T = T_\beta$. Thus, we get for the case of the squared loss, that

$$\text{div}(\hat{\mu})(y) = \text{tr} \left(X_T^\top (X_T X_T^\top + P_T \langle X\beta, y - X\beta \rangle)^+ X_T^\top \right).$$

6 Degrees of Freedom and Unbiased Risk Estimation

From now on, we will assume that

$$\text{the set } \mathcal{M} \text{ is finite.} \quad (C_{\mathcal{M}})$$

Assumption $(C_{\mathcal{M}})$ holds in many important cases, including the examples discussed in the paper: polyhedral penalties (e.g. the Lasso, general Lasso or ℓ^∞ -norm), as well as for the group Lasso and its general form.

Throughout this section, we use the same symbols to denote weak derivatives (whenever they exist) as for derivatives. Rigorously speaking, the identities have to be understood to hold Lebesgue-a.e. (Evans and Gariepy 1992).

So far, we have shown that outside $\mathcal{H} \cup \mathcal{G}$, the mapping $y \mapsto \hat{\mu}(y)$ enjoys (locally) nice smoothness properties, which in turn gives closed-form formula of its divergence. To establish that such formula holds Lebesgue a.e., a key argument that we need to show is that \mathcal{H} is of negligible Lebesgue measure. This is where o-minimal geometry enters the picture. In turn, for Y drawn from some appropriate probability measures with density with respect to the Lebesgue measure, this will allow us to establish unbiasedness of quadratic risk estimators.

6.1 O-minimal Geometry

Roughly speaking, to be able to control the size of \mathcal{H} , the function J cannot be too oscillating in order to prevent pathological behaviours. We now briefly recall here the definition. Some important properties of o-minimal structures that are relevant to our context together with their proofs are collected in Section A. The interested reader may refer to (van den Dries 1998; Coste 1999) for a comprehensive account and further details on o-minimal structures.

Definition 5 (Structure) A structure \mathcal{O} expanding \mathbb{R} is a sequence $(\mathcal{O}_k)_{k \in \mathbb{N}}$ which satisfies the following axioms:

1. Each \mathcal{O}_k is a Boolean algebra of subsets of \mathbb{R}^k , with $\mathbb{R}^k \in \mathcal{O}_k$.
2. Every semi-algebraic subset of \mathbb{R}^k is in \mathcal{O}_k .
3. If $A \in \mathcal{O}_k$ and $B \in \mathcal{O}_{k'}$, then $A \times B \in \mathcal{O}_{k+k'}$.
4. If $A \in \mathcal{O}_{k+1}$, then $\Pi_{k+1,k}(A) \in \mathcal{O}_k$, where $\Pi_{k+1,k} : \mathbb{R}^{k+1} \rightarrow \mathbb{R}^k$ is the projection on the first k components.

The structure \mathcal{O} is said to be o-minimal if, moreover, it satisfies

5. (*o-minimality*) Sets in \mathcal{O}_1 are precisely the finite unions of intervals and points of \mathbb{R} .

In the following, a set $A \in \mathcal{O}_k$ is said to be definable.

Definition 6 (Definable set and function) *Let \mathcal{O} be an o-minimal structure. The elements of \mathcal{O}_k are called the definable subsets of \mathbb{R}^p , i.e. $\Omega \subset \mathbb{R}^k$ is definable if $\Omega \in \mathcal{O}_k$. A map $f : \Omega \rightarrow \mathbb{R}^m$ is said to be definable if its graph $\mathcal{G}(f) = \{(x, u) \in \Omega \times \mathbb{R}^m : u = f(x)\} \subseteq \mathbb{R}^k \times \mathbb{R}^m$ is a definable subset of $\mathbb{R}^k \times \mathbb{R}^m$ (in which case m times applications of axiom 4 implies that Ω is definable).*

A fundamental class of o-minimal structures is the collection of semi-algebraic sets, in which case axiom 4 is actually a property known as the Tarski-Seidenberg theorem (Coste 2002). For example, in the special case where q is a rational number, $J = \|\cdot\|_q$ is semi-algebraic. When $q \in \mathbb{R}$ is not rational, $\|\cdot\|_q$ is not semi-algebraic, however, it can be shown to be definable in an o-minimal structure. To see this, we recall from (van den Dries and Miller 1996, Example 5 and Property 5.2) that there exists a (polynomially bounded) o-minimal structure that contains the family of functions $\{t > 0 : t^\gamma, \gamma \in \mathbb{R}\}$ and restricted analytic functions. Functions F_0 that correspond to the exponential family losses introduced in Example 3 are also definable.

Our o-minimality assumptions requires the existence of an o-minimal structure \mathcal{O} such that

$$F, J \text{ and } \mathcal{M}, \forall \mathcal{M} \in \mathcal{M}, \text{ are definable in } \mathcal{O}. \quad (C_{\mathcal{O}})$$

6.2 Degrees of Freedom and Unbiased Risk Estimation

We assume in this section that F takes the form (3) and that

$$\forall y \in \mathbb{R}^n, \quad F_0(\cdot, y) \text{ is strongly convex with modulus } \tau \quad (C_{\text{sconv}})$$

and

$$\exists L > 0, \quad \sup_{(\mu, y) \in \mathbb{R}^n \times \mathbb{R}^n} \|D(\nabla F_0)(\mu, y)\| \leq L. \quad (C_L)$$

Obviously, assumption (C_{sconv}) implies (C_{dp}) , and thus the claims of the previous section remain true. Moreover, this assumption holds for the squared loss, but also for some losses of the exponential family (4), possibly adding a small quadratic term in β . As far as assumption (C_L) is concerned, it is easy to check that it is fulfilled with $L = 1$ for any loss of the exponential family (4), since $D(\nabla F_0)(\mu, y) = -\text{Id}$.

Non-linear Gaussian regression. Assume that the observation model (1) specializes to $Y \sim \mathcal{N}(h(X\beta_0), \sigma^2 \text{Id}_n)$, where h is Lipschitz continuous.

Theorem 3 *The following holds.*

- (i) Under condition $(C_{\mathcal{O}})$, \mathcal{H} is of Lebesgue measure zero;
- (ii) Under conditions (C_{sconv}) and (C_L) , $h \circ \hat{\mu}$ is Lipschitz continuous, hence weakly differentiable, with an essentially bounded gradient.
- (iii) If conditions $(C_{\mathcal{O}})$, (C_{sconv}) , (C_F) and (C_L) hold, and \mathcal{G} is of zero-Lebesgue measure, then,
 - (a) $\widehat{df} = \text{tr}(\text{D}h(\hat{\mu}(Y))\Delta(Y))$ is an unbiased estimate of $df = \mathbb{E}(\text{div}(h \circ \hat{\mu}(Y)))$, where $\Delta(Y)$ is as given in Theorem 2.
 - (b) The SURE

$$\text{SURE}(h \circ \hat{\mu})(Y) = \|Y - h(\hat{\mu}(Y))\|^2 + 2\sigma^2 \widehat{df} - n\sigma^2 \quad (24)$$

is an unbiased estimator of the risk $\mathbb{E}(\|h(\hat{\mu}(Y)) - h(\mu_0)\|^2)$.

This theorem is proved in Section 8.7.

GLM with the continuous exponential family. Assume that the observation model (1) corresponds to the GLM with a distribution which belongs to a continuous standard exponential family as parameterized in (2). From the latter, we have

$$\nabla \log B(y) = \left(\frac{\partial \log B_i(y_i)}{\partial y_i} \right)_i.$$

Theorem 4 *Suppose that conditions $(C_{\mathcal{O}})$, (C_{sconv}) , (C_F) and (C_L) hold, and \mathcal{G} is of zero-Lebesgue measure. Then,*

- (i) $\widehat{df} = \text{tr}(\Delta(Y))$ is an unbiased estimate of $df = \mathbb{E}(\text{div}(\hat{\mu}(Y)))$.
- (ii) The SURE

$$\text{SURE}(\hat{\mu})(Y) = \|\nabla \log B(Y) - \hat{\mu}(Y)\|^2 + 2\widehat{df} - (\|\nabla \log B(Y)\|^2 - \|\mu_0\|^2) \quad (25)$$

is an unbiased estimator of the risk $\mathbb{E}(\|\hat{\mu}(Y) - \mu_0\|^2)$.

This theorem is proved in Section 8.7. Recall from Section 5 that there are many regularizers where \mathcal{G} is indeed empty, and for which Theorem 3 and 4 then apply.

Though $\text{SURE}(\hat{\mu})(Y)$ depends on μ_0 , which is obviously unknown, it is only through an additive constant, which makes it suitable for parameter selection by risk minimization. Moreover, even if it is not stated here explicitly, Theorem 4 can be extended to unbiasedly estimate other measures of the risk, including the *projection* risk, or the *estimation* risk (in the full rank case) through the Generalized Stein Unbiased Risk Estimator as proposed in (Eldar 2009, Section IV), see also (Vaiter et al 2013) in the Gaussian case.

7 Simulation results

Experimental setting. In this section, we illustrate the efficiency of the proposed DOF estimator on a parameter selection problem in the context of some imaging inverse problems. More precisely, we consider the linear Gaussian regression model $Y \sim \mathcal{N}(X\beta_0, \sigma^2 \text{Id}_n)$ where $\beta_0 \in \mathbb{R}^{p=p_1 \times p_2}$ is a column-vectorized version of an image defined on a 2-D discrete grid of size $p_1 \times p_2$. The estimation of β_0 is achieved by solving $(\mathcal{P}(y))$ with

$$F(\beta, y) = F_0(X\beta, y) = \|X\beta - y\|^2 \quad \text{and} \quad J(\beta) = \lambda \|D^* \beta\|_{1,2}$$

where $D^* \beta \in \mathbb{R}^{p \times 2}$ is the 2-D discrete gradient vector field of the image β , and $\lambda > 0$ is the regularization parameter. Clearly, J is the isotropic total variation regularization (Rudin et al 1992), which is a special case of the general group Lasso penalty (9) for blocks of size 2.

We aim at proposing an automatic and objective way to choose λ . This can be achieved typically by minimizing the SURE given in (24) with h being the identity, i.e.

$$\text{SURE}(\hat{\mu})(Y) = \|Y - \hat{\mu}(Y)\|^2 + 2\sigma^2 \widehat{df} - n\sigma^2$$

where $\widehat{df} = \text{tr}(\Delta(Y))$ according to Theorem 3(iii)-(a), and the expression of $\Delta(Y)$ is obtained from that of the general group Lasso in Example 32 with D^* the discrete 2-D gradient operator, and $-D$ is the discrete 2-D divergence operator. Owing to Proposition 1(ii) and Theorem 3(iii), the given SURE is indeed an unbiased estimator of the prediction risk.

As the image size p can be large, the exact computation of $\text{tr}(\Delta(y))$ can become computationally intractable. Instead, we devise an approach based on Monte-Carlo (MC) simulations (see, Vonesch et al 2008, for more details), that is

$$\widehat{df}^{\text{MC}}(z) = \langle z, \Delta(Y)z \rangle$$

with z a realization of $Z \sim \mathcal{N}(0, \text{Id}_n)$. It is clear that $\mathbb{E}_Z \left(\widehat{df}^{\text{MC}}(Z) \right) = \widehat{df}$.

It remains to compute the vector $\Delta(y)z$. This is achieved by taking $\Delta(y)z = X\nu$, where ν is a solution of

$$\left(X^\top X + \lambda D(\delta_{D^* \widehat{\beta}(y)} \circ Q_{(D^* \widehat{\beta}(y))^\perp}) D^* \right) \nu = X^\top z \quad \text{subject to} \quad \nu \in T,$$

where we recall that $T = \text{Ker}(D_{\Lambda^c}^*)$, $\Lambda = \text{supp}_{\mathcal{B}}(D^* \widehat{\beta}(y))$. Taking into account the constraint on T through its Lagrange multiplier ζ , solving for ν boils down to solving the following linear system with a symmetric and positive-definite matrix

$$\begin{pmatrix} X^\top X + \lambda D(\delta_{D^* \widehat{\beta}(y)} \circ Q_{(D^* \widehat{\beta}(y))^\perp}) D^* & D_{\Lambda^c} \\ D_{\Lambda^c}^* & 0 \end{pmatrix} \begin{pmatrix} \nu \\ \zeta \end{pmatrix} = \begin{pmatrix} X^\top z \\ 0 \end{pmatrix}. \quad (26)$$

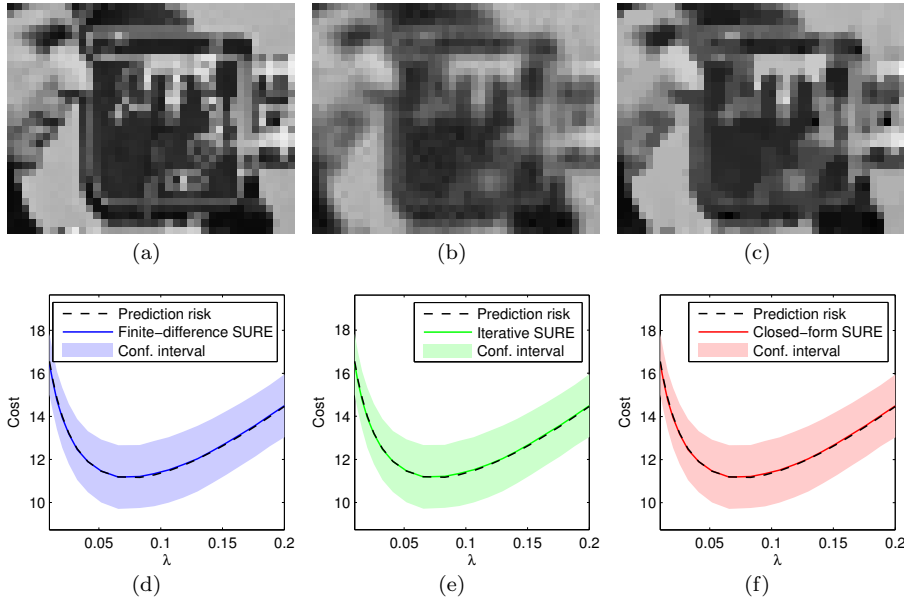


Fig. 1 (a) Original image β_0 . (b) Blurry observation y . (c) $\hat{\beta}(y)$ obtained for the value of λ minimizing the SURE estimate. (d-f) Prediction risk, average SURE and its confidence interval (\pm standard deviation) as a function of λ respectively for the finite difference approach (Ramani et al 2008), the iterative approach (Vonesch et al 2008), and our proposed approach.

Numerical solvers. In all experiments, optimization problem $(\mathcal{P}(y))$ was solved using Douglas-Rachford proximal splitting algorithm (Combettes and Pesquet 2007) with $2 \cdot 10^4$ iterations. Once the support Λ is identified with sufficiently high accuracy, the linear problem (26) is solved using the generalized minimal residual method (GMRES, Saad and Schultz 1986) with a relative accuracy of 10^{-7} .

Our proposed SURE estimator is compared for different values of λ with the approach of (Ramani et al 2008) based on finite difference approximations, as well as the approaches of (Vonesch et al 2008; Deledalle et al 2014) based on iterative chain rule differentiations. All curves are averaged on 40 independent realizations of Y and Z and their corresponding confidence intervals at \pm their standard deviation are displayed.

Deconvolution. We first consider an image of size $p = 34 \times 42$ with grayscale values ranging in $[0, 255]$ obtained from a close up of the standard *cameraman* image. X is a circulant matrix representing a periodic discrete convolution with a Gaussian kernel of width 1.5 pixel. The observation y is finally obtained by adding a zero-mean white Gaussian noise with $\sigma = 5$. Figure 1 depicts the evolution of the prediction risk and its SURE estimates as a function of λ .

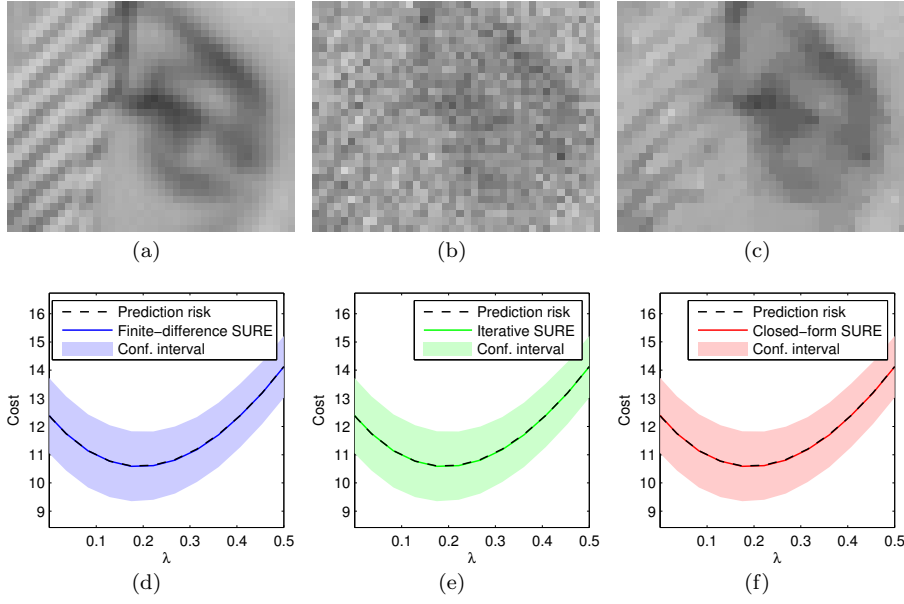


Fig. 2 (a) Original image β_0 . (b) Least squares estimate X^+y . (c) $\hat{\beta}(y)$ obtained for the value of λ minimizing the SURE estimate. (d-f) Prediction risk, average SURE and its confidence interval (\pm standard deviation) as a function of λ respectively for the finite difference approach (Ramani et al 2008), the iterative approach (Deledalle et al 2014), and our proposed approach.

Compressive sensing. We next consider an image of size $p = 34 \times 42$ with grayscale values ranging in $[0, 255]$ obtained from a close up of the standard *barbara* image. Now, X is a matrix corresponding to the composition of a periodic discrete convolution with a square kernel, and a random sub-sampling matrix with $n/p = 0.5$. The noise standard deviation is again $\sigma = 5$. Figure 2 shows the evolution of the prediction risk and its SURE estimates as a function of λ .

Discussion. The three approaches seem to provide the same results with average SURE curves that align very tightly with those of the prediction risk, with relatively small standard deviation compared to the range of variation of the prediction risk.

It is worth observing that the SURE obtained with finite differences (Ramani et al 2008) or with iterative differentiations (Vonesch et al 2008; Deledalle et al 2014) estimate the risk at the last iterate provided by the optimization algorithm to solve $(\mathcal{P}(y))$, which is not exactly $\hat{\beta}(y)$ in general. In fact, what is important is not $\hat{\beta}(y)$ by itself but rather its group support Λ . Thus, provided Λ has been perfectly identified, the three approaches provide, as observed, the same estimate of the risk up to machine precision. It may then be important to run the solver with a large number of iterations in order to provide an ac-

curate estimation of the risk. Even more important, solutions of (26) should be accurate enough to avoid bias in the estimation. The choice of $2 \cdot 10^4$ iterations for Douglas-Rachford and relative accuracy of 10^{-7} for GMRES appears in our simulations as a good trade-off between negligible bias and reasonable computational time.

8 Proofs

This section details the proofs of our results.

8.1 Preparatory lemma

By standard arguments of convex analysis, the following lemma gives the first-order sufficient and necessary optimality condition of a minimizer of $(\mathcal{P}(y))$.

Lemma 3 *A vector $\hat{\beta}(y) \in \mathbb{R}^p$ is a minimizer of $(\mathcal{P}(y))$ if, and only if,*

$$-\nabla F(\hat{\beta}(y), y) \in \partial J(\hat{\beta}(y)).$$

If J is partly smooth at $\hat{\beta}(y)$ relative to \mathcal{M} , then

$$-\nabla_{\mathcal{M}} F(\hat{\beta}(y), y) = \nabla_{\mathcal{M}} J(\hat{\beta}(y)) = e(\hat{\beta}(y)).$$

Proof The first monotone inclusion is just the first-order necessary and sufficient minimality condition for our convex program. The second claim follows from (12) and Fact 2. \square

8.2 Proof of Lemma 1

The equivalence is a consequence of simple arguments from linear algebra. Indeed, when both $\nabla_{\mathcal{M}}^2 F(\beta, y)$ and $\nabla_{\mathcal{M}}^2 J(\beta)$ are positive semidefinite on T , we have $\langle (\nabla_{\mathcal{M}}^2 F(\beta, y)\xi, \xi) \rangle \geq 0$ and $\langle (\nabla_{\mathcal{M}}^2 J(\beta)\xi, \xi) \rangle \geq 0, \forall \xi \in T$. Thus, for $(\mathcal{C}_{\beta, y})$ to hold, it is necessary and sufficient that $\nexists 0 \neq \xi \in T$ such that $\xi \in \text{Ker}(\nabla_{\mathcal{M}}^2 F(\beta, y))$ and $\xi \in \text{Ker}(\nabla_{\mathcal{M}}^2 J(\beta))$, which is exactly what we state.

When $\mathcal{M} = \beta + T$, the Riemannian Hessians $\nabla_{\mathcal{M}}^2 F(\beta, y)$ and $\nabla_{\mathcal{M}}^2 J(\beta)$ are given by (19) and (14). Convexity and smoothness of $F(\cdot, y)$ combined with (19) imply that $\nabla_{\mathcal{M}}^2 F(\beta, y)$ is positive semidefinite. Moreover, convexity and partial smoothness of J also yield that $\nabla_{\mathcal{M}}^2 J(\beta)$ is positive semidefinite, see (Liang et al 2014, Lemma 4.6). \square

8.3 Proof of Theorem 1

Let $y \notin \mathcal{H}$. To lighten the notation, we will drop the dependence of $\hat{\beta}$ on y , where $\hat{\beta}$ is a solution of $(\mathcal{P}(y))$ such that $(\mathcal{C}_{\hat{\beta},y})$ holds.

Let the constrained problem on \mathcal{M}

$$\min_{\beta \in \mathcal{M}} F(\beta, y) + J(\beta). \quad (\mathcal{P}(y)_{\mathcal{M}})$$

We define the notion of strong critical points that will play a pivotal role in our proof.

Definition 7 *A point $\hat{\beta}$ is a strong local minimizer of a function $f : \mathcal{M} \rightarrow \mathbb{R} \cup \{+\infty\}$ if f grows at least quadratically locally around $\hat{\beta}$ on \mathcal{M} , i.e. $\exists \delta > 0$ such that $f(\beta) \geq f(\hat{\beta}) + \delta \|\beta - \hat{\beta}\|^2$, $\forall \beta \in \mathcal{M}$ near $\hat{\beta}$.*

The following lemma gives an equivalent characterization of strong critical points that will be more convenient in our context.

Lemma 4 *Let $f \in C^2(\mathcal{M})$. A point $\hat{\beta}$ is a strong local minimizer of f if, and only if, it is a critical point of f , i.e. $\nabla_{\mathcal{M}} f(\hat{\beta}) = 0$, and satisfies the restricted positive definiteness condition*

$$\langle \nabla_{\mathcal{M}}^2 f(\hat{\beta}) \xi, \xi \rangle > 0 \quad \forall 0 \neq \xi \in \mathcal{T}_{\hat{\beta}}(\mathcal{M}).$$

Proof (of Lemma 4) The proof follows by combining the discussion after (Lewis 2003a, Definition 5.4) and (Miller and Malick 2005, Theorem 3.4). \square

We now define the following mapping

$$\Gamma : (\beta, y) \in \mathcal{M} \times \mathbb{R}^n \mapsto \nabla_{\mathcal{M}} F(\beta, y) + \nabla_{\mathcal{M}} J(\beta).$$

We split the proof of the theorem in three steps. We first show that there exists a continuously differentiable mapping $\bar{y} \mapsto \tilde{\beta}(\bar{y}) \in \mathcal{M}$ and an open neighborhood \mathcal{V}_y of y such that every element \bar{y} of \mathcal{V}_y satisfies $\Gamma(\tilde{\beta}(\bar{y}), \bar{y}) = 0$. Then, we prove that $\tilde{\beta}(\bar{y})$ is a solution of $(\mathcal{P}(\bar{y}))$ for any $\bar{y} \in \mathcal{V}_y$. Finally, we obtain (21) from the implicit function theorem.

Step 1: construction of $\tilde{\beta}(\bar{y})$. Using assumption (C_F) , the sum and smooth perturbation calculus rules of partial smoothness (Lewis 2003a, Corollary 4.6 and Corollary 4.7) entail that the function $(\beta, y) \mapsto F(\beta, y) + J(x)$ is partly smooth at $(\hat{\beta}, y)$ relative to $\mathcal{M} \times \mathbb{R}^m$, which is a C^2 -manifold of $\mathbb{R}^p \times \mathbb{R}^m$. Moreover, it is easy to see that $\mathcal{M} \times \mathbb{R}^m$ satisfies the transversality condition of (Lewis 2003a, Assumption 5.1). By assumption $(\mathcal{C}_{\hat{\beta},y})$, $\hat{\beta}$ is also a strong global minimizer of $(\mathcal{P}(y)_{\mathcal{M}})$, which implies in particular that $\Gamma(\hat{\beta}, y) = 0$; see Lemma 4. It then follows from (Lewis 2003a, Theorem 5.5) that there exist open neighborhoods $\tilde{\mathcal{V}}_y$ of y and $\tilde{\mathcal{V}}_{\hat{\beta}}$ of $\hat{\beta}$ and a continuously differentiable

mapping $\tilde{\beta} : \tilde{\mathcal{V}}_y \rightarrow \mathcal{M} \cap \tilde{\mathcal{V}}_{\tilde{\beta}}$ such that $\tilde{\beta}(y) = \hat{\beta}$, and $\forall \bar{y} \in \tilde{\mathcal{V}}_y$, $(\mathcal{P}(\bar{y})_{\mathcal{M}})$ has a *unique* strong local minimizer, i.e.

$$\Gamma(\tilde{\beta}(\bar{y}), \bar{y}) = 0 \quad \text{and} \quad (\mathcal{C}_{\tilde{\beta}(\bar{y}), \bar{y}}) \text{ holds,}$$

where we also used local normal sharpness property from partial smoothness of J ; see Fact 1.

Step 2: $\tilde{\beta}(\bar{y})$ is a solution of $(\mathcal{P}(\bar{y}))$. We now have to check the first-order optimality condition of $(\mathcal{P}(\bar{y}))$, i.e. that $-\nabla F(\tilde{\beta}(\bar{y}), \bar{y}) \in \partial J(\tilde{\beta}(\bar{y}))$; see Lemma 3. We distinguish two cases.

- Assume that $-\nabla F(\hat{\beta}, y) \in \text{ri } \partial J(\hat{\beta})$. The result then follows from (Lewis 2003a, Theorem 5.7(ii)) which, moreover, allows to assert in this case that $-\nabla F(\tilde{\beta}(\bar{y}), \bar{y}) \in \text{ri } \partial J(\tilde{\beta}(\bar{y}))$.
- We now turn to the case where $-\nabla F(\hat{\beta}, y) \in \text{rbd } \partial J(\hat{\beta})$. Observe that $(y, \hat{\beta}) \in \mathcal{A}_{\mathcal{M}}$. In particular $y \in \Pi_{n+p,n}(\mathcal{A}_{\mathcal{M}})$. Since by assumption $y \notin \mathcal{H}$, one has $y \notin \text{bd}(\Pi_{n+p,n}(\mathcal{A}_{\mathcal{M}}))$. Hence, there exists an open ball $\mathbb{B}(y, \varepsilon)$ for some $\varepsilon > 0$ such that $\mathbb{B}(y, \varepsilon) \subset \Pi_{n+p,n}(\mathcal{A}_{\mathcal{M}})$. Thus for every $\bar{y} \in \mathbb{B}(y, \varepsilon)$, there exists $\tilde{\beta} \in \widehat{\mathcal{M}}$ such that

$$-\nabla F(\tilde{\beta}, \bar{y}) \in \text{rbd } \partial J(\tilde{\beta}).$$

Since $\tilde{\beta} \in \mathcal{M}$, $\tilde{\beta}$ is also a critical point of $(\mathcal{P}(\bar{y})_{\mathcal{M}})$. But from Step 1, $\tilde{\beta}(\bar{y})$ is unique, whence we deduce that $\tilde{\beta}(\bar{y}) = \tilde{\beta}$. In turn, we conclude that

$$\forall \bar{y} \in \mathbb{B}(y, \varepsilon), \quad -\nabla F(\tilde{\beta}(\bar{y}), \bar{y}) \in \text{rbd } \partial J(\tilde{\beta}(\bar{y})) \subset \partial J(\tilde{\beta}(\bar{y})).$$

Step 3: Computing the differential. In summary, we have built a mapping $\tilde{\beta} \in \mathcal{C}^1(\mathcal{V})$, with $\mathcal{V} = \tilde{\mathcal{V}}_y \cap \mathbb{B}(y, \varepsilon)$, such that $\tilde{\beta}(\bar{y})$ is a solution of $(\mathcal{P}(\bar{y}))$ and fulfills $(\mathcal{C}_{\tilde{\beta}(\bar{y}), \bar{y}})$. We are then in position to apply the implicit function theorem to Γ , and we get the Jacobian of the mapping $\tilde{\beta}$ as

$$\mathcal{D}\tilde{\beta}(\bar{y}) = - \left(\nabla_{\mathcal{M}}^2 F(\tilde{\beta}(\bar{y}), \bar{y}) + \nabla_{\mathcal{M}}^2 J(\tilde{\beta}(\bar{y})) \right)^+ \mathcal{D}(\nabla_{\mathcal{M}} F)(\tilde{\beta}(\bar{y}), \bar{y})$$

where

$$\mathcal{D}(\nabla_{\mathcal{M}} F)(\beta, y) = \text{P}_{T_{\beta}} \mathcal{D}(\nabla F)(\beta, y),$$

where the equality is a consequence of (12) and linearity. \square

8.4 Proof of Lemma 2

(i) See (Vaiter et al 2015, Lemma 8).

(ii) This is a specialization of Lemma 1 using (C_{dp}) and (14). \square

8.5 Proof of Theorem 2

We can now prove Theorem 2. At any $y \notin \mathcal{H} \cup \mathcal{G}$, we consider $\widehat{\beta}(y)$ a solution of $(\mathcal{P}(y))$. By assumption, $(\mathcal{C}_{\widehat{\beta},y})$ holds. According to Theorem 1, one can construct a mapping $y \mapsto \tilde{\beta}(\bar{y})$ which is a solution to $(\mathcal{P}(\bar{y}))$, coincides with $\widehat{\beta}(y)$ at y , and is C^1 for \bar{y} in a neighborhood of y . Thus, by Lemma 2, $\widehat{\mu}(\bar{y}) = X\tilde{\beta}(\bar{y})$ is a single-valued mapping, which is also C^1 in a neighbourhood of y . Moreover, its differential is equal to $\Delta(y)$ as given, where we applied the chain rule in (18). \square

8.6 Proof of Proposition 1

The proofs of both statements are constructive.

- (i) Polyhedral penalty: any polyhedral convex J can be written as (Rockafellar 1996)

$$J(\beta) = \max_{i \in \{1, \dots, q\}} \{\langle d_i, \beta \rangle - b_i\} + \iota_{\mathcal{C}}(\beta),$$

$$\mathcal{C} = \{\beta \in \mathbb{R}^p : \langle a_k, \beta \rangle \leq c_k, k \in \{1, \dots, r\}\}.$$

It is straightforward to show that

$$\partial J(\beta) = \text{conv}\{d_i\}_{i \in I_\beta} + \text{cone}\{a_k\}_{k \in K_\beta}, \quad \text{where}$$

$$I_\beta = \{i : \langle d_i, \beta \rangle - b_i = J(\beta)\} \quad \text{and} \quad K_\beta = \{j : \langle a_j, \beta \rangle = c_j\},$$

and

$$T_\beta = \{h : \langle h, d_i \rangle = \langle h, d_j \rangle = \tau_\beta, \quad \forall i, j \in I_\beta\} \cap \{h : \langle h, a_k \rangle = 0, \quad \forall k \in K_\beta\}.$$

Let $\widehat{\beta}$ be a solution of $(\mathcal{P}(y))$ for J as above. Recall from Example 24 that $(\mathcal{C}_{\widehat{\beta},y})$ is equivalent to $\text{Ker}(X) \cap T_{\widehat{\beta}} = \{0\}$. Suppose that this condition does not. Thus, there exists a nonzero vector $h \in T_{\widehat{\beta}}$ such that the vector $v_t = \widehat{\beta} + th$, $t \in \mathbb{R}$, satisfies $Xv_t = X\widehat{\beta}$. Moreover,

$$\langle v_t, d_i \rangle - b_i = \begin{cases} J(\widehat{\beta}) + t\tau_{\widehat{\beta}}, & \text{if } i \in I_{\widehat{\beta}} \\ \langle \widehat{\beta}, d_i \rangle - b_i + t\langle h, d_i \rangle < J(\widehat{\beta}) + t\langle h, d_i \rangle & \text{otherwise.} \end{cases}$$

and

$$\langle v_t, a_k \rangle = \begin{cases} c_k, & \text{if } k \in K_{\widehat{\beta}} \\ \langle \widehat{\beta}, a_k \rangle + t\langle h, a_k \rangle < c_k + t\langle h, a_k \rangle & \text{otherwise.} \end{cases}$$

Thus, for $t \in]-t_0, t_0[$, where

$$t_0 = \min \left(\min_{i \notin I_{\widehat{\beta}}} \left\{ \frac{J(\widehat{\beta}) - \langle \widehat{\beta}, d_i \rangle + b_i}{|\langle h, d_i \rangle - \tau_{\widehat{\beta}}|} \right\}, \min_{k \notin K_{\widehat{\beta}}} \left\{ \frac{c_k - \langle \widehat{\beta}, a_k \rangle}{|\langle h, a_k \rangle|} \right\} \right),$$

we have $I_{v_t} = I_{\hat{\beta}}$ and $K_{v_t} = K_{\hat{\beta}}$. Moreover, $v_t \in \mathcal{C}$. Therefore, for all such t , we indeed have $\partial J(v_t) = \partial J(\hat{\beta})$ and $T_{v_t} = T_{\hat{\beta}}$. Altogether, we get that

$$-X^\top \nabla F_0(Xv_t, y) = -X^\top \nabla F_0(X\hat{\beta}, y) \in \partial J(\hat{\beta}) = \partial J(v_t),$$

i.e. v_t is a solution to $(\mathcal{P}(y))$. Thus, by Lemma 2, we deduce that $F_0(Xv_t, y) = F_0(X\hat{\beta}, y)$ and $J(v_t) = J(\hat{\beta})$. The continuity assumption (C_F) yields

$$F_0(Xv_{t_0}, y) = F_0(X\hat{\beta}, y).$$

Furthermore, since J is lsc and v_t is a minimizer of $(\mathcal{P}(y))$, we have

$$\liminf_{t \rightarrow t_0} J(v_t) \geq J(v_{t_0}) \geq \limsup_{t \rightarrow t_0} J(v_t) \iff J(v_{t_0}) = \lim_{t \rightarrow t_0} J(v_t) = J(\hat{\beta}).$$

Consequently, v_{t_0} is a solution of $(\mathcal{P}(y))$ such that $I_{\hat{\beta}} \subsetneq I_{v_{t_0}}$ or/and $K_{\hat{\beta}} \subsetneq K_{v_{t_0}}$, which in turn implies $T_{v_{t_0}} \subsetneq T_{\hat{\beta}}$. Iterating this argument, we conclude.

- (ii) General group Lasso: Let $\hat{\beta}$ be a solution of $(\mathcal{P}(y))$ for $J = \|D^* \cdot\|_{1,2}$, and $I_{\hat{\beta}} = \{i : b_i \in \mathcal{B} \text{ and } D_{b_i}^* \hat{\beta} \neq 0\}$, i.e. the set indexing the active blocks of $D^* \hat{\beta}$. We recall from Example 14 that the partial smoothness subspace $\mathcal{M} = T_{\hat{\beta}} = \text{Ker}(D_{\Lambda^c}^*)$, where $\Lambda = \text{supp}_{\mathcal{B}}(D^* \hat{\beta})$.

From Lemma 3 and the subdifferential of the group Lasso, $\hat{\beta}$ is indeed a minimizer if and only if there exists $\eta \in \mathbb{R}^p$ such that

$$-X^\top \nabla F_0(X\hat{\beta}, y) + \sum_{i \in I} D_{b_i} \eta_{b_i} = 0 \quad \text{and} \quad \begin{cases} \eta_{b_i} = \frac{D_{b_i}^* \hat{\beta}}{\|D_{b_i}^* \hat{\beta}\|} & \text{if } i \in I_{\hat{\beta}} \\ \|\eta_{b_i}\| \leq 1 & \text{otherwise.} \end{cases} \quad (27)$$

Suppose that $(\mathcal{C}_{\hat{\beta}, y})$ (or equivalently Lemma 2(ii)) does not hold at $\hat{\beta}$. This is equivalent to the existence of a nonzero vector $h \in \mathbb{R}^p$ in the set at the end of Example 26. Let $v_t = \hat{\beta} + th$, for $t \in \mathbb{R}$. By construction, v_t obeys

$$\begin{aligned} v_t \in T_{\hat{\beta}} &\iff \forall i \notin I_{\hat{\beta}}, D_{b_i}^* v_t = 0 \\ &\text{and } Xv_t = X\hat{\beta} \\ \text{and } \forall i \in I_{\hat{\beta}}, \exists \mu_i \in \mathbb{R}, D_{b_i}^* v_t &= (1 + t\mu_i) D_{b_i}^* \hat{\beta}. \end{aligned}$$

Let

$$t_0 = \min \{|t| : 1 + t\mu_i = 0, i \in I\} = \min_{i \in I_{\hat{\beta}}, \mu_i \neq 0} |\mu_i|^{-1}.$$

For all $t \in]-t_0, t_0[$, we have $1 + t\mu_i > 0$ for $i \in I_{\hat{\beta}}$ and $I_{v_t} = I_{\hat{\beta}}$ (in fact $T_{v_t} = T_{\hat{\beta}}$ by Fact 1), and thus

$$\frac{D_{b_i}^* v_t}{\|D_{b_i}^* v_t\|} = \frac{D_{b_i}^* \hat{\beta}}{\|D_{b_i}^* \hat{\beta}\|}, \quad \forall i \in I_{v_t}.$$

Moreover, $-X^\top \nabla F_0(Xv_t, y) = -X^\top \nabla F_0(X\hat{\beta}, y)$. Inserting the last statements in (27), we deduce that v_t is a solution of $(\mathcal{P}(y))$.

From Lemma 2(i), we get that $F_0(Xv_t, y) = F_0(X\hat{\beta}, y)$ and $\|D^*v_t\|_{1,2} = \|D^*\hat{\beta}\|_{1,2}$. By continuity of $F_0(\cdot, y)$ (assumption (C_F)), and of $\|\cdot\|_{1,2}$ one has

$$F_0(Xv_{t_0}) = F_0(X\hat{\beta}) \quad \text{and} \quad \|D^*v_{t_0}\|_{1,2} = \|D^*\hat{\beta}\|_{1,2}.$$

Clearly, we have constructed a solution v_{t_0} of $(\mathcal{P}(y))$ such that $I_{v_{t_0}} \subsetneq I_{\hat{\beta}}$, hence $\text{Ker}(\nabla_{\mathcal{M}}^2 J(v_{t_0})) \cap T_{v_{t_0}} \subsetneq \text{Ker}(\nabla_{\mathcal{M}}^2 J(\hat{\beta})) \cap T_{\hat{\beta}}$. Iterating this argument shows the result. \square

Remark 1 For the general group Lasso, the iterative construction is guaranteed to terminate at a non-trivial point. Indeed, if it were not the case, then eventually one would construct a solution such that $0 \neq h \in \text{Ker}(X) \cap \text{Ker}(D^*)$ leading to a contradiction with a classical condition in regularization theory. Moreover, $\text{Ker}(X) \cap \text{Ker}(D^*) = \{0\}$ is a sufficient (and necessary in our case) condition to ensure boundedness of the set of solutions to $(\mathcal{P}(y))$.

8.7 Proof of Theorem 3

(i) We obtain this assertion by proving that all $\mathcal{H}_{\mathcal{M}}$ are of zero measure for all \mathcal{M} , and that the union is over a finite set, because of $(C_{\mathcal{M}})$.

- Since J is definable by $(C_{\mathcal{O}})$, $\nabla F(\beta, y)$ is also definable by virtue of Proposition 2.
- Given $\mathcal{M} \in \mathcal{M}$ which is definable, $\widehat{\mathcal{M}}$ is also definable. Indeed, $\widehat{\mathcal{M}}$ can be equivalently written

$$\begin{aligned} \widehat{\mathcal{M}} &= \mathcal{M} \cap \{ \beta : \exists \varepsilon > 0, \forall \beta' \in \mathcal{M} \cap \mathbb{B}(\beta, \varepsilon), J \in \mathcal{C}^2(\beta') \} \\ &\cap \{ \beta : \forall (u, v) \in (\partial J(\beta))^2, \langle u - v, \beta' \rangle = 0, \forall \beta' \in \mathcal{T}_{\beta}(\mathcal{M}) \} \\ &\cap \{ \beta : \forall \beta_r \in \mathcal{M} \rightarrow \beta \text{ and } u \in \partial J(\beta), \exists u_r \rightarrow u \text{ s.t. } u_r \in \partial J(\beta_r) \} . \end{aligned}$$

Each of the four sets above capture a property of partial smoothness as introduced in Definition 1. $\widehat{\mathcal{M}}$ involves \mathcal{M} which is definable, its tangent space (which can be shown to be definable as a mapping of β using Proposition 2), ∂J whose graph is definable thanks to Proposition 3, continuity relations and algebraic equations, whence definability follows after interpreting the logical notations (conjunction, existence and universal quantifiers) in the first-order formula in terms of set operations, and using axioms 1-4 of definability in an o-minimal structure.

- Let $\mathbf{D} : \mathbb{R}^p \rightrightarrows \mathbb{R}^p$ the set-valued mapping whose graph is

$$\text{gph}(\mathbf{D}) = \{(\beta, \eta) : \eta \in \text{ri } \partial J(\beta)\} .$$

From Lemma 8, $\text{gph}(\mathbf{D})$ is definable. Since the graph ∂J is closed (Lemaréchal and Hiriart-Urruty 1996), and definable (Proposition 3), the set

$$\{(\beta, \eta) : \eta \in \text{rbd } \partial J(\beta)\} = \text{gph}(\partial J) \setminus \text{gph}(\mathbf{D}) ,$$

is also definable by axiom 1. This entails that $\mathcal{A}_{\mathcal{M}}$ is also a definable subset of $\mathbb{R}^n \times \widehat{\mathcal{M}}$ since

$$\begin{aligned} \mathcal{A}_{\mathcal{M}} = (\mathbb{R}^n \times \widehat{\mathcal{M}} \times \mathbb{R}^n) \cap \{ (y, \beta, \eta) : \eta = -\nabla F(\beta_T, y) \} \\ \cap (\mathbb{R}^n \times \{ (\beta, \eta) : \eta \in \text{rbd } \partial J(\beta) \}) . \end{aligned}$$

- By axiom 4, the canonical projection $\Pi_{n+p,n}(\mathcal{A}_{\mathcal{M}})$ is definable, and its boundary $\mathcal{H}_T = \text{bd}(\Pi_{n+p,n}(\mathcal{A}_{\mathcal{M}}))$ is also definable by (Coste 1999, Proposition 1.12) with a strictly smaller dimension than $\Pi_{n+p,n}(\mathcal{A}_{\mathcal{M}})$ (Coste 1999, Theorem 3.22).
- We recall now from (Coste 1999, Theorem 2.10) that any definable subset $A \subset \mathbb{R}^n$ in \mathcal{O} can be decomposed (stratified) in a disjoint finite union of q subsets C_i , definable in \mathcal{O} , called cells. The dimension of A is (Coste 1999, Proposition 3.17(4))

$$d = \max_{i \in \{1, \dots, q\}} d_i \leq n ,$$

where $d_i = \dim(C_i)$. Altogether we get that

$$\dim \mathcal{H}_{\mathcal{M}} = \dim \text{bd}(\Pi_{n+p,n}(\mathcal{A}_{\mathcal{M}})) < \dim \Pi_{n+p,n}(\mathcal{A}_{\mathcal{M}}) = d \leq n$$

whence we deduce that \mathcal{H} is of zero measure with respect to the Lebesgue measure on \mathbb{R}^n since the union is taken over the finite set \mathcal{M} by $(C_{\mathcal{M}})$.

(ii) $F_0(\cdot, y)$ is strongly convex with modulus τ if, and only if,

$$F_0(\mu, y) = G(\mu, y) + \frac{\tau}{2} \|\mu\|^2$$

where $G(\cdot, y)$ is convex and satisfies (C_F) , and in particular its domain in μ is full-dimensional. Thus, $(\mathcal{P}(y))$ amounts to solving

$$\min_{\beta \in \mathbb{R}^p} \frac{\tau}{2} \|X\beta\|^2 + G(X\beta, y) + J(\beta).$$

It can be recasted as a constrained optimization problem

$$\min_{\mu \in \mathbb{R}^n, \beta \in \mathbb{R}^p} \frac{\tau}{2} \|\mu\|^2 + G(\mu, y) + J(\beta) \text{ s.t. } \mu = X\beta.$$

Introducing the image (XJ) of J under the linear mapping X , it is equivalent to

$$\min_{\mu \in \mathbb{R}^n} \frac{\tau}{2} \|\mu\|^2 + G(\mu, y) + (XJ)(\mu) , \quad (28)$$

where $(XJ)(\mu) = \min_{\{\beta \in \mathbb{R}^p : \mu = X\beta\}} J(\beta)$ is the co-called pre-image of J under X . This is a proper closed convex function, which is finite on $\text{Span}(X)$. The minimization problem amounts to computing the proximal point at 0 of $G(\cdot, y) + (XJ)$, which is a proper closed and convex function. Thus this point exists and is unique.

Furthermore, by assumption (C_L) , the difference function

$$F_0(\cdot, y_1) - F_0(\cdot, y_2) = G(\cdot, y_1) - G(\cdot, y_2)$$

is Lipschitz continuous on \mathbb{R}^p with Lipschitz constant $L\|y_1 - y_2\|$. It then follows from (Bonnans and Shapiro 2000, Proposition 4.32) that $\widehat{\mu}(\cdot)$ is Lipschitz continuous with constant $2L/\tau$. Moreover, h is Lipschitz continuous, and thus so is the composed mapping $h \circ \widehat{\mu}(\cdot)$. From (Evans and Gariepy 1992, Theorem 5, Section 4.2.3), weak differentiability follows.

Rademacher theorem asserts that a Lipschitz continuous function is differentiable Lebesgue a.e. and its derivative and weak derivative coincide Lebesgue a.e., (Evans and Gariepy 1992, Theorem 2, Section 6.2). Its weak derivative, whenever it exists, is upper-bounded by the Lipschitz constant. Thus

$$\mathbb{E} \left(\left| \frac{\partial (h \circ \widehat{\mu})_i}{\partial y_i}(Y) \right| \right) < +\infty .$$

- (iii) Now, by the chain rule (Evans and Gariepy 1992, Remark, Section 4.2.2), the weak derivative of $h \circ \widehat{\mu}(\cdot)$ at y is precisely

$$D(h \circ \widehat{\mu})(y) = Dh(\widehat{\mu}(y)) \Delta(y) .$$

This formula is valid everywhere except on the set $\mathcal{H} \cup \mathcal{G}$ which is of Lebesgue measure zero as shown in (i). We conclude by invoking (ii) and Stein's lemma (Stein 1981) to establish unbiasedness of the estimator \widehat{df} of the DOF.

- (iv) Plugging the DOF expression (iii) into that of the SURE (Stein 1981, Theorem 1), the statement follows.
□

8.8 Proof of Theorem 4

For (i)-(iii), the proof is exactly the same as in Theorem 3. For (iv): combining the DOF expression (iii) and (Eldar 2009, Theorem 1), and rearranging the expression yields the stated result. □

9 Conclusion

In this paper, we proposed a detailed sensitivity analysis of a class of estimators obtained by minimizing a general convex optimization problem with a regularizing penalty encoding a low complexity prior. This was achieved through the concept of partial smoothness. This allowed us to derive an analytical expression of the local variations of these estimators to perturbations of the observations, and also to prove that the set where the estimator behaves non-smoothly as a function of the observations is of zero Lebesgue measure. Both results paved the way to derive unbiased estimators of the prediction risk in two random scenarios, one of which covers the continuous exponential family. This analysis covers a large set of convex variational estimators routinely used in statistics, machine learning and imaging (most notably group sparsity and multidimensional total variation penalty). The simulation results

confirm our theoretical findings and show that our risk estimator provides a viable way for automatic choice of the problem hyperparameters.

Despite its generality, there are still problems which do not fall within our settings. One can think for instance to the case of discrete (even exponential) distributions, risk estimation for non-canonical parameter of non-Gaussian distributions, non-convex regularizers, or the graphical Lasso.

Extension to the discrete case is far from obvious, even in the independent case. One can think for instance of using identities derived by (Hudson 1978; Hwang 1982), but so far, provably unbiased estimates of SURE (not generalized one) are only available for linear estimators.

If the distribution under consideration is from a continuous exponential family, so that our results apply, but one is interested in estimating the risk at a function of the canonical parameter. First, this function has to be Lipschitz continuous, and one has first to prove a formula of the corresponding SURE. So far, we are only aware of such results in the Gaussian case (hence our Theorem 3 which addresses this question precisely).

Strictly speaking, the ℓ^1 -penalized likelihood formulation of the graphical Lasso in (Yuan and Lin 2007) ((3) or (6) in that reference) does not fall within our framework. This is due to the fidelity/likelihood term which does not obey our assumptions. Note that the limitation due to fidelity/likelihood can be circumvented at the price of a quadratic approximation (Yuan and Lin 2007, Section 4) also used in (Meinshausen and Bühlmann 2006).

Extending our results to the non-convex case would be very interesting to handle penalties such as SCAD or MCP. This would however require more sophisticated material from variational analysis. Not to mention the other difficulties inherent to non-convexity, including handling critical points (that are not necessarily minimizers even local in general), and the fact that the mapping $y \mapsto \hat{\mu}(y)$ is no longer single-valued. All the above settings will be left to future work.

Acknowledgements This work has been supported by the European Research Council (ERC project SIGMA-Vision) and Institut Universitaire de France.

A Basic Properties of o-minimal Structures

In the following results, we collect some important stability properties of o-minimal structures. To be self-contained, we also provide proofs. To the best of our knowledge, these proofs, although simple, are not reported in the literature or some of them are left as exercises in the authoritative references van den Dries (1998); Coste (1999). Moreover, in most proofs, to show that a subset is definable, we could just write the appropriate first-order formula (see (Coste 1999, Page 12)(van den Dries 1998, Section Ch1.1.2)), and conclude using (Coste 1999, Theorem 1.13). Here, for the sake of clarity and avoid cryptic statements for the non-specialist, we will translate the first order formula into operations on the involved subsets, in particular projections, and invoke the above stability axioms of o-minimal structures. In the following, n denotes an arbitrary (finite) dimension which is not necessarily the number of observations used previously the paper.

Lemma 5 (Addition and multiplication) *Let $f : \Omega \subset \mathbb{R}^n \rightarrow \mathbb{R}^p$ and $g : \Omega \subset \mathbb{R}^n \subset \mathbb{R}^p$ be definable functions. Then their pointwise addition and multiplication is also definable.*

Proof Let $h = f + g$, and

$$B = (\Omega \times \mathbb{R} \times \Omega \times \mathbb{R} \times \Omega \times \mathbb{R}) \cap (\Omega \times \mathbb{R} \times \text{gph}(f) \times \text{gph}(h)) \cap S$$

where $S = \{(x, u, y, v, z, w) : x = y = z, u = v + w\}$ is obviously an algebraic (in fact linear) subset, hence definable by axiom 2. Axiom 1 and 2 then imply that B is also definable. Let $\Pi_{3n+3p, n+p} : \mathbb{R}^{3n+3p} \rightarrow \mathbb{R}^{n+p}$ be the projection on the first $n + p$ coordinates. We then have

$$\text{gph}(h) = \Pi_{3n+3p, n+p}(B)$$

whence we deduce that h is definable by applying $3n + 3p$ times axiom 4. Definability of the pointwise multiplication follows the same proof taking $u = v \cdot w$ in S . \square

Lemma 6 (Inequalities in definable sets) *Let $f : \Omega \subset \mathbb{R}^n \rightarrow \mathbb{R}$ be a definable function. Then $\{x \in \Omega : f(x) > 0\}$, is definable. The same holds when replacing $>$ with $<$.*

Clearly, inequalities involving definable functions are accepted when defining definable sets. There are many possible proofs of this statement.

Proof (1) Let $B = \{(x, y) \in \mathbb{R} \times \mathbb{R} : f(x) = y\} \cap (\Omega \times (0, +\infty))$, which is definable thanks to axioms 1 and 3, and that the level sets of a definable function are also definable. Thus

$$\{x \in \Omega : f(x) > 0\} = \{x \in \Omega : \exists y, f(x) = y, y > 0\} = \Pi_{n+1, n}(B),$$

and we conclude using again axiom 4. \square

Yet another (simpler) proof.

Proof (2) It is sufficient to remark that $\{x \in \Omega : f(x) > 0\}$ is the projection of the set $\{(x, t) \in \Omega \times \mathbb{R} : t^2 f(x) - 1 = 0\}$, where the latter is definable owing to Lemma 5. \square

Lemma 7 (Derivative) *Let $f : I \rightarrow \mathbb{R}$ be a definable differentiable function on an open interval I of \mathbb{R} . Then its derivative $f' : I \rightarrow \mathbb{R}$ is also definable.*

Proof Let $g : (x, t) \in I \times \mathbb{R} \mapsto g(x, t) = f(x+t) - f(x)$. Note that g is definable function on $I \times \mathbb{R}$ by Lemma 5. We now write the graph of f' as

$$\text{gph}(f') = \{(x, y) \in I \times \mathbb{R} : \forall \varepsilon > 0, \exists \delta > 0, \forall t \in \mathbb{R}, |t| < \delta, |g(x, t) - yt| < \varepsilon|t|\}.$$

Let $C = \{(x, y, v, t, \varepsilon, \delta) \in I \times \mathbb{R}^5 : ((x, t), v) \in \text{gph}(g)\}$, which is definable since g is definable and using axiom 3. Let

$$B = \{(x, y, v, t, \varepsilon, \delta) : t^2 < \delta^2, (v - ty)^2 < \varepsilon^2 t^2\} \cap C.$$

The first part in B is semi-algebraic, hence definable thanks to axiom 2. Thus B is also definable using axiom 1. We can now write

$$\text{gph}(f') = \mathbb{R}^3 \setminus \left(\Pi_{5,3} \left(\mathbb{R}^5 \setminus \Pi_{6,5}(B) \right) \right) \cap (I \times \mathbb{R}),$$

where the projectors and completions translate the actions of the existential and universal quantifiers. Using again axioms 4 and 1, we conclude. \square

With such a result at hand, this proposition follows immediately.

Proposition 2 (Differential and Jacobian) *Let $f = (f_1, \dots, f_p) : \Omega \rightarrow \mathbb{R}^p$ be a differentiable function on an open subset Ω of \mathbb{R}^n . If f is definable, then so its differential mapping and its Jacobian. In particular, for each $i = 1, \dots, n$ and $j = 1, \dots, p$, the partial derivative $\partial f_i / \partial x_j : \Omega \rightarrow \mathbb{R}$ is definable.*

We provide below some results concerning the subdifferential.

Proposition 3 (Subdifferential) *Suppose that f is a finite-valued convex definable function. Then for any $x \in \mathbb{R}^n$, the subdifferential $\partial f(x)$ is definable.*

Proof For every $x \in \mathbb{R}^n$, the subdifferential $\partial f(x)$ reads

$$\partial f(x) = \left\{ \eta \in \mathbb{R}^n : f(x') \geq f(x) + \langle \eta, x' - x \rangle \quad \forall x' \in \mathbb{R}^n \right\}.$$

Let $K = \{(\eta, x') \in \mathbb{R}^n \times \mathbb{R}^n : f(x') < f(x) + \langle \eta, x' - x \rangle\}$. Hence, $\partial f(x) = \mathbb{R}^n \setminus \Pi_{2n,n}(K)$. Since f is definable, the set K is also definable using Lemma 5 and 6, whence definability of $\partial f(x)$ follows using axiom 4. \square

Lemma 8 (Graph of the relative interior) *Suppose that f is a finite-valued convex definable function. Then, the set*

$$\{(x, \eta) : \eta \in \text{ri } \partial f(x)\}$$

is definable.

Proof Denote $C = \{(x, \eta) : \eta \in \text{ri } \partial f(x)\}$. Using the characterization of the relative interior of a convex set (Rockafellar 1996, Theorem 6.4), we rewrite C in the more convenient form

$$\begin{aligned} C = \{(x, \eta) : & \forall u \in \mathbb{R}^n, \forall z \in \mathbb{R}^n, f(z) - f(x) \geq \langle u, z - x \rangle, \\ & \exists t > 1, \forall x' \in \mathbb{R}^n, f(x') - f(x) \geq \langle (1-t)u + t\eta, x' - x \rangle\}. \end{aligned}$$

Let $D = \mathbb{R}^n \times \mathbb{R}^n \times \mathbb{R}^n \times \mathbb{R}^n \times (1, +\infty) \times \mathbb{R}^n$ and K defined as

$$K = \{(x, \eta, u, z, t, x') \in D : f(z) - f(x) \geq \langle u, z - x \rangle, f(x') - f(x) \geq \langle (1-t)u + t\eta, x' - x \rangle\}.$$

Thus,

$$C = \mathbb{R}^{2n} \setminus \Pi_{3n,2n} \left(\mathbb{R}^{3n} \setminus \Pi_{4n,3n} \left(\Pi_{4n+1,4n} \left(\mathbb{R}^{4n} \times (1, +\infty) \setminus \Pi_{5n+1,4n+1}(K) \right) \right) \right),$$

where the projectors and completions translate the actions of the existential and universal quantifiers. Using again axioms 4 and 1, we conclude. \square

References

- Absil PA, Mahony R, Trumf J (2013) An extrinsic look at the riemannian hessian. In: Geometric Science of Information, Lecture Notes in Computer Science, vol 8085, Springer Berlin Heidelberg, pp 361–368
- Bach F (2008) Consistency of the group lasso and multiple kernel learning. Journal of Machine Learning Research 9:1179–1225
- Bach F (2010) Self-concordant analysis for logistic regression. Electronic Journal of Statistics 4:384–414
- Bakin S (1999) Adaptive regression and model selection in data mining problems. Thesis (Ph.D.)—Australian National University, 1999
- Bickel PJ, Ritov Y, Tsybakov A (2009) Simultaneous analysis of lasso and Dantzig selector. Annals of Statistics 37(4):1705–1732
- Bolte J, Daniilidis A, Lewis AS (2011) Generic optimality conditions for semialgebraic convex programs. Mathematics of Operations Research 36(1):55–70
- Bonnans J, Shapiro A (2000) Perturbation analysis of optimization problems. Springer Series in Operations Research, Springer-Verlag, New York
- Brown LD (1986) Fundamentals of Statistical Exponential Families with Applications in Statistical Decision Theory, Monograph Series, vol 9. Institute of Mathematical Statistics Lecture Notes, IMS, Hayward, CA

- Bühlmann P, van de Geer S (2011) *Statistics for High-Dimensional Data: Methods, Theory and Applications*. Springer
- Bunea F (2008) Honest variable selection in linear and logistic regression models via ℓ_1 and $\ell_1 + \ell_2$ penalization. *Electronic Journal of Statistics* 2:1153–1194
- Candès E, Plan Y (2009) Near-ideal model selection by ℓ_1 minimization. *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? *J ACM* 58(3):11:1–11:37
- Candès EJ, Sing-Long CA, Trzasko JD (2012) Unbiased risk estimates for singular value thresholding and spectral estimators. *IEEE Transactions on Signal Processing* 61(19):4643–4657
- 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
- Chavel I (2006) *Riemannian geometry: a modern introduction*, Cambridge Studies in Advanced Mathematics, vol 98, 2nd edn. Cambridge University Press
- 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
- Combettes P, Pesquet J (2007) A Douglas–Rachford splitting approach to nonsmooth convex variational signal recovery. *IEEE Journal of Selected Topics in Signal Processing* 1(4):564–574
- Coste M (1999) An introduction to α -minimal geometry. Tech. rep., Institut de Recherche Mathématiques de Rennes
- Coste M (2002) An introduction to semialgebraic geometry. Tech. rep., Institut de Recherche Mathématiques de Rennes
- Daniilidis A, Hare W, Malick J (2009) Geometrical interpretation of the predictor-corrector type algorithms in structured optimization problems. *Optimization: A Journal of Mathematical Programming & Operations Research* 55(5-6):482–503
- Daniilidis A, Drusvyatskiy D, Lewis AS (2013) Orthogonal invariance and identifiability. Tech. rep., arXiv 1304.1198
- DasGupta A (2008) *Asymptotic Theory of Statistics and Probability*. Springer
- Deledalle CA, Vaïter S, Peyré G, Fadili M, Dossal C (2012) Risk estimation for matrix recovery with spectral regularization. In: *ICML’12 Workshop on Sparsity, Dictionaries and Projections in Machine Learning and Signal Processing*, (arXiv:1205.1482)
- Deledalle CA, Vaïter S, Peyré G, Fadili JM (2014) Stein unbiased gradient estimator of the risk (SUGAR) for multiple parameter selection. *SIAM J Imaging Sciences* 7(4):2448–2487
- Donoho D (2006) For most large underdetermined systems of linear equations the minimal ℓ^1 -norm solution is also the sparsest solution. *Communications on pure and applied mathematics* 59(6):797–829
- Dossal C, Kachour M, Fadili MJ, Peyré G, Chesneau C (2013) The degrees of freedom of penalized ℓ_1 minimization. *Statistica Sinica* 23(2):809–828
- Drusvyatskiy D, Lewis A (2011) Generic nondegeneracy in convex optimization. *Proc Amer Math Soc* 129:2519–2527
- Drusvyatskiy D, Ioffe A, Lewis A (2015) Generic minimizing behavior in semi-algebraic optimization. *SIAM J Optim* To appear
- Efron B (1986) How biased is the apparent error rate of a prediction rule? *Journal of the American Statistical Association* 81(394):461–470
- Eldar YC (2009) Generalized SURE for exponential families: Applications to regularization. *IEEE Transactions on Signal Processing* 57(2):471–481
- Evans LC, Gariepy RF (1992) *Measure theory and fine properties of functions*. CRC Press
- Fazel M, Hindi H, Boyd SP (2001) A rank minimization heuristic with application to minimum order system approximation. In: *American Control Conference, 2001. Proceedings of the 2001, IEEE*, vol 6, pp 4734–4739

- van de Geer SA (2008) High-dimensional generalized linear models and the lasso. *Annals of Statistics* 36:614–645
- de Geer SV (2008) High-dimensional generalized linear models and the lasso. *Annals of Statistics* 36(2):614–645
- Hansen NR, Sokol A (2014) Degrees of freedom for nonlinear least squares estimation. Tech. rep., arXiv preprint 1402.2997
- Hudson H (1978) A natural identity for exponential families with applications in multiparameter estimation. *The Annals of Statistics* 6(3):473–484
- Hwang JT (1982) Improving upon standard estimators in discrete exponential families with applications to poisson and negative binomial cases. *Ann Statist* 10(3):857–867
- Jacob L, Obozinski G, Vert JP (2009) Group lasso with overlap and graph lasso. In: Danyluk AP, Bottou L, Littman ML (eds) *Proc. ICML 2009*, vol 382, p 55
- Jégou H, Furon T, Fuchs JJ (2012) Anti-sparse coding for approximate nearest neighbor search. In: *Acoustics, Speech and Signal Processing (ICASSP), 2012 IEEE International Conference on*, IEEE, pp 2029–2032
- Kakade SM, Shamir O, Sridharan K, Tewari A (2010) Learning exponential families in high-dimensions: Strong convexity and sparsity. In: *AISTATS*
- Kato K (2009) On the degrees of freedom in shrinkage estimation. *Journal of Multivariate Analysis* 100(7):1338–1352
- Lee JM (2003) *Smooth manifolds*. Springer
- Lemaréchal C, Hiriart-Urruty J (1996) *Convex analysis and minimization algorithms: Fundamentals*, vol 305. Springer-Verlag
- Lemaréchal C, Oustry F, Sagastizábal C (2000) The \mathcal{U} -lagrangian of a convex function. *Trans Amer Math Soc* 352(2):711–729
- Lewis A (1995) The convex analysis of unitarily invariant matrix functions. *Journal of Convex Analysis* 2:173–183
- Lewis A, Sendov H (2001) Twice differentiable spectral functions. *SIAM Journal on Matrix Analysis on Matrix Analysis and Applications* 23:368–386
- Lewis AS (2003a) Active sets, nonsmoothness, and sensitivity. *SIAM Journal on Optimization* 13(3):702–725
- Lewis AS (2003b) The mathematics of eigenvalue optimization. *Mathematical Programming* 97(1–2):155–176
- Lewis AS, Zhang S (2013) Partial smoothness, tilt stability, and generalized Hessians. *SIAM Journal on Optimization* 23(1):74–94
- Liang J, Fadili MJ, Peyré G, Luke R (2014) Activity Identification and Local Linear Convergence of Douglas–Rachford/ADMM under Partial Smoothness. arXiv:14126858
- Liu H, Zhang J (2009) Estimation consistency of the group lasso and its applications. *Journal of Machine Learning Research* 5:376–383
- Lyubarskii Y, Vershynin R (2010) Uncertainty principles and vector quantization. *Information Theory, IEEE Transactions on* 56(7):3491–3501
- McCullagh P, Nelder JA (1989) *Generalized Linear Models*, second edition edn. *Monographs on Statistics & Applied Probability*, Chapman & Hall/CRC, URL <http://www.worldcat.org/isbn/0412317605>
- Meier L, Geer SVD, Bühlmann P (2008) The group lasso for logistic regression. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 70(1):51–71
- Meinshausen N, Bühlmann P (2006) High-dimensional graphs and variable selection with the lasso. *Annals of Statistics* 34:1436–1462
- Meyer M, Woodroffe M (2000) On the degrees of freedom in shape-restricted regression. *Annals of Statistics* 28(4):1083–1104
- Miller SA, Malick J (2005) Newton methods for nonsmooth convex minimization: connections among-lagrangian, riemannian newton and sqp methods. *Mathematical programming* 104(2–3):609–633
- Mordukhovich B (1992) Sensitivity analysis in nonsmooth optimization. *Theoretical Aspects of Industrial Design* (D A Field and V Komkov, eds), *SIAM Volumes in Applied Mathematics* 58:32–46
- 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

- 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: *Proc. EUSIPCO 2011, EURASIP, Barcelona, Espagne*, URL <http://hal.archives-ouvertes.fr/hal-00597772>
- Ramani S, Blu T, Unser M (2008) Monte-Carlo SURE: a black-box optimization of regularization parameters for general denoising algorithms. *IEEE Trans Image Process* 17(9):1540–1554
- 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
- Rockafellar RT (1996) *Convex Analysis*. Princeton Landmarks in Mathematics and Physics, Princeton University Press
- Rudin L, Osher S, Fatemi E (1992) Nonlinear total variation based noise removal algorithms. *Physica D: Nonlinear Phenomena* 60(1-4):259–268
- Saad Y, Schultz MH (1986) Gmres: A generalized minimal residual algorithm for solving nonsymmetric linear systems. *SIAM Journal on scientific and statistical computing* 7(3):856–869
- Solo V, Ulfarsson M (2010) Threshold selection for group sparsity. In: *Acoustics Speech and Signal Processing (ICASSP), 2010 IEEE International Conference on*, IEEE, pp 3754–3757
- Stein C (1981) Estimation of the mean of a multivariate normal distribution. *The Annals of Statistics* 9(6):1135–1151
- Studer C, Yin W, Baraniuk RG (2012) Signal representations with minimum ℓ_∞ -norm. In: *Communication, Control, and Computing, Proc. 50th Ann. Allerton Conf. on*
- 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
- Tibshirani RJ, Taylor J (2012) Degrees of freedom in Lasso problems. *Ann Statist* 40(2):639–1284
- Tikhonov AN, Arsenin VY (1997) *Solutions of Ill-posed Problems*. V. H. Winston and Sons
- Vaïter S, Deledalle C, Peyré G, Fadili MJ, Dossal C (2012) Degrees of freedom of the group Lasso. In: *ICML’12 Workshops*, pp 89–92
- Vaïter S, Deledalle C, Peyré G, Dossal C, Fadili MJ (2013) Local behavior of sparse analysis regularization: Applications to risk estimation. *Applied and Computational Harmonic Analysis* 35(3):433–451
- Vaïter S, Peyré G, Fadili MJ (2014) Model Consistency of Partly Smooth Regularizers. *arXiv:1405.1004*
- Vaïter S, Golbabaee M, Fadili MJ, Peyré G (2015) Model selection with low complexity priors. *Information and Inference: A Journal of the IMA (IMAIAI)*
- van den Dries L (1998) Tame topology and o-minimal structures, *Math. Soc. Lecture Note*, vol 248. Cambridge Univ Press
- van den Dries L, Miller C (1996) Geometric categories and o-minimal structures. *Duke Math J* 84:497–540
- Vonesch C, Ramani S, Unser M (2008) Recursive risk estimation for non-linear image deconvolution with a wavelet-domain sparsity constraint. In: *ICIP, IEEE*, pp 665–668
- Wei F, Huang J (2010) Consistent group selection in high-dimensional linear regression. *Bernoulli* 16(4):1369–1384
- Wright SJ (1993) Identifiable surfaces in constrained optimization. *SIAM Journal on Control and Optimization* 31(4):1063–1079
- Yuan M, Lin Y (2006) Model selection and estimation in regression with grouped variables. *J of The Roy Stat Soc B* 68(1):49–67
- Yuan M, Lin Y (2007) Model selection and estimation in the gaussian graphical model. *Biometrika* 94(1):19–35
- Zou H, Hastie T, Tibshirani R (2007) On the “degrees of freedom” of the Lasso. *The Annals of Statistics* 35(5):2173–2192