



**HAL**  
open science

# Model Consistency of Partly Smooth Regularizers

Samuel Vaïter, Gabriel Peyré, Jalal M. Fadili

► **To cite this version:**

Samuel Vaïter, Gabriel Peyré, Jalal M. Fadili. Model Consistency of Partly Smooth Regularizers. IEEE Transactions on Information Theory, 2018, 64 (3), pp.1725-1737. 10.1109/TIT.2017.2713822 . hal-01658847

**HAL Id: hal-01658847**

**<https://hal.science/hal-01658847>**

Submitted on 7 Dec 2017

**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.

# Model Consistency of Partly Smooth Regularizers

Samuel Vaïter, Gabriel Peyré and Jalal Fadili

**Abstract**—This paper studies least-square regression penalized with partly smooth convex regularizers. This class of penalty functions is very large and versatile, and allows to promote solutions conforming to some notion of low-complexity. Indeed, such penalties/regularizers force the corresponding solutions to belong to a low-dimensional manifold (the so-called model) which remains stable when the penalty function undergoes small perturbations. Such a good sensitivity property is crucial to make the underlying low-complexity (manifold) model robust to small noise. In a deterministic setting, we show that a generalized “irrepresentable condition” implies stable model selection under small noise perturbations in the observations and the design matrix, when the regularization parameter is tuned proportionally to the noise level. We also prove that this condition is almost necessary for stable model recovery. We then turn to the random setting where the design matrix and the noise are random, and the number of observations grows large. We show that under our generalized “irrepresentable condition”, and a proper scaling of the regularization parameter, the regularized estimator is model consistent. In plain words, with a probability tending to one as the number of measurements tends to infinity, the regularized estimator belongs to the correct low-dimensional model manifold. This work unifies and generalizes a large body of literature, where model consistency was known to hold, for instance for the Lasso, group Lasso, total variation (fused Lasso) and nuclear/trace norm regularizers. We show that under the deterministic model selection conditions, the forward-backward proximal splitting algorithm used to solve the penalized least-square regression problem, is guaranteed to identify the model manifold after a finite number of iterations. Lastly, we detail how our results extend from the quadratic loss to an arbitrary smooth and strictly convex loss function. We illustrate the usefulness of our results on the problem of low-rank matrix recovery from random measurements using nuclear norm minimization.

**Index Terms**—Regularization, regression, inverse problems, model consistency, partial smoothness, sensitivity analysis, sparsity, low-rank.

## I. INTRODUCTION

### A. Problem Statement

We consider the following observation model

$$y = \Phi x_0 + w,$$

where  $\Phi \in \mathbb{R}^{p \times n}$  is the design matrix (in statistics or machine learning) or the forward operator (in signal and imaging sciences),  $x_0 \in \mathbb{R}^n$  is the vector to recover and  $w \in \mathbb{R}^p$  is the noise. The design can be either deterministic or random, and similarly for the noise  $w$ .

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

Samuel Vaïter is with CNRS and IMB, Université de Bourgogne and was affiliated with CNRS and CEREMADE, Université Paris Dauphine when this work was completed. Gabriel Peyré is with CNRS and CEREMADE, Université Paris Dauphine. Jalal Fadili is with the Normandie Univ, ENSICAEN, CNRS, GREYC.

some prior structure on the object  $x_0$  to be estimated. We therefore consider a proper, lower-semicontinuous (lsc) and convex function  $J : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$  to promote such a prior. Without loss of generality, we also assume that  $J$  is non-negative. This then leads to solving the following convex optimization problem

$$\min_{x \in \mathbb{R}^n} \left\{ J(x) + \frac{1}{2\lambda} \|\Phi x - y\|^2 \right\}, \quad (1)$$

where  $\lambda > 0$  is the so-called regularization parameter used to balance the amount of regularization and loss.

To simplify the notations, we introduce the following “canonical” parameters

$$\theta = (\mu, u, \Gamma) = \left( \frac{\lambda}{p}, \frac{\Phi^* y}{p}, \frac{\Phi^* \Phi}{p} \right) \in \Theta = \mathbb{R}^+ \times \mathbb{R}^n \times \mathbb{R}^{n \times n}$$

and we denote

$$\varepsilon = \frac{\Phi^* w}{p} = u - \Gamma x_0,$$

where  $\Phi^*$  is the adjoint operator to the linear operator  $\Phi$ . In the following, we assume that  $y \in \text{Im}(\Phi)$  and thus  $u \in \text{Im}(\Gamma)$ . Obviously, this does not entail any loss of generality, as the loss term can always be written as  $\frac{1}{2\lambda} \|\Phi x - P_{\text{Im}(\Phi)} y\|^2 + \frac{1}{2\lambda} \|P_{\text{Im}(\Phi)^\perp} y\|^2$ , where  $P_T$  is the orthogonal projection on  $T$ .

With these new parameters, the original problem (1) now reads

$$\min_{x \in \mathbb{R}^n} E(x, \theta) \quad (\mathcal{P}_\theta)$$

where

$$E(x, \theta) = J(x) + \frac{1}{2\mu} \langle \Gamma x, x \rangle - \frac{1}{\mu} \langle x, u \rangle + \frac{1}{2\mu} \langle \Gamma^+ u, u \rangle,$$

and  $A^+$  stands for the Moore-Penrose pseudo-inverse of a matrix  $A$ . With these notations,  $E$  is a function on  $\mathbb{R}^n \times \Theta$ .

We also consider the constrained problem

$$\min_{x \in \mathbb{R}^n} \{ E(x, \theta_0) = J(x) + \iota_{\mathcal{H}_u}(x) \} \quad (\mathcal{P}_{\theta_0})$$

where

$$\mathcal{H}_u = \{ x \in \mathbb{R}^n ; \Gamma x = u \},$$

$\theta_0 = (0, u, \Gamma)$ , and  $\iota_{\mathcal{C}}$  is the indicator function of the non-empty closed convex set  $\mathcal{C}$ , i.e.  $\iota_{\mathcal{C}}(x) = 0$  if  $x \in \mathcal{C}$  and  $\iota_{\mathcal{C}}(x) = +\infty$  otherwise. Problem  $(\mathcal{P}_{\theta_0})$  can be viewed as a limit of  $(\mathcal{P}_\theta)$  as  $\mu \rightarrow 0^+$ .

At this stage, it is worth mentioning that though we focus here, for simplicity of exposition, on the squared loss  $x \mapsto \frac{1}{2} \|y - \Phi x\|^2$ , our results generalize to more general smooth losses, see Section III-E for further details.

The goal of this paper is to assess the recovery performance of  $(\mathcal{P}_\theta)$ , i.e. to understand how close are the properties of the recovered solution of  $(\mathcal{P}_\theta)$  to those of  $x_0$ . More precisely,

we focus here on the low-noise regime, i.e. when  $\varepsilon$  is small enough, and we investigate stability in  $\ell^2$  sense, but also, and more importantly, the identifiability of the correct low-dimensional manifold associated to  $x_0$ . This unifies and extends a large body of literature, including sparsity and low-rank regularization, which turn to be a very special case of the powerful theory of partly smooth regularization.

## B. Notations

We recall some basic ingredients from differential geometry and convex analysis that are essential to our exposition. For a function  $J$ ,  $\text{dom}(J) = \{x \in \mathbb{R}^p ; J(x) < +\infty\}$  is its domain. We denote  $\partial J(x)$  the subdifferential at  $x$  of the proper, lsc and convex function  $J$ . Geometrically, when  $x \in \text{dom}(J)$ ,  $\partial J(x)$  (if non-empty) is the set of gradients of the affine minorants of  $J$  supporting it at  $x$ . The subdifferential  $\partial J(x)$  is a closed convex set. We denote  $\text{ri}(\mathcal{C})$  (resp.  $\text{rbd}(\mathcal{C})$ ) the relative interior of  $\mathcal{C}$  (resp. relative boundary), i.e. its interior (boundary) for the topology of its affine hull (the smallest affine space containing  $\mathcal{C}$ ).

A good source on smooth manifold theory is [32]. A set  $\mathcal{M} \subset \mathbb{R}^n$  is a  $C^2$ -smooth manifold around a point  $x \in \mathbb{R}^n$ , if  $x \in \mathcal{M}$  and  $\mathcal{M}$  consists locally around  $x$  of the solutions of some  $C^2$ -smooth equations with linearly independent gradients. In this case, the tangent space of  $\mathcal{M}$  at  $x$  is denoted  $\mathcal{T}_x(\mathcal{M})$ . We define the tangent model subspace as

$$T_x = \text{par}(\partial J(x))^\perp,$$

where  $\text{par}(\mathcal{C}) = \mathbb{R}(\mathcal{C} - \mathcal{C})$  is the subspace parallel to the set  $\mathcal{C} \subset \mathbb{R}^n$ . For a linear space  $T$ , we denote  $P_T$  the orthogonal projection on  $T$  and for a matrix  $\Gamma \in \mathbb{R}^{n \times n}$ ,  $\Gamma_T = P_T \Gamma P_T$ .

We use the abbreviation  $O(a_k, b_k)$  to mean  $O(\max(a_k, b_k))$ .

## II. PARTLY-SMOOTH FUNCTIONS

Toward the goal of studying the recovery guarantees of problem  $(P_\theta)$ , our central assumption will be that  $J$  is a partly smooth function. Partial smoothness of functions was originally defined by [34]. Our definition hereafter specializes it to the case of proper, lsc and convex functions. Rigorously speaking, in the following, one should speak of  $C^2$ -smooth embedded submanifold of  $\mathbb{R}^p$ . Nevertheless, to lighten terminology, we shall state  $C^2$ -manifold. For a smooth manifold  $\mathcal{M}$  around  $x \in \mathcal{M}$ ,  $\mathcal{T}_{x'}(\mathcal{M})$  will denote the tangent space to  $\mathcal{M}$  at any point  $x'$  near  $x$  in  $\mathcal{M}$ .

**Definition 1.** Let  $J$  be a proper, lsc convex function, and  $x \in \mathbb{R}^p$  such that  $\partial J(x) \neq \emptyset$ .  $J$  is partly smooth at  $x$  relative to a set  $\mathcal{M}$  containing  $x$  if

- (i) (Smoothness)  $\mathcal{M}$  is a  $C^2$ -manifold around  $x$  and  $J$  restricted to  $\mathcal{M}$  is  $C^2$  around  $x$ .
- (ii) (Sharpness) The tangent space  $\mathcal{T}_x(\mathcal{M})$  is  $T_x$ .
- (iii) (Continuity) The set-valued mapping  $\partial J$  is continuous at  $x$  relative to  $\mathcal{M}$ .

$J$  is said to be partly smooth relative to a set  $\mathcal{M}$  if  $\mathcal{M}$  is a manifold and  $J$  is partly smooth at each point  $x \in \mathcal{M}$  relative to  $\mathcal{M}$ .  $J$  is said to be locally partly smooth at  $x$  relative to a

set  $\mathcal{M}$  if  $\mathcal{M}$  is a manifold and there exists a neighbourhood  $U$  of  $x$  such that  $J$  is partly smooth at each point  $x' \in \mathcal{M} \cap U$  relative to  $\mathcal{M}$ .

Note that in the previous definition,  $\mathcal{M}$  needs only to be defined locally around  $x$ , and it can be shown to be locally unique thanks to prox-regularity of proper, lsc and convex functions, see [27, Corollary 4.2].

Loosely speaking, a partly smooth function behaves smoothly as we move on the identifiable manifold, and sharply if we move normal to the manifold. Examples showing the importance of properties (i)–(iii) and why their individual lack will cause issue are provided in [26, Section 2.2].

**Remark 1** (Discussion of the properties). Since  $J$  is proper, lsc and convex, it is subdifferentially regular at any point in its domain, and in particular at  $x$ . Therefore, the regularity property [34, Definition 2.7(ii)] is automatically verified. In view of [34, Proposition 2.4(i)–(iii)], the sharpness property (ii) is equivalent to [34, Definition 2.7(iii)]. The continuity property (iii) is equivalent to the fact that  $\partial J$  is inner semicontinuous at  $x$  relative to  $\mathcal{M}$ , that is: for any sequence  $x_k$  in  $\mathcal{M}$  converging to  $x$  and any  $\eta \in \partial J(x)$ , there exists a sequence of subgradients  $\eta_k \in \partial J(x_k)$  converging to  $\eta$ . This equivalent characterization will be essential in the proof of our main result.

## A. Examples in Imaging and Machine Learning

We describe below some popular examples of partly smooth regularizers that are routinely used in machine learning, statistics, signal and image processing. We first expose basic building examples (sparsity, group sparsity) and then show how the machinery of partial smoothness enables a powerful calculus to create new priors (using post-composition with a linear operator, spectral lifting, positive linear combinations and separable priors). It turns out that all these examples are partly smooth functions as has been shown in [54].

a)  $\ell^1$  sparsity.: One of the most popular non-quadratic convex regularization is the  $\ell^1$  norm  $J(x) = \sum_{i=1}^n |x_i|$ , which promotes sparsity. Indeed, it is easy to check that  $J$  is partly smooth at  $x$  relative to the subspace

$$\mathcal{M} = T_x = \{u \in \mathbb{R}^n ; \text{supp}(u) \subseteq \text{supp}(x)\}.$$

For any  $x \in \mathbb{R}^n$ ,  $\mathcal{M}$  is a linear subspace, which is obviously a  $C^2$ -manifold. Moreover, on a neighborhood of  $x$  in  $\mathcal{M}$ , the  $\ell^1$  norm is locally linear and thus  $C^2$ . These two facts prove property (i). As far as property (ii) is concerned, since again  $\mathcal{M}$  is a linear subspace, its tangent space  $\mathcal{T}_x(\mathcal{M})$  is nothing but  $T_x$ . Finally, the subdifferential of the  $\ell^1$  norm is a constant set locally around  $x$  along  $\mathcal{M}$ , which in turns shows property (iii). The use of sparse regularization has been popularized in the signal processing literature under the name of basis pursuit [13], and in the statistics literature under the name of Lasso [52].

b)  $\ell^1 - \ell^2$  group sparsity.: To better capture the sparsity pattern of natural signals and images, it is useful to structure the sparsity into non-overlapping blocks/groups  $\mathcal{B}$  such that  $\bigcup_{b \in \mathcal{B}} b = \{1, \dots, n\}$ . This group structure is enforced by

using typically the mixed  $\ell^1 - \ell^2$  norm  $J(x) = \sum_{b \in \mathcal{B}} \|x_b\|$ , where  $x_b = (x_i)_{i \in b} \in \mathbb{R}^{|b|}$ . We refer to [58, 3] and references therein for more details. Unlike the  $\ell^1$  norm, and except the case  $|b| = 1$ , the  $\ell^1 - \ell^2$  norm is not polyhedral, but can be still be shown to be partly smooth at  $x$  relative to the linear manifold

$$\mathcal{M} = T_x = \{x' ; \text{supp}_{\mathcal{B}}(x') \subseteq \text{supp}_{\mathcal{B}}(x)\},$$

where

$$\text{supp}_{\mathcal{B}}(x) = \bigcup \{b ; x_b \neq 0\}.$$

See [54].

*c) Spectral functions.*: The natural spectral extension of sparsity to matrix-valued data  $x \in \mathbb{R}^{n_0 \times n_0}$  (where  $n = n_0^2$ ) is to impose a low-rank prior, which should be understood as sparsity of the singular values. Denote  $x = U_x \text{diag}(\Lambda_x) V_x^*$  an SVD decomposition of  $x$ , where  $\Lambda_x \in \mathbb{R}_+^{n_0}$ . The nuclear norm is defined as

$$J(x) = \|x\|_* = \|\Lambda_x\|_1. \quad (2)$$

It has been used for instance in machine learning applications [3], matrix completion [45, 7] and phase retrieval [11]. The nuclear norm can be shown to be partly smooth at  $x$  relative to the manifold [36, Example 2]

$$\mathcal{M} = \{x' ; \text{rank}(x') = \text{rank}(x)\}. \quad (3)$$

More generally, if  $j : \mathbb{R}^{n_0} \rightarrow \mathbb{R}$  is a permutation-invariant closed convex function, then one can consider the function  $J(x) = j(\Lambda_x)$  which can be shown to be a convex function as well [35]. When restricted to the linear space of symmetric matrices,  $j$  is partly smooth at  $\Lambda_x$  for a manifold  $m_{\Lambda_x}$ , if and only if  $J$  is partly smooth at  $x$  relative to the manifold

$$\mathcal{M} = \{U \text{diag}(\Lambda) U^* ; \Lambda \in m_{\Lambda_x}, U \in \mathcal{O}_{n_0}\},$$

where  $\mathcal{O}_{n_0} \subset \mathbb{R}^{n_0 \times n_0}$  is the group of orthogonal matrices. This result is proved in [14, Theorem 3.19], building upon the work of [15] on manifold smoothness transfer under spectral lifting. This result can be extended to non-symmetric (possibly rectangular) matrices by requiring that  $j$  is an absolutely permutation-invariant closed convex function, see [14, Theorem 5.3]. The nuclear norm  $\|\cdot\|_*$  is a special case where  $j(\Lambda) = \|\Lambda\|_1$ .

*d) Analysis regularizers.*: If  $J_0 : \mathbb{R}^q \rightarrow \mathbb{R} \cup \{+\infty\}$  is a proper lsc convex function and  $D \in \mathbb{R}^{n \times q}$  is a linear operator, an analysis regularizer (following the terminology introduced in [18]) is of the form

$$J(x) = J_0(D^*x).$$

Such a prior controls the low-complexity (as measured by  $J_0$ ) of the correlations between the columns of  $D$  and  $x$ . A popular example is when taking  $J_0 = \|\cdot\|_1$  and  $D^*$  a finite-difference approximation of the gradient of an image. This defines the (anisotropic) total variation, which promotes piecewise constant images, and is popular in image processing [49]. The fused Lasso [53] corresponds to  $J_0$  being the  $\ell^1$ -norm and  $D^*$  is the concatenation of the identity and a finite-difference operator. To cope with correlated covariates in linear

regression, it was devised in [25, 46] to use a family of analysis-type priors where  $J_0 = \|\cdot\|_*$  is the nuclear norm.

If  $J_0$  is partly smooth at  $\alpha = D^*x$  for the manifold  $\mathcal{M}_\alpha^0$ , then it is shown in [34, Theorem 4.2] that  $J$  is partly smooth at  $x$  relative to the manifold

$$\mathcal{M} = \{x' \in \mathbb{R}^n ; D^*x' \in \mathcal{M}_\alpha^0\}.$$

provided that the following transversality condition holds [32, Theorem 6.30(a)]

$$\text{Ker}(D) \cap \mathcal{T}_\alpha(\mathcal{M}_\alpha^0)^\perp = \{0\} \iff \text{Im}(D^*) + \mathcal{T}_\alpha(\mathcal{M}_\alpha^0) = \mathbb{R}^N.$$

Moreover, the co-dimension of  $\mathcal{M}$  in  $\mathbb{R}^p$  equals the co-dimension of  $\mathcal{M}_\alpha^0$  in  $\mathbb{R}^q$ .

*e) Mixed regularization.*: Starting from a family of proper, lsc and convex functions  $\{J_\ell\}_{\ell \in \mathcal{L}}$ ,  $\mathcal{L} = \{1, \dots, L\}$ , it is possible to design a convex function as  $J(x) = \sum_{\ell \in \mathcal{L}} \rho_\ell J_\ell(x)$ , where  $\rho_\ell > 0$  are weights. A popular example is to impose both sparsity and low rank of a matrix, by using  $J_1 = \|\cdot\|_1$  and  $J_2 = \|\cdot\|_*$ , see for instance [22, 42].

Suppose that  $\bigcap_{\ell \in \mathcal{L}} \text{ri}(\text{dom}(J_\ell)) \neq \emptyset$ . Let  $\mathcal{S} \subseteq \mathbb{R}^p$  be a  $C^2$ -manifold. If each  $J_\ell$  is partly smooth at  $x$  relative to a manifold  $\mathcal{M}^\ell \subseteq \mathcal{S}$ , then it can be shown that  $J$  is also partly smooth at  $x$  for

$$\mathcal{M} = \bigcap_{\ell \in \mathcal{L}} \mathcal{M}^\ell,$$

with the proviso that the manifolds  $\mathcal{M}^\ell$  intersect transversally, i.e.

$$\sum_{\ell \in \mathcal{L}} z_\ell = 0$$

and

$$\forall \ell \in \mathcal{L}, z_\ell \in \mathcal{T}_x(\mathcal{M}^\ell)^\perp \Rightarrow z_\ell \in \mathcal{T}_x(\mathcal{S})^\perp \text{ for each } \ell \in \mathcal{L}.$$

Moreover, the co-dimension of  $\mathcal{M}$  (in  $\mathcal{S}$ ) equals the sum of the co-dimensions of  $\mathcal{M}^\ell$ . This assertion is a weaker version of [34, Corollary 4.8], since we use convexity and closedness of the functions  $J_\ell$ . For the case where  $\mathcal{L} = 2$ , the above transversality condition reads [32, Theorem 6.30(b)]

$$\begin{aligned} \mathcal{T}_x(\mathcal{M}_1)^\perp \cap \mathcal{T}_x(\mathcal{M}_2)^\perp &= \mathcal{T}_x(\mathcal{S})^\perp \\ \iff \mathcal{T}_x(\mathcal{M}_1) + \mathcal{T}_x(\mathcal{M}_2) &= \mathcal{T}_x(\mathcal{S}). \end{aligned} \quad (4)$$

*f) Separable Regularization.*: Let  $\{J_\ell\}_{\ell \in \mathcal{L}}$ ,  $\mathcal{L} = \{1, \dots, L\}$ , be a family of proper lsc convex functions. If  $J_\ell$  is partly smooth at  $x_\ell$  relative to a manifold  $\mathcal{M}_{x_\ell}^\ell$ , then the separable function

$$J(\{x_\ell\}_{\ell \in \mathcal{L}}) = \sum_{\ell \in \mathcal{L}} J_\ell(x_\ell)$$

is partly smooth at  $(x_1, \dots, x_L)$  relative to  $\mathcal{M}_{x_1}^1 \times \dots \times \mathcal{M}_{x_L}^L$  [34, Proposition 4.5].

One fundamental problem that has attracted a lot of interest in the recent years in data processing involves decomposing an observed object into a linear combination of components/constituents  $x_\ell$ ,  $\ell \in \mathcal{L}$ . One instance of such a problem is image decomposition into texture and piece-wise-smooth (cartoon) parts, see e.g. [51, 1, 43] and references therein. Another example of decomposition is principal component

pursuit, proposed in [8], to decompose a matrix which is the superposition of a low-rank component and a sparse component. In this case  $J_1 = \|\cdot\|_1$  and  $J_2 = \|\cdot\|_*$ .

### III. MAIN RESULTS

In the following, we denote  $T = T_{x_0}$ ,  $e = P_T(\partial J(x_0)) \in \mathbb{R}^n$ .

#### A. Linearized pre-certificate and minimal norm dual certificate

Before stating our main contributions, we first introduce a central object of this paper, which controls the stability of  $\mathcal{M}$  when the signal to noise ratio is large enough.

**Definition 2** (Linearized pre-certificate). *For some matrix  $\Gamma \in \mathbb{R}^{n \times n}$ , assuming  $\ker(\Gamma) \cap T = \{0\}$ , we define  $\eta_\Gamma = \Gamma \Gamma_T^+ e$ .*

Recall that  $J$  is proper lsc convex function, and we suppose that  $\text{Im}(\tilde{\Gamma}) \cap \partial J(x_0) \neq \emptyset$  (so-called range or source condition in the inverse problem community). The latter is equivalent to the fact that  $x_0$  is a minimizer of  $(\mathcal{P}_{\theta_0})$ . This is straightforward to see by writing the first-order optimality condition of this convex program.

**Definition 3** (Minimal norm certificate). *The minimal norm certificate is the vector*

$$\hat{\eta}_{\tilde{\Gamma}} = \tilde{\Gamma} \hat{z}_{\tilde{\Gamma}}, \quad \text{where } \hat{z}_{\tilde{\Gamma}} = \underset{\tilde{\Gamma} z \in \partial J(x_0)}{\text{argmin}} \|z\|. \quad (5)$$

This certificate is uniquely defined as the constraint set is non-empty closed and convex, and the solution of the minimization problem, which is the projection of the origin on it, is obviously unique.

**Proposition 1.** *Assume that  $\ker(\tilde{\Gamma}) \cap T = \{0\}$ . Then,*

$$\eta_{\tilde{\Gamma}} \in \text{ri}(\partial J(x_0)) \implies \hat{\eta}_{\tilde{\Gamma}} = \eta_{\tilde{\Gamma}}, \quad (6)$$

$$\hat{\eta}_{\tilde{\Gamma}} \in \text{ri}(\partial J(x_0)) \implies \eta_{\tilde{\Gamma}} = \hat{\eta}_{\tilde{\Gamma}}. \quad (7)$$

*Under either of these conditions,  $x_0$  is the unique minimizer to  $(\mathcal{P}_{0, \tilde{\Gamma} x_0, \tilde{\Gamma}})$ .*

The proof is postponed to Section V-B.

#### B. Deterministic model consistency

We first consider the case where  $\Phi$  and  $w$  (or equivalently  $\Gamma$  and  $u$ ) are fixed and deterministic. Our main contribution is the following *model consistency* theorem, which shows the robustness of the manifold  $\mathcal{M}$  associated to  $x_0$  to small perturbations on both the observations and the design matrix, provided that  $\mu$  (or equivalently  $\lambda$ ) is well chosen. As a product, we also get  $\ell^2$  stability.

**Theorem 1.** *Assume that  $J$  is locally partly smooth at  $x_0$  relative to  $\mathcal{M}$  and that there exists  $\tilde{\Gamma} \in \mathbb{R}^{n \times n}$  such that*

$$\ker(\tilde{\Gamma}) \cap T = \{0\}, \quad \text{and} \quad \eta_{\tilde{\Gamma}} \in \text{ri}(\partial J(x_0)). \quad (8)$$

*Then, there exists a constant  $C > 0$  such that if*

$$\max\left(\|\Gamma - \tilde{\Gamma}\|, \|\varepsilon\| \mu^{-1}, \mu\right) \leq C, \quad (9)$$

*the solution  $x_\theta$  of  $(\mathcal{P}_\theta)$  is unique and satisfies*

$$x_\theta \in \mathcal{M} \quad \text{and} \quad \|x_\theta - x_0\| = O(\|\varepsilon\|). \quad (10)$$

This theorem is proved in Section V-C.

**Remark 2** (Stability constants). *Observe that the non-degeneracy and restricted injectivity conditions in (8) can be viewed as a geometric generalization of the so-called irrepresentable condition in statistics (see Section III-F for further details). They guarantee in particular that  $x_0$  is identifiable in the exact case, i.e. a unique solution to  $(\mathcal{P}_{0, \Gamma x_0, \Gamma})$ . Theorem 1 ensures that under (8), solving  $(\mathcal{P}_\theta)$  indeed recovers a unique solution  $x_\theta$  having the correct model (i.e.  $x_\theta \in \mathcal{M}$ ). In our way of proving model consistency, we also get  $\ell^2$  stability to such a small noise. When the underlying model has a low complexity (typically  $\mathcal{M}$  has a small dimension), this means that the recovery will be highly stable, which is reflected both in the constant  $C$  in (9) and the (local) Lipschitz constant hidden behind the bound  $\|x_\theta - x_0\| = O(\|\varepsilon\|)$ . Obtaining sharp estimates of these constants for general low complexity models (manifolds) is rather challenging and will necessitate even more involved arguments from differential geometry. The case of regularizers  $J$  where the partial smoothness manifold is affine was deeply investigated in [54]. There, the authors derived explicit (though quite involved) formulae for  $C$  and the Lipschitz constant. It is however easier to see how these constants behave in the small noise limit, i.e. when  $\|\varepsilon\| \rightarrow 0$  with  $\mu = C_0 \|\varepsilon\|$  for some large enough constant  $C_0$  ( $C_0 \geq 1/C$  from (9)), and where we assume for simplicity that  $\Gamma = \tilde{\Gamma}$  for simplicity. Indeed, as shown in the proof of Theorem 1 (see in particular (30)), for enough small noise level, the recovery error behaves as*

$$\|x_\theta - x_0\| \sim \|\tilde{\Gamma}_T^+(\varepsilon - \mu e_{x_0})\| \leq (1 + C_0 \|e_{x_0}\|) \|\tilde{\Gamma}_T^+\| \|\varepsilon\|.$$

*The operator  $\tilde{\Gamma}_T^+$  amounts to inverting the restriction of the operator  $\tilde{\Gamma}$  to the low dimensional subspace  $T$ . The constant  $\|\tilde{\Gamma}_T^+\|$  thus captures the complexity of the model at  $x_0$ , and this term increases as the dimension of  $\mathcal{M}$  (equal to that of  $T$  by the sharpness property ((ii))) increases.*

**Remark 3** ("Distance" to degeneracy). *It is worth emphasizing that the "distance" of  $\eta_{\tilde{\Gamma}}$  to degeneracy affects both stability of the model selection and  $\ell^2$  stability. Again, this can be quantified precisely for linear manifolds as in [54]. In the general case, it is much more difficult. Nevertheless, we can still give some clues in the low-noise regime. Let  $\eta_\theta = \frac{u - \Gamma x_\theta}{\mu}$ . From the first-order condition in the proof of Theorem 1, we have  $\eta_\theta \in \text{Aff}(\partial J(x_\theta))$ . Thus*

$$\eta_\theta \in \text{ri}(\partial J(x_\theta)) \iff \text{dist}(P_{\partial J(x_\theta)}(\eta_\theta), \text{rbd}(\partial J(x_\theta))) > 0.$$

*The triangle inequality, (8), convexity and closedness of  $\partial J(x)$*

as well as its continuity along  $\mathcal{M}$  then yield

$$\begin{aligned} & \text{dist}(\eta_{\tilde{\Gamma}}, \text{rbd}(\partial J(x_\theta))) \\ & \leq \text{dist}(\mathbb{P}_{\partial J(x_\theta)}(\eta_\theta), \text{rbd}(\partial J(x_\theta))) \\ & \quad + \|\mathbb{P}_{\partial J(x_\theta)}(\eta_\theta) - \mathbb{P}_{\partial J(x_\theta)}(\eta_{\tilde{\Gamma}})\| \\ & \quad + \|\mathbb{P}_{\partial J(x_\theta)}(\eta_{\tilde{\Gamma}}) - \mathbb{P}_{\partial J(x_0)}(\eta_{\tilde{\Gamma}})\| \\ & \leq \text{dist}(\mathbb{P}_{\partial J(x_\theta)}(\eta_\theta), \text{rbd}(\partial J(x_\theta))) + \|\eta_\theta - \eta_{\tilde{\Gamma}}\| \\ & \quad + c\|x_\theta - x_0\|\|\eta_{\tilde{\Gamma}}\| \end{aligned}$$

for  $c \geq 0$ . To simplify the discussion, we suppose that  $\text{rbd}(\partial J(x_\theta)) \subset \text{rbd}(\partial J(x_0))$  for small enough noise (this is true for instance if  $J$  is locally polyhedral around  $x_0$ ). This implies  $\text{dist}(\eta_{\tilde{\Gamma}}, \text{rbd}(\partial J(x_0))) \leq \text{dist}(\eta_{\tilde{\Gamma}}, \text{rbd}(\partial J(x_\theta)))$ . Consequently, to have  $\eta_\theta \in \text{ri}(\partial J(x_\theta))$ , as required to show uniqueness of  $x_\theta$  in Theorem 1, it is sufficient that

$$\|\eta_\theta - \eta_{\tilde{\Gamma}}\| + c\|x_\theta - x_0\|\|\eta_{\tilde{\Gamma}}\| < \text{dist}(\eta_{\tilde{\Gamma}}, \text{rbd}(\partial J(x_0)))$$

which is valid since  $\text{dist}(\eta_{\tilde{\Gamma}}, \text{rbd}(\partial J(x_0))) > 0$  owing to (8). In view of (31) and (35), this holds true if

$$C_1 \frac{\|\varepsilon\|}{\text{dist}(\eta_{\tilde{\Gamma}}, \text{rbd}(\partial J(x_0)))} \leq \mu \leq C_2 \text{dist}(\eta_{\tilde{\Gamma}}, \text{rbd}(\partial J(x_0)))$$

and

$$\|\Gamma - \tilde{\Gamma}\| \leq C_3 \text{dist}(\eta_{\tilde{\Gamma}}, \text{rbd}(\partial J(x_0))),$$

where  $C_1$ ,  $C_2$  and  $C_3$  are positive constants that depend on  $\eta_{\tilde{\Gamma}}$ ,  $c$  and on the constants in the  $O(\cdot)$  terms in (31) and (35). These constants encode again the complexity of  $\mathcal{M}$  and  $x_0$  in them. One can clearly see that the closer  $\eta_{\tilde{\Gamma}}$  to  $\text{rbd}(\partial J(x_0))$ , the smaller the noise that can be tolerated for model stability. Moreover,  $\mu$  must be chosen large enough but not too large. In particular, if  $\mu = C_1 \frac{\|\varepsilon\|}{\text{dist}(\eta_{\tilde{\Gamma}}, \text{rbd}(\partial J(x_0)))}$ , we have from the previous remark that

$$\|x_\theta - x_0\| \sim \left(1 + \frac{C_1 \|e_{x_0}\|}{\text{dist}(\eta_{\tilde{\Gamma}}, \text{rbd}(\partial J(x_0)))}\right) \|\tilde{\Gamma}_T^+ \|\varepsilon\|,$$

which in turn shows the influence of non-degeneracy on  $\ell^2$  stability. The same reasoning can be extended to the scenario where  $J$  is the separable sum of sublinear regularizers, in which case the assumption  $\text{rbd}(\partial J(x_\theta)) \subset \text{rbd}(\partial J(x_0))$  can be removed.

**Remark 4** (Inverse problems). A typical case of application of this result is in inverse problems that are encountered in various disciplines in science and engineering, such as in signal and image processing. In such a setting, the forward operator  $\Phi$  is generally fixed and known, and one then takes  $\tilde{\Gamma} = \Gamma = \Phi^* \Phi / p$ .

**Remark 5** (Uncertain design/forward operator). If only a noisy version of the forward operator (in inverse problems) or the design (in regression) is available then this can also be handled by Theorem 1. This scenario has been considered for sparse recovery (i.e.  $J$  the  $\ell^1$ -norm) by several authors for sparse linear regression and compressed sensing, see e.g. [28, 48, 39].

**Remark 6** (Random setting). In statistics or machine learning, one considers a regression problem where the design  $\Phi$  and the noise  $w$  are random, under the asymptotic regime

where the number of observations  $p$ , i.e. number of rows  $\Phi$ , grows large, so that  $\Gamma$  only reach  $\tilde{\Gamma}$  in the limit  $p \rightarrow +\infty$ . See Theorem 2 for details.

**Remark 7** (Identification of the manifold). Theorem 1 guarantees that, under some hypotheses on  $x_0$  and  $\theta$ ,  $x_\theta$  belongs to  $\mathcal{M}$ . For all the regularizations considered in Section II-A, it turns out that actually  $\mathcal{M}_{x_\theta} = \mathcal{M}$ . This is because, for any  $(x, x')$  with  $x' \in \mathcal{M}_x$  close enough to  $x$ , one has  $\mathcal{M}_{x'} = \mathcal{M}_x$ .

The following proposition, proved in Section V-F, shows that Theorem 1 is in some sense sharp, since the hypothesis  $\eta_\Gamma \in \text{ri}(\partial J(x_0))$  (almost) characterizes the stability of  $\mathcal{M}$ .

**Proposition 2.** Suppose that  $x_0$  is the unique solution of  $(\mathcal{P}_{(0, \tilde{\Gamma}_{x_0, \tilde{\Gamma}})})$  and that

$$\ker(\tilde{\Gamma}) \cap T = \{0\}, \quad \text{and} \quad \eta_{\tilde{\Gamma}} \notin \partial J(x_0). \quad (11)$$

Then there exists  $C > 0$  such that if (9) holds, then any solution  $x_\theta$  of  $(\mathcal{P}_\theta)$  for  $\mu > 0$  satisfies  $x_\theta \notin \mathcal{M}$ .

In the particular case where  $\varepsilon = 0$  (no noise) and  $\tilde{\Gamma} = \Gamma$ , this result shows that the manifold  $\mathcal{M}$  cannot be correctly identified by solving  $(\mathcal{P}_{(\mu, \Gamma_{x_0, \Gamma})})$  for any  $\mu > 0$  small enough.

**Remark 8** (Critical case). The only case not covered by either Theorem 1 or Proposition 2 is when  $\eta_{\tilde{\Gamma}} \in \text{rbd}(\partial J(x_0))$ . In this case, one cannot conclude in general, since depending on the noise  $w$ , one can have either stability or non-stability of  $\mathcal{M}$ . We refer to [55] where an example illustrates this situation for the 1-D total variation  $J = \|D_{\text{DIF}}^* \cdot\|_1$  (here  $D_{\text{DIF}}^*$  is a discretization of the 1-D derivative operator).

### C. Probabilistic model consistency

We now turn to study consistency of our estimator. In this section, we work under the classical setting where  $n$  and  $x_0$  are fixed as the number of observations  $p \rightarrow \infty$ . We consider that the design matrix and the noise are random. More precisely, the data  $(\varphi_i, w_i)$  are random vectors in  $\mathbb{R}^n \times \mathbb{R}$ ,  $i = 1, \dots, n$ , where  $\varphi_i$  is the  $i$ -th row of  $\Phi$ , are assumed independent and identically distributed (i.i.d.) samples from a joint probability distribution such that  $\mathbb{E}(w_i | \varphi_i) = 0$ , finite fourth-order moments, i.e.  $\mathbb{E}(w_i^4) < +\infty$  and  $\mathbb{E}(\|\varphi_i\|^4) < +\infty$ . Note that in general,  $w_i$  and  $\varphi_i$  are not necessarily independent. It is possible to extend our result to other distribution models by weakening some of the assumptions and strengthening others, see e.g. [31, 59, 3]. Let's denote  $\tilde{\Gamma} = \mathbb{E}(\xi^* \xi) \in \mathbb{R}^{n \times n}$ , where  $\xi$  is any row of  $\Phi$ . We do not make any assumption on invertibility of  $\tilde{\Gamma}$ .

To make the discussion clearer, the canonical parameters  $\theta$  will be indexed by  $p$ . The estimator  $x_{\theta_p}$  obtained by solving  $(\mathcal{P}_{\theta_p})$  for a sequence  $\theta_p$  is said to be consistent for  $x_0$  if,  $\lim_{p \rightarrow +\infty} \Pr(x_{\theta_p} \text{ is unique}) \rightarrow 1$  and  $x_{\theta_p}$  converges to  $x_0$  in probability. The estimator is said to be model consistent if  $\lim_{p \rightarrow +\infty} \Pr(x_{\theta_p} \in \mathcal{M}) \rightarrow 1$ , where  $\mathcal{M}$  is the manifold associated to  $x_0$ .

The following result ensures model consistency for certain scaling of  $\mu_p$ . It is proved in Section V-E

**Theorem 2.** *If conditions (8) hold and*

$$\mu_p = o(1) \quad \text{and} \quad \mu_p^{-1} = o(p^{1/2}). \quad (12)$$

*Then the estimator  $x_{\theta_p}$  of  $x_0$  obtained by solving  $(\mathcal{P}_{\theta_p})$  is model consistent.*

**Remark 9** (Sharpness of the criterion). *One can also state a probabilistic equivalent to Proposition 2. That is, if  $x_0$  is the unique solution of  $(\mathcal{P}_{0, \tilde{\Gamma}x_0, \tilde{\Gamma}})$ , and conditions (11) and (12) hold, then the estimator  $x_{\theta_p}$  of  $x_0$  defined by solving  $(\mathcal{P}_{\theta_n})$  cannot be model consistent.*

#### D. Algorithmic Implications

A popular iterative scheme to compute a solution of  $(\mathcal{P}_\theta)$  is the Forward-Backward (F-B) splitting algorithm. A comprehensive treatment of the convergence properties of this algorithm, and other proximal splitting schemes, can be found in the monograph [4]. Starting from some  $x^0 \in \mathbb{R}^n$ , the algorithm implements the following iteration

$$x^{k+1} = \text{Prox}_{\tau_k \mu J} (x^k + \tau_k (u - \Gamma x^k)), \quad (13)$$

where the step size satisfies  $0 < \underline{\tau} \leq \tau_k \leq \bar{\tau} < 2/\|\Gamma\|$ , and the proximity operator is defined, for  $\gamma > 0$ , as

$$\text{Prox}_{\gamma J}(x) = \underset{x' \in \mathbb{R}^n}{\text{argmin}} \frac{1}{2} \|x - x'\|^2 + \gamma J(x'). \quad (14)$$

The following theorem shows that the F-B algorithm correctly identifies the manifold  $\mathcal{M}$  after a finite number of iterations.

**Theorem 3.** *Suppose that the assumptions of Theorem 1 hold. Then, there exists  $k_0$  large enough, such that for all  $k \geq k_0$ , the F-B iterates satisfy  $x^k \in \mathcal{M}$ .*

*Proof.* Inspection of the proof of Theorem 1 shows that the solution  $x_\theta$  of  $(\mathcal{P}_\theta)$ , which is unique, is such that the vector  $\eta_\theta = \frac{u - \Gamma x_\theta}{\lambda}$  satisfies  $\eta_\theta \in \text{ri}(\partial J(x_\theta))$  when (8) and (9) hold. Moreover, as  $x_\theta \in \mathcal{M}$ ,  $x_\theta$  is near  $x_0$ , and  $J$  is locally partly smooth at  $x_0$  relative to  $\mathcal{M}$ , it is also partly smooth at  $x_\theta$  relative to the same manifold  $\mathcal{M}$ . Altogether, this implies that the assumptions of [38, Theorem 3.1] are fulfilled and the manifold identification claim follows.  $\square$

This result sheds light on the convergence behaviour of this algorithm in the favourable case where condition (8) holds and  $(\|\Gamma - \tilde{\Gamma}\|, \|\varepsilon\|/\mu, \mu)$  are sufficiently small.

#### E. General Loss Functions

For the sake of simplicity, we have described our contributions with the squared loss function  $u \in \mathbb{R}^p \mapsto \frac{1}{2} \|y - u\|^2$ . Our results, however, extend readily to the case of more general loss functions of the form  $F(u, y)$ . In the following  $\nabla_1^2 F(u, y)$  denotes the Hessian of  $F$  with respect to the first variable evaluated at  $(u, y)$ .

We thus consider the variational problem

$$\min_x F(\Phi x, y) + \lambda J(x),$$

where the loss function  $F : \mathbb{R}^p \times \mathbb{R}^p \rightarrow \mathbb{R}$  fulfills the following assumptions:

**(A.1)** For any  $y \in \mathbb{R}^p$ ,  $F(\cdot, y) \in C^2(\mathbb{R}^p)$  is  $\sigma_m$ -strongly convex and  $\nabla_1 F(\cdot, y)$  is  $\sigma_M$ -Lipschitz continuous,  $\sigma_m > 0$  and  $\sigma_M > 0$ .

**(A.2)** The gradient of  $F$  with respect to the first variable,  $\nabla_1 F(u, y)$ , is such that  $\nabla_1 F(u, u) = 0$ .

Any loss function of the form  $F(u, y) = G(u) - \langle u, y \rangle$ , where  $G$  is a  $C^2$  strongly convex function and its gradient is Lipschitz-continuous, satisfies assumptions **(A.1)**. Assumption **(A.2)** is quite natural for a data fidelity term, and is fulfilled for instance for some losses in the exponential family.

In this setting, Theorem 1 (and in a similar way our other contributions) remains valid, and one simply needs to replace condition (9) by

$$\max \left( \|\tilde{\Gamma} - \tilde{\Gamma}\|, \|\varepsilon\| \mu^{-1}, \mu \right) \leq C, \quad (15)$$

where now

$$\tilde{\Gamma} = \frac{1}{p} \Phi^* \nabla_1^2 F(y, y) \Phi \quad \text{and} \quad \varepsilon = \frac{1}{p} \Phi^* \nabla_1^2 F(y, y) w,$$

where  $\nabla_1^2 F(y, y)$  is the Hessian with respect to the first variable (assumed to be positive definite by assumption **(A.1)**) taken at  $(y, y)$ . A detailed treatment on the way to adapt the proofs to handle such a generic loss is provided in Section V-D.

#### F. Relation to Previous Works

a) *Works on linear convergence rates.*: Following the pioneer work [6] (who study convergence in terms of Bregman divergence), there is a large amount of works on the study conditions under which  $\|x_\theta - x_0\| = O(\|\varepsilon\|)$  (so-called linear convergence rate) where  $x_\theta$  is any solution of  $(\mathcal{P}_\theta)$ , see for instance the book [50] for an overview of these results. The initial work of [24] proves a sharp criteria to ensure linear convergence rate for the  $\ell^1$  norm, and this approach is further extended to arbitrary convex functions by [23] and [19], who respectively proved linear convergence rates in terms of the penalty  $J$  and  $\ell^2$ -norm.

These works show that if

$$\ker(\Gamma) \cap T = \{0\} \quad \text{and} \quad \exists \eta \in \text{Im}(\Gamma) \cap \text{ri}(\partial J(x_0)) \quad (16)$$

(which is often called the source condition), then linear convergence rate holds. Note that condition (8) implies (16), but it is stronger. Indeed, condition (16) does not ensure model consistency (10), which is a stronger requirement. Model consistency requires, as we show in our work, the use of a special certificate, the minimal norm certificate  $\hat{\eta}_\Gamma$ , which is equal to  $\eta_\Gamma$  if  $\eta_\Gamma \in \text{ri}(\partial J(x_0))$  (see Proposition 1).

b) *Works on model consistency.*: Theorem 1 is a generalization of a large body of results in the literature. For the Lasso, i.e.  $J = \|\cdot\|_1$ , and when  $\Gamma = \tilde{\Gamma}$ , to the best of our knowledge, this result was initially stated in [21]. In this setting, the result (10) corresponds to the correct identification of the support, i.e.  $\text{supp}(x_\theta) = \text{supp}(x_0)$ . Condition (8) for  $J = \|\cdot\|_1$  is known in the statistics literature under the name ‘‘irrepresentable condition’’, see e.g. [59]. [31] have shown estimation consistency for Lasso for fixed  $n$  and  $x_0$

and asymptotic normality of the estimates. The authors in [59] proved Theorem 2 for  $J = \|\cdot\|_1$ , though under slightly different assumptions on the covariance and noise distribution. A similar result was established in [30] for the elastic net, i.e.  $J = \|\cdot\|_1 + \rho\|\cdot\|_2^2$  for  $\rho > 0$ . In [2] and [3], the author has shown Theorem 2 for two special cases, namely the group Lasso nuclear/trace norm minimization, under a specialization of (8) to these two penalties and in an asymptotic setting. Note that these previous works assume that the asymptotic covariance  $\tilde{\Gamma}$  is invertible. We do not impose such an assumption, and only require the weaker restricted injectivity condition  $\ker(\tilde{\Gamma}) \cap T = \{0\}$ . In a previous work [55], we have proved an instance of Theorem 1 when  $\Gamma = \tilde{\Gamma}$  and  $J(x) = \|D^*x\|_1$ , where  $D \in \mathbb{R}^{n \times q}$  is an arbitrary linear operator. This covers as special cases the discrete anisotropic total variation or the fused Lasso. This result was further generalized in [54] when  $\Gamma = \tilde{\Gamma}$ , and  $J$  belongs to the class of partly smooth functions relative to affine manifolds  $\mathcal{M}$ , i.e.  $\mathcal{M} = x + T_x$ . Typical instances encompassed in this class are the  $\ell^1 - \ell^2$  norm, or its analysis version, as well as non-negative polyhedral functions including the  $\ell^\infty$  norm. Note that the nuclear norm (and composition of it with linear operators as studied for instance in [25, 46]), whose manifold is not affine, does not fit into the framework of [54], while it is covered by Theorem 1. [29] investigated a class of geometrically decomposable penalties, for which they formulated the irrepresentable condition and used it to establish  $\ell^2$ -consistency and model consistency. This class of penalties turns out to be a very special case of ours. When the noise is i.i.f. zero-mean subgaussian, These authors also derived rank consistency of the nuclear norm in a high-dimensional setting which can be viewed as non-asymptotic form of [3]. Lastly, a similar result was proved in [17] for an infinite dimensional sparse recovery problem over space of Radon measures, when  $J$  is the total variation of a measure (not to be confused with the total variation semi-norm mentioned above). In this setting, an interesting finding is that, when  $\hat{\eta}_{\Phi^* \Phi} \in \text{ri}(\partial J(x_0))$ ,  $\hat{\eta}_{\Phi^* \Phi}$  is not equal to  $\eta_{\Phi^* \Phi}$  but to a different certificate (called ‘‘vanishing derivative’’ certificate by [17]) that can also be computed by solving a linear system.

*c) Compressed sensing:* Condition (8) is often used when  $\Phi$  is drawn from the Gaussian matrix ensemble to asses the performance of compressed sensing recovery with  $\ell^1$  norm, see [56, 16]. This is extended to a more general family of decomposable norms (including in particular  $\ell^1 - \ell^2$  norms and the nuclear norm) in [10], but only in the noiseless setting. Our result shows that this analysis extends to the noisy setting as well, and ensures model consistency at high signal to low noise levels. The same condition is used to asses the performance of matrix completion (i.e. the operator  $\Phi$  is a random masking operator) in a noiseless setting [7, 12]. It was also used to ensure  $\ell^2$  robustness of matrix completion in a noisy setting [9], and our findings shows that these results also ensure rank consistency for matrix completion at high signal to low noise levels.

*d) Sensitivity analysis.:* Sensitivity analysis is a central theme in variational analysis. Theorems 1 can be understood as a sensitivity analysis of the minimizers of  $E$  at the point

$(x_0, \theta_0)$ . Classical sensitivity analysis of non-smooth optimization problems seeks conditions to ensure smoothness of the mapping  $\theta \mapsto x_\theta$  where  $x_\theta$  is a minimizer of  $f(\cdot, \theta)$ , see for instance [40, 47, 5].

This is usually guaranteed by the non-degenerate source condition and restricted injectivity condition (16), which, as already reviewed above, ensures linear convergence rate, and hence Lipschitz behaviour of this mapping. The result captured by Theorem 1 goes one step further, by assessing that  $\mathcal{M}_{x_0}$  is a stable manifold (in the sense of [57]), since the minimizer  $x_\theta$  is unique and remains in  $\mathcal{M}_{x_0}$  for  $\theta$  close to  $\theta_0$ . Our starting point for establishing Theorem 1 is the inspiring work of [34] who first introduced the notion of partial smoothness and showed that this broad class of functions enjoys a powerful calculus and sensitivity theory. For convex functions (which is the setting considered in our work), partial smoothness is closely related to  $\mathcal{U} - \mathcal{V}$ -decompositions developed in [33]. In fact, the behaviour of a partly smooth 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 [34, 37]. A important and distinctive feature of Theorem 1 is that, partial smoothness of  $J$  at  $x_0$  relative to  $\mathcal{M}$  transfers to  $E(\cdot, \theta)$  for  $\lambda > 0$ , but not when  $\lambda = 0$  in general. In particular, [34, Theorem 5.7] does not apply to prove our claim.

#### IV. CASE STUDY: NUCLEAR NORM REGULARIZATION

In this section, we illustrate the usefulness of our model consistency results to derive a sharp manifold stability analysis for the nuclear norm (a.k.a trace norm) regularization. As detailed in Section III-F, previous consistency results due to [3] only apply to the overdetermined setting, while our result tackles arbitrary design  $\Phi$  by only requiring the weaker injectivity condition (8). For simplicity of exposition, we consider recovery of square matrices of size  $n = n_0 \times n_0$ , but the same holds for arbitrary rectangular matrices.

##### A. Irrepresentability Criterion IC

The nuclear norm, defined in (2), turns out to be the tightest convex relation of the rank function on the spectral ball. It is then the best convex candidate to enforce a low-rank prior [20]. It is moreover partly smooth at any  $x_0 \in \mathbb{R}^{n_0 \times n_0}$  relative to the manifold  $\mathcal{M}$  of fixed rank  $r = \text{rank}(x_0)$  defined in (3).

Let  $x_0 = U \text{diag}(\sigma(x_0))V^*$  be a reduced rank- $r$  SVD decomposition of  $x_0$ , where  $V, U \in \mathbb{R}^{n_0 \times r}$  have orthonormal columns and  $\sigma(x_0) \in (\mathbb{R}_+^*)^r$  is the vector of singular values of  $x_0$ . The subdifferential of the nuclear norm at  $x_0$  reads (see for instance [10])

$$\partial \|\cdot\|_*(x_0) = \left\{ \eta \in \mathbb{R}^{n_0 \times n_0} ; \eta_T = e \quad \text{and} \quad \|\eta_S\| \leq 1 \right\}, \quad (17)$$

where  $\|\eta\|$  is the operator norm,  $T = \mathcal{T}_x(\mathcal{M})$ ,  $S = T^\perp$  and  $e = P_T(\partial J(x_0))$ , with

$$T = \{UA^* + BV^* ; A, B \in \mathbb{R}^{n_0 \times r}\} \quad \text{and} \quad e = UV^*$$



and  $S$  is the subspace of matrices spanned by the family  $(wz^*)$ , where  $w$  (resp.  $z$ ) is any vector orthogonal to  $U$  (resp.  $V$ ).

The relative interior of  $\partial\|\cdot\|_*(x_0)$  is formed by subgradients  $\eta$  satisfying the inequality in (17) strictly. Thus, condition (8) in the case  $\tilde{\Gamma} = \Gamma$  takes the analytical form

$$\eta_\Gamma \in \text{ri}(\partial J(x_0)) \iff \mathbf{IC}(x_0) < 1, \quad (18)$$

where

$$\mathbf{IC}(x_0) = \|\text{P}_S \Gamma \Gamma_T^\dagger e\|.$$

The value of  $\mathbf{IC}(x_0)$  can then be easily computed. Loosely speaking, the smaller the quantity  $1 - \mathbf{IC}(x_0)$  is, the further  $\eta_\Gamma$  is from the relative boundary of  $\partial J(x_0)$ , and in turn the smaller the stability constant controlling  $\|x_\theta - x_0\|/\|\varepsilon\|$  in (10) is.

### B. Recovery from Gaussian Measurements

Bounding  $\mathbf{IC}$  for an arbitrary operator  $\Phi$  and matrix  $x_0$  is in general difficult. It is however possible to leverage tools from random matrix theory to obtain sharp upper-bounds when  $\Phi$  is drawn from certain matrix ensembles. This strategy has been deployed to study matrix completion problems, see for instance [7, 12]. Another problem on which we now focus is when  $\Phi$  is drawn from the standard Gaussian ensemble, i.e. its entries are independent identically distributed from  $\mathcal{N}(0, 1)$ . The following result, proved by [10], shows that  $\mathbf{IC}(x_0) < 1$  with high probability as soon as  $p$  is larger than  $6rn_0$  (up to negligible terms).

**Proposition 3** ([10], Theorem 1.2). *Let  $x_0 \in \mathbb{R}^{n_0 \times n_0}$  such that  $\text{rank}(x_0) = r$ . If*

$$p \geq \delta r(6n_0 - 5r) \quad (19)$$

for some  $\delta > 1$ , then  $\mathbf{IC}(x_0) < 1$  with probability at least  $1 - 2e^{-(1-\delta)n_0/8}$ .

Combining this result with Theorem 1, this shows that under the scaling (19) of  $(p, n_0, r)$ , one obtains with high probability on the design matrix a rank-consistent estimation of the unknown matrix  $x_0$ , which is (to the best of our knowledge) a novel result.

Figure 1 illustrates this result by computing the average (over 25 Monte Carlo replications) values of  $\mathbf{IC}(x_0)$  for either a varying  $p$  or rank  $r$ . The shaded area corresponds to  $\pm 3 \times$  standard deviation across the 25 replications, and the dashed vertical line indicates the transition predicted by (19). This suggests numerically that the upper-bound (19) is indeed sharp.

### C. Forward-Backward Model Consistency and Unconsistency

As detailed in Section III-D, our theoretical analysis of model consistency also sheds light on the behavior of proximal splitting algorithms, and in particular of the celebrated F-B scheme (13). In the special case  $J = \|\cdot\|_*$  considered in

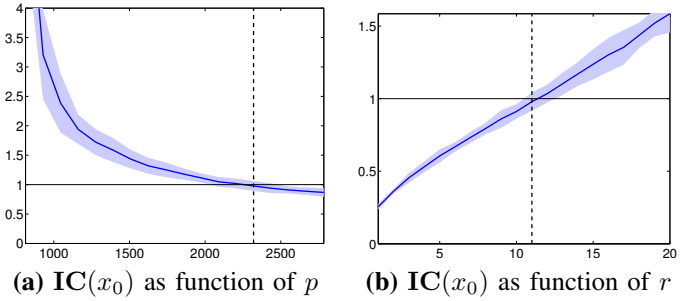


Fig. 1. Curves of  $\mathbf{IC}(x_0)$  (central solid line: average, blue shaded band:  $\pm 3 \times$  standard deviation) for 25 realizations of  $\Phi \in \mathbb{R}^{p \times n_0^2}$  from the standard Gaussian ensemble, where  $n_0 = 10^3$ , and a random  $x_0 = AB^*$  of rank  $r$  where  $A, B \in \mathbb{R}^{n_0 \times r}$  are Gaussian matrices. **(a)**  $\mathbf{IC}(x_0)$  as a function of  $p$  for a fixed  $r = 4$ . The vertical dashed line shows the threshold  $p = r(6n_0 - 5r)$  indicated by (19). **(b)**  $\mathbf{IC}(x_0)$  as a function of  $r$  for a fixed  $p = 0.6n_0^2$ . The vertical dashed line shows the threshold  $r = (3n_0 - \sqrt{9n_0^2 - 5p})/5$  indicated by (19).

this section, the proximal mapping (14) at  $x \in \mathbb{R}^{n_0 \times n_0}$  is computed by simply soft-thresholding the singular values

$$\text{Prox}_{\gamma\|\cdot\|_*}(x) = U \text{diag}(\text{Prox}_{\gamma\|\cdot\|_1}(\sigma(x)))V^*, \quad (20)$$

$$\text{where } \text{Prox}_{\gamma\|\cdot\|_1}(s) = (\text{sign}(s_i) \max(0, |s_i| - \gamma))_i, \quad (21)$$

where  $x = U \text{diag}(\sigma(x))V^*$  is a reduced SVD decomposition of  $x$ .

As in the previous section, we consider here a compressed sensing scenario, where again  $\Phi \in \mathbb{R}^{p \times n_0^2}$  is drawn from the standard Gaussian ensemble, and  $x_0 \in \mathbb{R}^{n_0 \times n_0}$  has low-rank. The observations  $y = \Phi x_0 + w \in \mathbb{R}^p$  are generated with an additive zero-mean white Gaussian noise  $w$  of standard deviation  $10^{-3}\|\Phi x_0\|$  (but the same conclusion holds for a noise of arbitrary small amplitude). We then compute (approximately) a minimizer of (1) using the F-B iterations (13), tuning the regularization parameter  $\mu = C_0\|w\|$  in accordance to the noise level, as prescribed by Theorems 1 and 3. As detailed in these theorems, the value of  $C_0$  is chosen large enough to obtain the desired denoising effect (otherwise the solution does not have a low complexity), but its precise value does not affect the observed identification results we describe below. In the numerical results reported hereafter, we used  $n_0 = 20$ ,  $P = 3n_0^2/4 = 300$  and  $\text{rank}(x_0) = 3$ .

Figure 2 shows how the F-B iterations behave radically differently depending on whether the non-degeneracy condition  $\mathbf{IC}(x_0) < 1$  holds or not. Each curve shows the evolution of  $\text{rank}(x^k)$  during the course of iterations, for different randomized instances of the low-rank matrix  $x_0$  to recover. As predicted by Theorem 3, one can see that for those  $x_0$  where model consistency holds (i.e.  $\mathbf{IC}(x_0) < 1$ , plotted in red), F-B converges with the correct rank, i.e.  $\text{rank}(x^k) = \text{rank}(x_0)$  for  $k$  large enough. In sharp contrast, when the model is not stable (i.e.  $\mathbf{IC}(x_0) > 1$ , displayed in blue), one observes numerically that  $\text{rank}(x^k) > \text{rank}(x_0)$ , and that the correct rank is never selected by the algorithm (although of course one still has  $x^k \rightarrow x_0$ ).

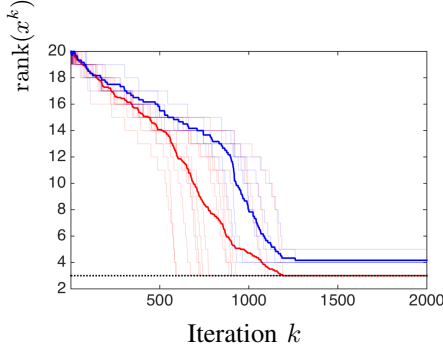


Fig. 2. Evolution of  $\text{rank}(x^k)$  as a function of  $k$  during the progress of Forward-Backward iterations (13) to solve (1) using observations  $y = \Phi x_0 + w$ . The light red (resp. blue) curves represent the evolution for an input  $x_0$  satisfying  $\text{IC}(x_0) < 1$  (resp.  $\text{IC}(x_0) > 1$ ). The bold red (resp. blue) curve is the average of the light red (resp. blue) curves.

## V. PROOFS

### A. Uniqueness sufficient condition

**Proposition 4.** *Let  $J$  be a proper lsc convex function. For a point  $x$ , assume that*

$$\ker(\Gamma) \cap T_x = \{0\}, \quad \text{and} \quad \eta_\Gamma \in \text{ri}(\partial J(x)).$$

*Then  $x$  is the unique minimizer of  $(\mathcal{P}_\theta)$  (resp.  $(\mathcal{P}_{0, \Gamma x, \Gamma})$ ).*

*Proof.* This is a consequence of [54, Corollary 1]. Though their result was stated for  $J$  finite-valued convex, it remains valid when it is proper lsc and convex. Indeed, in this case,  $J$  is subdifferentially regular at  $x$  [47, Example 7.27]. Moreover,  $\partial J(x) \neq \emptyset$  by assumption, and thus the directional derivative at  $x$  is proper, sublinear and closed, and it is the support of  $\partial J(x)$  [47, Theorem 8.30]. Continuing the proof as in [54, Corollary 1] shows the claim.  $\square$

### B. Proof of Proposition 1

*Proof of (6)* Under condition  $\ker(\Phi) \cap T = \{0\}$ , we have from the definition of  $\tilde{\Gamma}_T^+$ , that

$$z_{\tilde{\Gamma}} = \tilde{\Gamma}_T^+ e = \underset{z}{\text{argmin}} \|z\| \quad \text{subject to} \quad \tilde{\Gamma}_T z = e \quad (22)$$

and thus

$$\eta_{\tilde{\Gamma}} = \tilde{\Gamma} z_{\tilde{\Gamma}}.$$

Clearly, the constraint set of problem (22) includes that of (5), which entails

$$\|z_{\tilde{\Gamma}}\| \leq \|\tilde{z}_{\tilde{\Gamma}}\|.$$

If  $\eta_{\tilde{\Gamma}} \in \text{ri}(\partial J(x_0))$ , then  $z_{\tilde{\Gamma}}$  is also a feasible point of problem (5) and thus

$$\|\tilde{z}_{\tilde{\Gamma}}\| \leq \|z_{\tilde{\Gamma}}\|.$$

Altogether, we get that  $\|\tilde{z}_{\tilde{\Gamma}}\| = \|z_{\tilde{\Gamma}}\|$  and, since  $\tilde{z}_{\tilde{\Gamma}}$  is the unique minimizer of (5), we get that  $\tilde{z}_{\tilde{\Gamma}} = z_{\tilde{\Gamma}}$ , which implies that  $\hat{\eta}_{\tilde{\Gamma}} = \eta_{\tilde{\Gamma}}$ .

*Proof of (7)* Let  $S = T^\perp$ . Problem (5) can be conveniently rewritten as

$$\tilde{z}_{\tilde{\Gamma}} = \underset{z}{\text{argmin}} \|z\| \quad \text{subject to} \quad \begin{cases} \tilde{\Gamma}_T z = e \\ \tilde{\Gamma}_S z \in P_S(\partial J(x_0)) \end{cases}.$$

The fact that  $\hat{\eta}_{\tilde{\Gamma}} = \tilde{\Gamma} \tilde{z}_{\tilde{\Gamma}} \in \text{ri}(\partial J(x_0))$  implies  $P_S \hat{\eta}_{\tilde{\Gamma}} = P_S \tilde{\Gamma} \tilde{z}_{\tilde{\Gamma}} \in \text{ri}(P_S \partial J(x_0))$ , and thus, the second constraint in the last problem is inactive. We then recover problem (22), which in turn implies that  $\hat{\eta}_{\tilde{\Gamma}} = \eta_{\tilde{\Gamma}}$ .

*Proof of uniqueness.* See Proposition 4.

### C. Proof of Theorem 1

In order to prove Theorem 1, we consider any sequence  $\theta_k = (\mu_k, u_k = \Gamma_k x_0 + \varepsilon_k, \Gamma_k)_k$  where  $\Phi_k \in \mathbb{R}^{p_k \times n}$ . Assume that

$$(\Gamma_k, \varepsilon_k \mu_k^{-1}, \mu_k) \longrightarrow (\tilde{\Gamma}, 0, 0). \quad (23)$$

Then proving Theorem 1 boils down to showing that for  $k$  large enough, the solution  $x_k$  of  $(\mathcal{P}_{\theta_k})$  is unique and satisfies  $x_k \in \mathcal{M}$ .

*a) Constrained problem.* We consider the following non-smooth, in general non-convex, constrained minimization problem

$$x_k \in \underset{x \in \mathcal{M} \cap \mathcal{K}}{\text{Argmin}} E(x, \theta_k) \quad (24)$$

where  $\mathcal{K}$  is an arbitrary fixed convex compact neighbourhood of  $x_0$ .

The following key lemma establishes the convergence of  $x_k$  to  $x_0$ .

**Lemma 1.** *Under conditions (8) and (23),  $x_k \rightarrow x_0$ .*

*Proof.* We denote  $\|u\|_\Gamma^2 = \langle \Gamma u, u \rangle$  for any positive semidefinite matrix  $\Gamma$ . Under condition (8), Proposition 4 implies that  $x_0$  is the unique solution of  $(\mathcal{P}_{0, \tilde{\Gamma} x_0, \tilde{\Gamma}})$ . By optimality of  $x_k$  one has  $E(x_k, \theta_k) \leq E(x_0, \theta_k)$  and hence

$$\begin{aligned} & \frac{1}{2} \|x_k\|_{\Gamma_k}^2 - \langle x_k, \Gamma_k x_0 + \varepsilon_k \rangle + \mu_k J(x_k) \\ & \leq \frac{1}{2} \|x_0\|_{\Gamma_k}^2 - \langle x_0, \Gamma_k x_0 + \varepsilon_k \rangle + \mu_k J(x_0) \end{aligned}$$

which is equivalently stated as

$$\frac{1}{2} \|x_k - x_0\|_{\Gamma_k}^2 - \langle x_k - x_0, \varepsilon_k \rangle + \mu_k J(x_k) \leq \mu_k J(x_0). \quad (25)$$

Since  $x_k \in \mathcal{K}$ , the sequence  $(x_k)_k$  is bounded, and we let  $x^*$  be any cluster point. Using (23), that  $J$  is non-negative and lsc, and  $J(x_k)$  are bounded, we have

$$\begin{aligned} \limsup_{k \rightarrow \infty} (\mu_k J(x_k)) & \leq \lim_{k \rightarrow \infty} \mu_k \limsup_{k \rightarrow \infty} J(x_k) = 0 \quad \text{and} \\ \liminf_{k \rightarrow \infty} (\mu_k J(x_k)) & \geq \lim_{k \rightarrow \infty} \mu_k \liminf_{k \rightarrow \infty} J(x_k) \\ & \geq J(x^*) \lim_{k \rightarrow \infty} \mu_k = 0, \end{aligned}$$

and thus  $\lim_{k \rightarrow \infty} (\mu_k J(x_k)) = 0$ . Consequently, passing to the limit in (25), using (23), and continuity of the inner product and the norm, shows that  $\|x^* - x_0\|_\Gamma^2 \leq 0$ , or equivalently

$\tilde{\Gamma}x^* = \tilde{\Gamma}x_0$ , i.e.  $x^*$  is a feasible point of  $(\mathcal{P}_{0,\tilde{\Gamma}x_0,\tilde{\Gamma}})$ . Furthermore, since  $\frac{1}{2}\|x_k - x_0\|_{\tilde{\Gamma}_k}^2 \geq 0$ , (25) yields

$$-\langle x_k - x_0, \frac{\varepsilon_k}{\mu_k} \rangle + J(x_k) \leq J(x_0).$$

Passing again to the limit, using lower semicontinuity of  $J$ , (23) and continuity of the inner product, we then get

$$\begin{aligned} J(x^*) &\leq \liminf_{k \rightarrow \infty} J(x_k) \\ &= \liminf_{k \rightarrow \infty} \left( -\langle x_k - x_0, \frac{\varepsilon_k}{\mu_k} \rangle + J(x_k) \right) \\ &\leq \limsup_{k \rightarrow \infty} \left( -\langle x_k - x_0, \frac{\varepsilon_k}{\mu_k} \rangle + J(x_k) \right) \\ &= \limsup_{k \rightarrow \infty} J(x_k) \leq J(x_0). \end{aligned}$$

Combining this with the previous claim on feasibility of  $x^*$  for  $(\mathcal{P}_{0,\tilde{\Gamma}x_0,\tilde{\Gamma}})$  allows to conclude that  $x^*$  is a solution of  $(\mathcal{P}_{0,\tilde{\Gamma}x_0,\tilde{\Gamma}})$ . Since  $x_0$  is unique, this leads to  $x^* = x_0$ .  $\square$

We now aim at showing that for  $k$  large enough,  $x_k$  is the unique solution of  $(\mathcal{P}_{\theta_k})$ .

*b) Convergence of the tangent model subspace.:* By definition of the constrained problem (24),  $x_k \in \mathcal{M}$ . Moreover, since  $E(\cdot, \theta_k)$  is partly smooth at  $x_0$  relative to  $\mathcal{M}$ , the sharpness property Definition 1(ii) holds at all nearby points in the manifold  $\mathcal{M}$  [34, Proposition 2.10]. Thus as soon as  $k$  is large enough, we have  $T_k = \mathcal{T}_{x_k}(\mathcal{M})$ . Using the fact that  $\mathcal{M}$  is of class  $C^2$ , we get

$$T_k = \mathcal{T}_{x_k}(\mathcal{M}) \longrightarrow \mathcal{T}_{x_0}(\mathcal{M}) = T \quad (26)$$

when (23) holds, where the convergence should be understood over the Grassmannian of linear subspaces with the same dimension (or equivalently, as the convergence of the projection operators  $P_{T_k} \rightarrow P_T$ ). Since  $\ker(\tilde{\Gamma}) \cap T = \{0\}$ , (26) implies that for  $k$  large enough, when (23) holds,

$$\ker(\Gamma_k) \cap T_k = \{0\}, \quad (27)$$

which we assume from now on.

*c) First order condition.:* Let  $\mathbb{B}$  be the Euclidean unit ball in  $\mathbb{R}^p$ . Take  $\mathcal{K} = x_0 + r\mathbb{B}$  for  $r > 0$  sufficiently large. For any  $\delta > 0$ ,  $\exists k_\delta$  such that  $\forall k \geq k_\delta$ ,  $x_k \in x_0 + \delta\mathbb{B}$  according to Lemma 1. Thus, for  $k$  large enough, i.e.  $\delta$  sufficiently small, we indeed have  $x_k \in \text{int}(\mathcal{K})$ . Furthermore, it is easy to see that  $\iota_{\mathcal{K}}$  is locally partly smooth at  $x_0$  relative to  $\mathcal{K}$ , and thus is partly smooth at  $x_k$  relative to  $\mathcal{K}$  for  $k$  large enough. Moreover, local partial smoothness of  $J$  at  $x_0$  relative to  $\mathcal{M}$  entails that  $J$  is also partly smooth at  $x_k$  relative to  $\mathcal{M}$ . Therefore, the sum rule [34, Corollary 4.6] (the transversality condition is satisfied as  $\mathcal{K}$  is full-dimensional and  $x_k \in \text{int}(\mathcal{K})$ , see (4)) shows that, for all sufficiently large  $k$ ,  $J + \iota_{\mathcal{K}}$  is locally partly smooth at  $x_k$  relative to  $\mathcal{M} \cap \mathcal{K}$ , and then so is  $E(\cdot, \theta_k) + \iota_{\mathcal{K}}$  by the smooth perturbation rule [34, Corollary 4.7]. Therefore,

[34, Proposition 2.4(a)-(b)] applies, and it follows that  $x_k$  is a critical point of (24) if, and only if,

$$\begin{aligned} 0 &\in \text{Aff}(\partial E(x_k, \theta_k) + N_{\mathcal{K}}(x_k)) \\ &= \frac{\Gamma_k x_k - u_k}{\mu_k} + \text{Aff}(\partial J(x_k)) \\ &= \frac{\Gamma_k x_k - u_k}{\mu_k} + e_{x_k} + T_k^\perp, \end{aligned}$$

where  $e_{x_k} = P_{T_k}(\partial J(x_k))$ . The first equality comes from the fact that  $E(\cdot, \theta)$  is a closed convex function, and that the normal cone of  $\mathcal{K}$  at  $x_k$  vanishes on the interior points of  $\mathcal{K}$ , and the second one from the decomposability of the subdifferential. Projecting this relation onto  $T_k$ , we get, since  $e_{x_k} \in T_k$ ,

$$P_{T_k}(\Gamma_k x_k - u_k) + \mu_k e_{x_k} = 0. \quad (28)$$

*d) Convergence of the primal variables.:* Since both  $x_k$  and  $x_0$  belong to  $\mathcal{M}$ , and partial smoothness implies that  $\mathcal{M}$  is a manifold of class  $C^2$  around each of them, we deduce that each point in their respective neighbourhoods has a unique projection on  $\mathcal{M}$  [44]. In particular,  $x_k = P_{\mathcal{M}}(x_k)$  and  $x_0 = P_{\mathcal{M}}(x_0)$ . Moreover,  $P_{\mathcal{M}}$  is of class  $C^1$  near  $x_k$  [36, Lemma 4]. Thus,  $C^2$  differentiability shows that

$$x_k - x_0 = P_{\mathcal{M}}(x_k) - P_{\mathcal{M}}(x_0) = \text{D}P_{\mathcal{M}}(x_k)(x_k - x_0) + R(x_k)$$

where  $R(x_k) = O(\|x_k - x_0\|^2)$  and where  $\text{D}P_{\mathcal{M}}(x_k)$  is the derivative of  $P_{\mathcal{M}}$  at  $x_k$ . Using [36, Lemma 4], and recalling that  $T_k = \mathcal{T}_{x_k}(\mathcal{M})$  by the sharpness property, the derivative  $\text{D}P_{\mathcal{M}}(x_k)$  is given by  $\text{D}P_{\mathcal{M}}(x_k) = P_{T_k}$ . Inserting this in (28), we get

$$P_{T_k} \Gamma_k (P_{T_k}(x_k - x_0) + R(x_k)) - P_{T_k} \varepsilon_k + \mu_k e_{x_k} = 0. \quad (29)$$

Using (27),  $\Gamma_{k,T_k}$  has full rank, and thus

$$x_k - x_0 = \Gamma_{k,T_k}^+ (\varepsilon_k - \mu_k e_{x_k} - \Gamma_k R(x_k)), \quad (30)$$

where we also used that  $T_k^\perp \subset \ker(\Gamma_{k,T_k}^+)$ . One has  $\Gamma_{k,T_k}^+ \rightarrow \tilde{\Gamma}_T^+$  so that  $\Gamma_{k,T_k}^+ \Gamma_k = O(1)$  and  $\Gamma_{k,T_k}^+ = O(1)$ . Altogether, we thus obtain the bound

$$\|x_k - x_0\| = O(\|\varepsilon_k\|, \mu_k). \quad (31)$$

*e) Convergence of the dual variables.:* We define  $\eta_k = \frac{u_k - \Gamma_k x_k}{\mu_k}$ . Arguing as above, and using (30) we have

$$\begin{aligned} \mu_k \eta_k &= \varepsilon_k + \Gamma_k(x_0 - x_k) \\ &= \varepsilon_k - \Gamma_k \Gamma_{k,T_k}^+ (\varepsilon_k - \mu_k e_{x_k} - \Gamma_k R(x_k)) \\ &= \varepsilon_k - \Gamma_k P_{T_k} \Gamma_{k,T_k}^+ (\varepsilon_k - \mu_k e_{x_k} - \Gamma_k R(x_k)) \\ &= P_{V_{T_k}^\perp} \varepsilon_k + P_{V_{T_k}} \Gamma_k R(x_k) + \mu_k \Gamma_k \Gamma_{k,T_k}^+ e_{x_k}, \end{aligned}$$

where we denoted  $V_{T_k} = \text{Im}(\Gamma_k P_{T_k})$ , and used that  $\text{Im}(\Gamma_{k,T_k}^+) \subset T_k$ . We thus arrive at

$$\begin{aligned} \|\eta_k - \eta_{\tilde{\Gamma}}\| &= O\left(\|\varepsilon_k\| \mu_k^{-1}, \|\Gamma_k \Gamma_{k,T_k}^+ e_{x_k} - \eta_{\tilde{\Gamma}}\|, \right. \\ &\quad \left. \|\Gamma_k\| \|x_k - x_0\|^2 \mu_k^{-1}\right). \end{aligned}$$

Since  $\mathcal{M}$  is a  $C^2$  manifold, and by partial smoothness ( $J$  is  $C^2$  on  $\mathcal{M}$ ), we have  $x \mapsto e_x$  is  $C^1$  on  $\mathcal{M}$ , one has

$$\|e_{x_k} - e\| = O(\|x_k - x_0\|). \quad (32)$$

Using the triangle inequality, we get

$$\|\Gamma_k \Gamma_{k,T_k}^+ - \tilde{\Gamma}_T^+\| \leq \|\Gamma_{k,T_k}^+\| \|\Gamma_k - \tilde{\Gamma}\| + \|\tilde{\Gamma}\| \|\Gamma_{k,T_k}^+ - \tilde{\Gamma}_T^+\|.$$

Again, since  $\Gamma_{k,T_k}^+ \rightarrow \tilde{\Gamma}_T^+$ , we have  $\|\Gamma_{k,T_k}^+\| = O(1)$ . Moreover,  $A \mapsto A^+$  is smooth at  $A = \Gamma_T$  along the manifold of matrices of constant rank, and  $\mathcal{M}$  is a  $C^2$  manifold near  $x_0$ . Thus

$$\begin{aligned} \|\Gamma_{k,T_k}^+ - \tilde{\Gamma}_T^+\| &= O(\|\Gamma_{k,T_k} - \tilde{\Gamma}_T\|) \\ &= O(\|\Gamma_k - \tilde{\Gamma}\|, \|\mathbb{P}_{T_k} - \mathbb{P}_T\|) \\ &= O(\|\Gamma_k - \tilde{\Gamma}\|, \|x_k - x_0\|). \end{aligned}$$

This shows that

$$\|\Gamma_k \Gamma_{k,T_k}^+ - \tilde{\Gamma}_T^+\| = O(\|\Gamma_k - \tilde{\Gamma}\|, \|x_k - x_0\|). \quad (33)$$

Putting (32) and (33) together implies

$$\|\Gamma_k \Gamma_{k,T_k}^+ e_{x_k} - \eta_{\tilde{\Gamma}}\| = O(\|\Gamma_k - \tilde{\Gamma}\|, \|x_k - x_0\|).$$

Altogether, we get the bound

$$\begin{aligned} \|\eta_k - \eta_{\tilde{\Gamma}}\| &= O\left(\|\varepsilon_k\| \mu_k^{-1}, \|x_k - x_0\|, \|\Gamma_k - \tilde{\Gamma}\|, \right. \\ &\quad \left. \|\Gamma_k\| \|x_k - x_0\|^2 \mu_k^{-1}\right). \end{aligned} \quad (34)$$

Since  $\|x_k - x_0\|$  is bounded according to (31), we arrive at

$$\|\eta_k - \eta_{\tilde{\Gamma}}\| = O\left(\|\Gamma_k - \tilde{\Gamma}\|, \|\varepsilon_k\| \mu_k^{-1}, \mu_k\right). \quad (35)$$

f) *Convergence inside the relative interior.*: Using the hypothesis that  $\eta_{\tilde{\Gamma}} \in \text{ri}(\partial J(x_0))$ , we will show that for  $k$  large enough,

$$\eta_k \in \text{ri}(\partial J(x_k)). \quad (36)$$

Let us suppose this does not hold. Then there exists a sub-sequence of  $\eta_k$ , that we do not relabel for the sake of readability of the proof, such that

$$\eta_k \in \text{rbd}(\partial J(x_k)). \quad (37)$$

According to (35) and Lemma 1, under (23),  $(x_k, \eta_k) \rightarrow (x_0, \eta_{\tilde{\Gamma}})$ . Condition (37) is equivalently stated as, for each  $k$

$$\exists z_k \in T_{x_k}^\perp, \quad \forall \eta \in \partial J(x_k), \quad \langle z_k, \eta - \eta_k \rangle \geq 0, \quad (38)$$

where one can impose the normalization  $\|z_k\| = 1$  by positive-homogeneity. Up to a sub-sequence (that for simplicity we still denote  $z_k$  with a slight abuse of notation), since  $z_k$  is in a compact set, we can assume  $z_k$  approaches a non-zero cluster point  $z^*$ .

Since  $T_{x_k}^\perp \rightarrow T^\perp$  because  $\mathcal{M}$  is a  $C^2$  manifold, one has that  $z^* \in T^\perp$ . We now show that

$$\forall v \in \partial J(x_0), \quad \langle z^*, v - \eta_{\tilde{\Gamma}} \rangle \geq 0. \quad (39)$$

Indeed, let us consider any  $v \in \partial J(x_0)$ . In view of the continuity property in Definition 1(iii)  $\partial J$  is continuous at  $x_0$  along  $\mathcal{M}$ , so that since  $x_k \rightarrow x_0$  there exists  $v_k \in \partial J(x_k)$  with  $v_k \rightarrow v$ . Applying (38) with  $\eta = v_k$  gives  $\langle z_k, v_k - \eta_k \rangle \geq 0$ . Taking the limit  $k \rightarrow +\infty$  in this inequality leads to (39), which contradicts the fact that  $\eta_{\tilde{\Gamma}} \in \text{ri}(\partial J(x_0))$ . In view of (36) and (27), using Proposition 4 shows that  $x_k$  is the unique solution of  $(\mathcal{P}_{\theta_k})$ .

## D. General Loss Function

We now detail the necessary arguments to adapt the proof of Theorem 1 to a generic loss function satisfying assumptions **(A.1)**-**(A.2)**.

a) *Proof of Proposition 4*: It follows from assumption **(A.1)** that  $F(\cdot, y)$  is strictly convex, and the uniqueness follows from [38, Theorem A.1].

b) *Proof of Lemma 1*: Problem (24) now reads

$$x_k \in \underset{x \in \mathcal{M} \cap \mathcal{K}}{\text{Argmin}} F(\Phi_k x, y_k) + \lambda_k J(x).$$

Optimality of  $x_k$  entails

$$F(\Phi_k x_k, y_k) + \lambda_k J(x_k) \leq F(\Phi_k x_0, y_k) + \lambda_k J(x_0).$$

By assumptions **(A.1)**-**(A.2)**, we have the following useful inequalities for any  $u \in \mathbb{R}^p$ , see e.g. [41, p. 57 and 64]

$$\begin{aligned} \frac{\sigma_m}{2} \|y - u\|^2 &\leq F(u, y) - F(y, y) \\ &= F(u, y) - F(y, y) - \langle \nabla F(y, y), u - y \rangle \\ &\leq \frac{\sigma_M}{2} \|y - u\|^2. \end{aligned}$$

It then follows that

$$F(\Phi_k x_k, y_k) - F(\Phi_k x_0, y_k) \geq \frac{\sigma_m}{2} \|y_k - \Phi_k x_k\|^2 - \frac{\sigma_M}{2} \|w_k\|^2$$

and therefore

$$\begin{aligned} \lambda_k J(x_0) &\geq \frac{\sigma_m}{2} \|y_k - \Phi_k x_k\|^2 - \frac{\sigma_M}{2} \|w_k\|^2 + \lambda_k J(x_k) \\ &\geq \frac{\sigma_m}{2\sigma_M} \|y_k - \Phi_k x_k\|_{\nabla_1^2 F(y_k, y_k)}^2 \\ &\quad - \frac{\sigma_M}{2} \|w_k\|^2 + \lambda_k J(x_k) \\ &\geq \frac{\sigma_m}{2\sigma_M} \|x_k - x_0\|_{\Phi_k^* \nabla_1^2 F(y_k, y_k) \Phi_k}^2 \\ &\quad - \frac{\sigma_m}{\sigma_M} \langle x_k - x_0, \Phi_k^* \nabla_1^2 F(y_k, y_k) w_k \rangle \\ &\quad - \frac{\sigma_M}{2} \left(1 - \frac{\sigma_m^2}{\sigma_M^2}\right) \|w_k\|^2 + \lambda_k J(x_k), \end{aligned}$$

where we used strong convexity of assumption **(A.1)** in the second and third inequalities. Dividing both sides by  $1/P$  we obtain

$$\begin{aligned} \frac{\sigma_m}{2\sigma_M} \|x_k - x_0\|_{\tilde{\Gamma}_k}^2 - \frac{\sigma_m}{\sigma_M} \langle x_k - x_0, \check{\varepsilon}_k \rangle \\ - \frac{\sigma_M}{2} \left(1 - \frac{\sigma_m^2}{\sigma_M^2}\right) \|n^{-1/2} w_k\|^2 + \mu_k J(x_k) \leq \mu_k J(x_0), \end{aligned}$$

where now

$$\tilde{\Gamma}_k = \frac{1}{p} \Phi_k^* \nabla_1^2 F(y_k, y_k) \Phi_k,$$

and

$$\check{\varepsilon}_k = \frac{1}{p} \Phi_k^* \nabla_{1,2} F(y_k, y_k) w_k.$$

Changing (23) to  $(\tilde{\Gamma}_k, \check{\varepsilon}_k \mu_k^{-1}, \mu_k) \rightarrow (\tilde{\Gamma}, 0, 0)$ , which entails implicitly that  $n^{-1/2} w_k \rightarrow 0$ , and arguing as in the rest of the proof of the lemma allows to conclude that  $x_k \rightarrow x_0$ .

c) *Proof of Theorem 1:*  $C^2$ -continuity of  $F$  allows to use the smooth perturbation rule to conclude that partial smoothness of  $J$  is preserved upon adding  $F$ . Condition (28) now becomes

$$P_{T_k} \Phi_k^* \nabla_1 F(\Phi_k x_k, y_k) + \lambda_k e_{x_k} = 0.$$

Using again assumptions **(A.1)**-**(A.2)** and expanding  $\nabla_1 F(\Phi_k x_k, y_k)$  at  $(y_k, y_k)$  to the first order, we obtain

$$\begin{aligned} & \nabla_1 F(\Phi_k x_k, y_k) \\ &= \nabla_1^2 F(y_k, y_k) \Phi_k (x_k - x_0) - \nabla_1^2 F(y_k, y_k) w_k \\ & \quad + O(\|x_k - x_0\|^2) + O(\|w_k\|^2) \\ &= \nabla_1^2 F(y_k, y_k) \Phi_k (P_{T_k}(x_k - x_0) + R(x_k)) - \nabla_1^2 F(y_k, y_k) w_k \\ & \quad + O(\|x_k - x_0\|^2) + O(\|w_k\|^2). \end{aligned}$$

Dividing by  $p$ , plugging this expansion back into the above first-order (criticality) condition, and grouping the  $O(\cdot)$  terms, condition (29) becomes

$$\begin{aligned} & \check{\Gamma}_{k, T_k}(x_k - x_0) - P_{T_k} \check{\varepsilon}_k + \mu_k e_{x_k} + P_{T_k}(n^{-1} \Phi_k^* + \check{\Gamma}_k) R(x_k) \\ & \quad + P_{T_k} \Phi_k^* Q(n^{-1/2} w_k) = 0, \end{aligned}$$

where  $Q(n^{-1/2} w_k) = O(\|n^{-1/2} w_k\|^2)$ . Then with the new notations  $(\check{\Gamma}_k, \check{\varepsilon}_k)$  in place of  $(\Gamma_k, \varepsilon_k)$ , one sees that the proof continues unchanged.

### E. Proof of Theorem 2

It is sufficient to check that (9) is in force with probability 1 as  $p \rightarrow +\infty$ . Owing to classical results on convergence of sample covariances, which apply thanks to the assumption that the fourth order moments are finite, we get  $\Gamma_p - \tilde{\Gamma} = O_P(p^{-1/2})$  and  $\frac{1}{p} \langle \xi_i, w \rangle = O_P(p^{-1/2})$ , where we used the assumption that  $\mathbb{E}(\langle \xi_i, w \rangle) = 0$ . As  $p$  is fixed, it follows that  $\|\Gamma_p - \tilde{\Gamma}\| = O_P(p^{-1/2})$  and  $\|\varepsilon_p\| = O_P(p^{-1/2})$ . Thus under the scaling (12), we get

$$\begin{aligned} & \left( \|\Gamma_p - \tilde{\Gamma}\|, \|\varepsilon_p\| \mu_p^{-1}, \mu_p \right) \\ &= \left( O_P(p^{-1/2}), \frac{1}{\mu_p p^{1/2}} O_P(1), o(1) \right) \\ &= \left( O_P(p^{-1/2}), o(1) O_P(1), o(1) \right) \\ &= \left( O_P(p^{-1/2}), o(1), o(1) \right), \end{aligned}$$

which indeed converges to 0 in probability. This concludes the proof.

### F. Proof of Proposition 2

Let  $(x_k)_k$  be a sequence of solutions to the constrained problem (24). Since  $x_0$  is the unique minimizer to  $(\mathcal{P}_{(0, \tilde{\Gamma} x_0, \tilde{\Gamma})})$  and (8) is satisfied,  $\eta_{\tilde{\Gamma}}$  is well-defined. Moreover, arguing as in the proof of Lemma 1 and Theorem 1, under condition (9), we have  $(x_k, \eta_k) \rightarrow (x_0, \eta_{\tilde{\Gamma}})$ , and  $\eta_k \in \eta_{\tilde{\Gamma}} + C\mathbb{B}$ .

Let  $\tau = \text{dist}(\eta_{\tilde{\Gamma}}, \partial J(x_0)) = \inf_{\eta \in \partial J(x_0)} \|\eta - \eta_{\tilde{\Gamma}}\|$ . Since  $\partial J(x_0)$  is a non-empty, closed and convex set, the infimum is attained and one has  $\tau > 0$  since  $\eta_{\tilde{\Gamma}} \notin \partial J(x_0)$ .

We now prove the claim by contradiction. Let  $x_j$  be a solution of  $(\mathcal{P}_{\theta_j})$  such that (9) holds at  $\theta_j$  for  $j$  sufficiently

large (taking  $C$  smaller if necessary so that  $C < \tau$ ), and suppose that  $x_j \in \mathcal{M}$ . Thus,  $x_j$  is also a solution of (24) for  $\theta_j$ , whence it follows that  $\eta_j \in \eta_{\tilde{\Gamma}} + C\mathbb{B}$ . Using the triangle inequality, we then get

$$\text{dist}(\eta_j, \partial J(x_0)) > \tau - C > 0. \quad (40)$$

Now, in view of the continuity property in Definition 1((iii)), we have  $\partial J(x_k) \rightarrow \partial J(x_0)$  along  $\mathcal{M}$ . This is equivalent, since  $\partial J(x_0)$  is closed and using [47, Corollary 4.7], to  $\text{dist}(\eta, \partial J(x_k)) \rightarrow \text{dist}(\eta, \partial J(x_0))$  for every  $\eta \in \mathbb{R}^p$ , i.e.

$$\begin{aligned} & \forall \delta > 0, \exists k_0, \forall k \geq k_0, \forall \eta \in \mathbb{R}^p, \\ & \quad |\text{dist}(\eta, \partial J(x_k)) - \text{dist}(\eta, \partial J(x_0))| < \delta. \end{aligned}$$

In particular, as  $x_j$  is a minimizer of  $(\mathcal{P}_{\theta_j})$  for  $j$  large enough, we have  $\eta_j \in \partial J(x_j)$ , and thus  $\text{dist}(\eta_j, \partial J(x_0)) < \delta$ , leading to a contradiction with (40). Hence,  $x_j \notin \mathcal{M}$ .

## VI. CONCLUSION

In this paper, we provided a very general and principled analysis of the recovery performance when partly smooth functions are used to regularize linear inverse/regression problems. This class of functions encompass all popular regularizers used in the literature. The generality of our results is unprecedented since for the first time, a unified analysis is provided together with a generalized “irrepresentable condition” to guarantee consistent identification of the low-complexity manifold underlying the original object. Our work also shows that model consistency is not only of theoretical interest, but also has algorithmic and practical consequences. Indeed, after a finite number of iterations, the iterates of the proximal splitting algorithm used to solve the original optimization problem (here the Forward-Backward), are guaranteed to lie on the original manifold. This opens the door to acceleration by switching to a higher-order smooth optimization method, exploiting the smoothness of the partly smooth objective function along the identified smooth model manifold.

## ACKNOWLEDGEMENTS

The authors would like to thank Vincent Duval and Jérôme Malick for fruitful discussions. This work has been supported by the European Research Council (ERC project SIGMA-Vision).

## REFERENCES

- [1] J.-F. Aujol et al. “Image decomposition into a bounded variation component and an oscillating component”. In: *Journal of Mathematical Imaging and Vision* 22 (2005), pp. 71–88.
- [2] F.R. Bach. “Consistency of the Group Lasso and Multiple Kernel Learning”. In: *Journal of Machine Learning Research* 9 (2008), pp. 1179–1225.
- [3] F.R. Bach. “Consistency of Trace Norm Minimization”. In: *Journal of Machine Learning Research* 9 (2008), pp. 1019–1048.
- [4] H. H. Bauschke and P.L. Combettes. *Convex Analysis and Monotone Operator Theory in Hilbert Spaces*. Springer, 2011.

- [5] J. F. Bonnans and A. Shapiro. *Perturbation analysis of optimization problems*. Springer Series in Operations Research and Financial Engineering. Springer Verlag, 2000.
- [6] M. Burger and S. Osher. “Convergence rates of convex variational regularization”. In: *Inverse Problems* 20.5 (2004), p. 1411.
- [7] E. Candès and B. Recht. “Exact Matrix Completion via Convex Optimization”. English. In: *Foundations of Computational Mathematics* 9.6 (2009), pp. 717–772.
- [8] E. J. Candès et al. “Robust Principal Component Analysis?” In: *J. ACM* 58.3 (June 2011), 11:1–11:37.
- [9] E.J. Candès and Y. Plan. “Matrix Completion With Noise”. In: *Proceedings of the IEEE* 98.6 (2010), pp. 925–936.
- [10] E.J. Candès and B. Recht. “Simple bounds for recovering low-complexity models”. In: *Math. Program* 141.1-2 (2013), pp. 577–589.
- [11] E.J. Candès, T. Strohmer, and V. Voroninski. “PhaseLift: Exact and Stable Signal Recovery from Magnitude Measurements via Convex Programming”. In: *Communications on Pure and Applied Mathematics* 66.8 (2013), pp. 1241–1274. ISSN: 1097-0312.
- [12] E.J. Candès and T. Tao. “The power of convex relaxation: Near-optimal matrix completion”. In: *IEEE Transactions on Information Theory* 56.5 (2009), pp. 2053–2080.
- [13] S.S. Chen, D.L. Donoho, and M.A. Saunders. “Atomic decomposition by basis pursuit”. In: *SIAM journal on scientific computing* 20.1 (1999), pp. 33–61.
- [14] A. Daniilidis, D. Drusvyatskiy, and A. S. Lewis. “Orthogonal Invariance and Identifiability”. In: *to appear in SIAM J. Matrix Anal. Appl.* (2014).
- [15] A. Daniilidis, J. Malick, and H. Sendov. “Spectral (Isotropic) Manifolds and Their Dimension”. In: *to appear in Journal d’Analyse Mathématique* (2014).
- [16] C. Dossal et al. “Sharp Support Recovery from Noisy Random Measurements by L1 minimization”. In: *Applied and Computational Harmonic Analysis* 33.1 (2012), pp. 24–43. DOI: 10.1016/j.acha.2011.09.003.
- [17] V. Duval and G. Peyré. *Exact Support Recovery for Sparse Spikes Deconvolution*. Tech. rep. Preprint hal-00839635, 2013.
- [18] M. Elad, P. Milanfar, and R. Rubinstein. “Analysis versus synthesis in signal priors”. In: *Inverse problems* 23.3 (2007), p. 947. DOI: 10.1088/0266-5611/23/3/007.
- [19] J. Fadili et al. “Stable Recovery with Analysis Decomposable Priors”. In: *Proc. Sampta’13*. 2013, pp. 113–116.
- [20] M. Fazel. “Matrix rank minimization with applications”. PhD thesis. Stanford University, 2002.
- [21] J.J. Fuchs. “On sparse representations in arbitrary redundant bases”. In: *IEEE Transactions on Information Theory* 50.6 (2004), pp. 1341–1344.
- [22] M. Golbabaee and P. Vandergheynst. “Hyperspectral Image Compressed Sensing Via Low-Rank And Joint-Sparse Matrix Recovery”. In: *2012 IEEE International Conference On Acoustics, Speech And Signal Processing (ICASSP)*. IEEE, 2012, pp. 2741–2744.
- [23] M. Grasmair. “Linear convergence rates for Tikhonov regularization with positively homogeneous functionals”. In: *Inverse Problems* 27 (2011), p. 075014.
- [24] M. Grasmair, O. Scherzer, and M. Haltmeier. “Necessary and sufficient conditions for linear convergence of  $\ell_1$ -regularization”. In: *Communications on Pure and Applied Mathematics* 64.2 (2011), pp. 161–182.
- [25] E. Grave, G. Obozinski, and F. Bach. “Trace Lasso: a trace norm regularization for correlated designs”. In: *Proc. NIPS*. Ed. by John Shawe-Taylor et al. 2011, pp. 2187–2195.
- [26] W. L. Hare. “Nonsmooth optimization with smooth substructure”. PhD thesis. Simon Fraser University, 2005.
- [27] W.L. Hare and A.S. Lewis. “Identifying active constraints via partial smoothness and prox-regularity”. In: *J. Convex Anal.* 11.2 (2004), pp. 251–266.
- [28] M.A. Herman and T. Strohmer. “General Deviants: An Analysis of Perturbations in Compressed Sensing”. In: *Selected Topics in Signal Processing, IEEE Journal of* 4.2 (2010), pp. 342–349.
- [29] Y. Sun J. D. Lee and Y. E. Taylor. “On model selection consistency of regularized  $M$ -estimators”. In: *Electronic Journal of Statistics* 9 (2015), pp. 608–642.
- [30] J. Jia and B. Yu. “On model selection consistency of the elastic net when  $p \gg n$ ”. In: *Statistica Sinica* 20 (2010), pp. 595–611.
- [31] K. Knight and W. Fu. “Asymptotics for Lasso-Type Estimators”. In: *The Annals of Statistics* 28.5 (2000), pp. 1356–1378.
- [32] J. M. Lee. *Smooth manifolds*. Springer, 2003.
- [33] C. Lemaréchal, F. Oustry, and C. Sagastizábal. “The  $U$ -Lagrangian of a convex function”. In: *Trans. Amer. Math. Soc.* 352.2 (2000), pp. 711–729.
- [34] A. S. Lewis. “Active sets, nonsmoothness, and sensitivity”. In: *SIAM Journal on Optimization* 13.3 (2003), pp. 702–725.
- [35] A. S. Lewis. “The mathematics of eigenvalue optimization”. In: *Mathematical Programming* 97.1–2 (2003), pp. 155–176.
- [36] A. S. Lewis and J. Malick. “Alternating Projections on Manifolds”. In: *Mathematics of Operations Research* 33.1 (2008), pp. 216–234.
- [37] A. S. Lewis and S. Zhang. “Partial Smoothness, Tilt Stability, and Generalized Hessians”. In: *SIAM Journal on Optimization* 23.1 (2013), pp. 74–94.
- [38] J. Liang, M.J Fadili, and G. Peyré. *Local Linear Convergence of Forward-Backward under Partial Smoothness*. Tech. rep. appeared in NIPS 2014. arxiv preprint arXiv:1407.5611, 2014.
- [39] Po-Ling Loh and Martin J. Wainwright. “High-dimensional regression with noisy and missing data: Provable guarantees with nonconvexity”. In: *The Annals of Statistics* 40.3 (June 2012), pp. 1637–1664. DOI: 10.1214/12-AOS1018.
- [40] B.S. Mordukhovich. “Sensitivity analysis in nonsmooth optimization”. In: *Theoretical Aspects of Industrial De-*

- sign (D. A. Field and V. Komkov, eds.), *SIAM Volumes in Applied Mathematics* 58 (1992), pp. 32–46.
- [41] Y. Nesterov. *Introductory Lectures on Convex Optimization: A Basic Course*. Science & Business Media. Springer, 2004.
- [42] S. Oymak et al. “Simultaneously Structured Models with Application to Sparse and Low-rank Matrices”. In: *arXiv preprint arXiv:1212.3753* (2012).
- [43] G. Peyré, M.J. Fadili, and J.-L. Starck. “Learning the Morphological Diversity”. In: *SIAM Journal on Imaging Sciences* 3.3 (2010), pp. 646–669.
- [44] R.A. Poliquin, R.T. Rockafellar, and L. Thibault. “Local differentiability of distance functions”. In: *Trans. Amer. Math. Soc.* 352 (2000), pp. 5231–5249.
- [45] B. Recht, M. Fazel, and P.A. Parrilo. “Guaranteed minimum-rank solutions of linear matrix equations via nuclear norm minimization”. In: *SIAM review* 52.3 (2010), pp. 471–501.
- [46] E. Richard, F.R. Bach, and J-P. Vert. “Intersecting singularities for multi-structured estimation”. In: *Proc. ICML*. Vol. 28. JMLR Proceedings. JMLR.org, 2013, pp. 1157–1165.
- [47] R.T. Rockafellar and R. Wets. *Variational analysis*. Vol. 317. Springer Verlag, 1998.
- [48] Mathieu Rosenbaum and Alexandre B. Tsybakov. “Sparse recovery under matrix uncertainty”. In: *The Annals of Statistics* 38.5 (Oct. 2010), pp. 2620–2651. DOI: 10.1214/10-AOS793.
- [49] L.I. Rudin, S. Osher, and E. Fatemi. “Nonlinear total variation based noise removal algorithms”. In: *Physica D: Nonlinear Phenomena* 60.1 (1992), pp. 259–268.
- [50] O. Scherzer et al. *Variational Methods in Imaging*. 1st. Applied Mathematical Sciences. Springer, 2009. ISBN: 0387309314.
- [51] J.-L. Starck, M. Elad, and D.L. Donoho. “Image Decomposition Via The Combination of Sparse Representations and Variational Approach”. In: *IEEE Trans. Image Processing* 14.10 (2005), pp. 1570–1582.
- [52] R. Tibshirani. “Regression shrinkage and selection via the Lasso”. In: *Journal of the Royal Statistical Society. Series B. Methodological* 58.1 (1996), pp. 267–288.
- [53] R. Tibshirani et al. “Sparsity and smoothness via the fused Lasso”. In: *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 67.1 (2005), pp. 91–108.
- [54] S. Vaïter et al. *Model Selection with Low Complexity Priors*. Tech. rep. arXiv preprint arXiv:1307.2342, 2013.
- [55] S. Vaïter et al. “Robust Sparse Analysis Regularization”. In: *IEEE Transactions on Information Theory* 59.4 (2013), pp. 2001–2016.
- [56] M. J. Wainwright. “Sharp Thresholds for High-Dimensional and Noisy Sparsity Recovery Using  $\ell_1$ -Constrained Quadratic Programming (Lasso)”. In: *IEEE Transactions on Information Theory* 55.5 (2009), pp. 2183–2202.
- [57] S. J. Wright. “Identifiable Surfaces in Constrained Optimization”. In: *SIAM Journal on Control and Optimization* 31.4 (1993), pp. 1063–1079.
- [58] M. Yuan and Y. Lin. “Model selection and estimation in regression with grouped variables”. In: *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 68.1 (2005), pp. 49–67.
- [59] P. Zhao and B. Yu. “On Model Selection Consistency of Lasso”. In: *J. Mach. Learn. Res.* 7 (Dec. 2006), pp. 2541–2563. ISSN: 1532-4435.

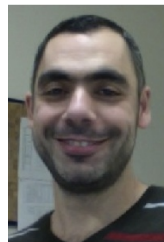


**Samuel Vaïter** is junior researcher (CR2) at the Centre Nationale de Recherche Scientifique (CNRS), working in IMB, Université de Bourgogne, Dijon. His current research interests focus on variational regularization in signal and image processing, convex analysis, sparsity and risk estimation. He has studied Applied Mathematics and Theoretical Computer Science in Lyon and Paris. He has obtained in 2014 a PhD in Applied Mathematics from Université Paris-Dauphine. He has then worked as a post-doc in CMAP, Ecole Polytechnique, Palaiseau.



**Gabriel Peyré** is senior researcher at the Centre Nationale de Recherche Scientifique (CNRS), working in DMA, Ecole Normale Supérieure, Paris. His research is focused on developing mathematical and numerical tools in sparse regularization and optimal transport, with applications in computer vision, graphics and neurosciences. Since 2005 Gabriel Peyré has co-authored 65 papers in international journals, 70 conference proceedings in top vision and image processing conferences, and two books.

He is the creator of the “Numerical tour of signal processing” ([www.numerical-tours.com](http://www.numerical-tours.com)), a popular online repository of Matlab/Python/Julia resources to teach modern signal and image processing. His research was supported by a ERC starting grant (SIGMA-Vision, 2010-2015) and is now supported by a ERC consolidator grant (NORIA 2017-2021).



**Jalal Fadili** is a Full Professor at Ecole National Supérieure d’Ingénieurs de Caen, and Junior member of Institut Universitaire de France since Oct. 2013. He holds several scientific management positions (editorial activities, national excellence research networks). He also held visiting positions at several universities (QUT-Australia, Stanford, Cal-Tech, EPFL-Switzerland, MIT). In the last decade, he has been an invited or plenary speaker at various international events. His research interests include mathematical signal and image processing, mathematical statistics, inverse problems, variational methods and regularization theory, and non-smooth optimization. His areas of application include medical and astronomical imaging. He has published more than 170 papers in the leading journals and conferences of these fields, 7 book chapters and 2 books.