

## Consistent Learning by Composite Proximal Thresholding

Patrick L. Combettes · Saverio Salzo · Silvia  
Villa

Received: date / Accepted: date

**Abstract** We investigate the modeling and the numerical solution of machine learning problems with prediction functions which are linear combinations of elements of a possibly infinite dictionary of functions. We propose a novel flexible composite regularization model, which makes it possible to incorporate various priors on the coefficients of the prediction function, including sparsity and hard constraints. We show that the estimators obtained by minimizing the regularized empirical risk are consistent in a statistical sense, and we design an error-tolerant composite proximal thresholding algorithm for computing such estimators. New results on the asymptotic behavior of the proximal forward-backward splitting method are derived and exploited to establish the convergence properties of the proposed algorithm. In particular, our method features a  $o(1/m)$  convergence rate in objective values.

---

The work of P. L. Combettes was partially supported by the CNRS MASTODONS project under grant 2016TABASCO and by the CNRS Imag'in project under grant 2015OPTIMISME. S. Villa is a member of the Gruppo Nazionale per l'Analisi Matematica, la Probabilità e le loro Applicazioni (GNAMPA) of the Istituto Nazionale di Alta Matematica (INdAM).

---

P. L. Combettes  
North Carolina State University  
Department of Mathematics  
Raleigh, NC 27695-8205, USA  
E-mail: plc@math.ncsu.edu

S. Salzo  
Massachusetts Institute of Technology and Istituto Italiano di Tecnologia  
Laboratory for Computational and Statistical Learning  
Cambridge, MA 02139, USA  
E-mail: saverio.salzo@iit.it

S. Villa  
Politecnico di Milano  
Dipartimento di Matematica  
20133 Milano, Italy  
E-mail: silvia.villa@polimi.it

**Keywords** Consistent estimator · convex optimization · forward-backward splitting · proximal algorithm · sparse data representation

**Mathematics Subject Classification (2000)** MSC 68T05 · 65K10 · 90C25 · 62G08

## 1 Introduction

A central task in data science is to extract information from collected observations. Optimization procedures play a central role in the modeling and the numerical solution of data-driven information extraction problems. In the present paper, we consider the problem of learning from examples within the framework of linear models [5, 21, 23]. The goal is to estimate a functional relation  $f$  from an input set  $\mathcal{X}$  into an output set  $\mathcal{Y} \subset \mathbb{R}$ . The data set consists of the observation of a finite number of realizations  $z_n = (x_i, y_i)_{1 \leq i \leq n}$  in  $\mathcal{X} \times \mathcal{Y}$  of independent input/output random pairs with an unknown common distribution  $P$ . We adopt a linear model, i.e., we assume that the target function  $f$  can be approximated by estimators of the form

$$f_u: \mathcal{X} \rightarrow \mathbb{R}: x \mapsto \sum_{k \in \mathbb{K}} \mu_k \phi_k(x), \quad (1.1)$$

where  $\mathbb{K}$  is at most countable,  $u = (\mu_k)_{k \in \mathbb{K}} \in \ell^2(\mathbb{K})$ , and  $(\phi_k)_{k \in \mathbb{K}}$  is a family of bounded measurable functions from  $\mathcal{X}$  to  $\mathbb{R}$ ; such a family is called a *dictionary*, and its elements are called *features*. Ideally, one could measure the performance of an estimator  $f_u$  by the quadratic risk

$$R(f_u) = \int_{\mathcal{X} \times \mathcal{Y}} |f_u(x) - y|^2 dP(x, y), \quad (1.2)$$

but this risk is not accessible as  $P$  is unknown. Thus, based on the available data,  $R$  is replaced by the *empirical risk*  $R_n(f_u) = (1/n) \sum_{i=1}^n |f_u(x_i) - y_i|^2$ . However, the direct minimization of  $R_n(f_u)$  with respect to  $u \in \ell^2(\mathbb{K})$  leads in general to estimators that may not be consistent, that is they do not approach the minimizer of the true risk (1.2) as the sample size  $n$  becomes arbitrarily large. Therefore, regularization is needed to restore consistency [19, 21]. In our approach, the estimator  $f_{\hat{u}_{n,\lambda}}$  is computed via the approximate minimization of the convex regularized empirical risk

$$\hat{u}_{n,\lambda} \in \operatorname{Argmin}_{u \in \ell^2(\mathbb{K})}^{\varepsilon_n} \left( \frac{1}{n} \sum_{i=1}^n |f_u(x_i) - y_i|^2 + \lambda \sum_{k \in \mathbb{K}} g_k(\mu_k) \right), \quad (1.3)$$

where  $\varepsilon_n \in \mathbb{R}_{++}$  accounts for the precision with which the minimization is performed,  $\lambda \in \mathbb{R}_{++}$  is the regularization parameter, and the convex functions  $(g_k)_{k \in \mathbb{K}}$  enforce or promote prior knowledge on the coefficients  $(\mu_k)_{k \in \mathbb{K}}$  of the decomposition of the target function  $f$  with respect to the dictionary. Our objective is to select a family of regularizers  $(g_k)_{k \in \mathbb{K}}$  that model a broad range of prior knowledge and, at the same time, lead to implementable solution algorithms that produce consistent estimators. To satisfy this dual objective, we shall focus our attention on the following flexible composite model: each function  $g_k: \mathbb{R} \rightarrow ]-\infty, +\infty]$  is of the form

$$g_k = \iota_{C_k} + \sigma_{D_k} + h_k, \quad h_k - \eta |\cdot|^r \in \Gamma_0^+(\mathbb{R}), \quad r \in ]1, 2], \quad \eta \in \mathbb{R}_{++}, \quad (1.4)$$

where  $\iota_{C_k}$  is the indicator function of a closed interval  $C_k \subset \mathbb{R}$ ,  $\sigma_{D_k}$  is the support function of a closed bounded interval  $D_k \subset \mathbb{R}$ ,  $\eta \in \mathbb{R}_{++}$ , and  $h_k: \mathbb{R} \rightarrow \mathbb{R}_+$  is convex and such that  $h_k(0) = 0$ . In (1.4), the role of  $C_k$  is to explicitly enforce hard constraints on the coefficients and the role of  $D_k = [\underline{\omega}_k, \overline{\omega}_k]$  is to select the thresholding interval in which the coefficients are set to zero. Note that we have no restriction on the end points  $\underline{\omega}_k$  and  $\overline{\omega}_k$ , and thus sparsity can for instance be activated only for positive coefficients by setting  $\underline{\omega}_k = 0$  and  $\overline{\omega}_k > 0$  (see Figure 2 and Remark 6(v)). Finally,  $h_k$  provides stability and will be seen to be instrumental in guaranteeing consistency. This function plays a role similar to that of the square function in elastic net formulations [38, 16]. In particular it can assume the form of an  $\ell^r$  ( $1 < r < 2$ ) term in the regularizer, which provides stability [16, Remark 1], has proved to be effective in sparsity-based regularization [25], and reduces the shrinkage of the nonzero coefficients with respect to  $\ell^2$  [38] (see Figure 1). Note that model (1.3)–(1.4) refines that considered in [9], where the  $C_k$ 's are not explicitly considered, the  $D_k$ 's are assumed to satisfy the condition  $\bigcap_{k \in \mathbb{K}} D_k \supset [-\omega, \omega]$ , with  $\omega > 0$ , and the  $h_k$ 's are assumed differentiable. This flexible model unifies several statistical estimation techniques, such as ridge regression [23, 24], elastic net [16, 38], bridge regression [22], and generalized Gaussian models [1]. Applications that may benefit from the special composite structure of (1.3)–(1.4) are those based on feature selection, for instance in genomic data analysis, see [17, 30, 38].

The main objective of our paper is to investigate statistical and algorithmic aspects of the estimators based on (1.3)–(1.4). Our main contributions are the following:

- For suitable sequences of vanishing regularization parameters  $(\lambda_n)_{n \in \mathbb{N}}$ , we prove the consistency of the estimators  $(f_{\hat{u}_{n, \lambda_n}})_{n \in \mathbb{N}}$  as  $n \rightarrow +\infty$ , as well as the convergence of the coefficients  $(\hat{u}_{n, \lambda_n})_{n \in \mathbb{N}}$  in  $\ell^r(\mathbb{K})$ , meaning that they converge to the corresponding minimizers of the true integral risk over the constraint set. This generalizes in particular the results pertaining to the elastic net framework [16], possibly obtaining, by a suitable choice of the  $h_k$ 's, a sparser pattern of features and a reduced shrinkage effect on the nonzero coefficients. Moreover, our statistical model has the following additional new features: (a) it allows for hard constraints on the coefficients; and (b) the thresholding operation can be performed over any bounded interval.
- We establish new minimizing properties for an error-tolerant forward-backward splitting algorithm in Hilbert spaces. In particular, we establish a  $o(1/m)$  rate of convergence for the objective function values in the presence of variable proximal parameters, relaxations, and computational errors. These results improve on the state of the art, which, when dealing with convergence in objective function values, considers the non-relaxed version and covers either the error free-case [4, 15], or convergence only in an ergodic sense [32].
- We provide a procedure for the inexact computation of the proximity operators of regularizers of the form given in (1.3)–(1.4), which generates approximations amenable by the proposed error-tolerant forward-backward algorithm. This leads to a fully implementable proximal thresholding gradient algorithm for the computation of the estimators  $f_{\hat{u}_{n, \lambda}}$  that features a worst-case  $o(1/m)$  rate of convergence of the objective values.

The paper is organized as follows. In Section 2, we set the problem formally and present the main results concerning the statistical and algorithmic issues pertaining to the proposed estimators. Section 3 is devoted to proving the consistency of the estimators, which is stated in Theorem 2. In Section 4, we prove Theorem 3, which concerns the asymptotic behavior of an error-tolerant proximal forward-backward splitting algorithm, and Theorem 5, which specifically deals with the structure considered in (1.3)–(1.4). Additional properties of the regularizers defined in (1.3) are studied in Appendices A and B.

**Notation.**  $\mathbb{N}^* = \mathbb{N} \setminus \{0\}$ ,  $\mathbb{R}_+ = [0, +\infty[$ , and  $\mathbb{R}_{++} = ]0, +\infty[$ . Throughout,  $\mathbb{K}$  is an at most countably infinite index set. We denote by  $(e_k)_{k \in \mathbb{K}}$  the canonical orthonormal basis of  $\ell^2(\mathbb{K})$ . The canonical norm of  $\ell^r(\mathbb{K})$  is denoted by  $\|\cdot\|_r$ . Moreover, if  $u$  and  $v$  denote sequences in  $\ell^r(\mathbb{K})$ , their  $k$ th components are respectively denoted by the Greek letters  $\mu_k$  and  $\nu_k$ . Let  $\mathcal{H}$  be a real Hilbert space. We denote by  $\langle \cdot | \cdot \rangle$  and  $\|\cdot\|$  the scalar product and the associated norm of  $\mathcal{H}$ . The set of proper lower semicontinuous convex functions from  $\mathcal{H}$  to  $] -\infty, +\infty]$  is denoted by  $\Gamma_0(\mathcal{H})$ , and the subset of  $\Gamma_0(\mathcal{H})$  of functions valued in  $[0, +\infty]$  by  $\Gamma_0^+(\mathcal{H})$ . Let  $\varphi \in \Gamma_0(\mathcal{H})$ . The subdifferential of  $\varphi$  at  $u \in \mathcal{H}$  is  $\partial\varphi(u) = \{u^* \in \mathcal{H} \mid (\forall v \in \mathcal{H}) \varphi(u) + \langle v - u | u^* \rangle \leq \varphi(v)\}$  and, for every  $\varepsilon \in \mathbb{R}_{++}$ ,  $\text{Argmin}_{\mathcal{H}}^{\varepsilon} \varphi = \{u \in \mathcal{H} \mid \varphi(u) \leq \inf \varphi(\mathcal{H}) + \varepsilon\}$ . Let  $\delta \in \mathbb{R}_{++}$ . The  $\delta$ -subdifferential of  $\varphi$  at  $u \in \mathcal{H}$  is

$$\partial_{\delta}\varphi(u) = \{u^* \in \mathcal{H} \mid (\forall v \in \mathcal{H}) \varphi(u) + \langle v - u | u^* \rangle \leq \varphi(v) + \delta\}. \quad (1.5)$$

Let  $\mathcal{D}$  be a nonempty subset of  $\mathcal{H}$ . The indicator function of  $\mathcal{D}$  is denoted by  $\iota_{\mathcal{D}}$  and the support function of  $\mathcal{D}$  is  $\sigma_{\mathcal{D}}: \mathcal{H} \rightarrow ] -\infty, +\infty]$ :  $u \mapsto \sup_{v \in \mathcal{D}} \langle v | u \rangle$ . Let  $u \in \mathcal{H}$ . Then  $\text{prox}_{\varphi} u = \text{argmin}_{v \in \mathcal{H}} (\varphi(v) + (1/2)\|u - v\|^2)$  [27]. Suppose that  $\mathcal{D}$  is a nonempty, closed, and convex subset of  $\mathcal{H}$ . Then  $\text{prox}_{\iota_{\mathcal{D}}} = \text{proj}_{\mathcal{D}}$  is the projection operator onto  $\mathcal{D}$ , and  $\text{prox}_{\sigma_{\mathcal{D}}} = \text{Id} - \text{proj}_{\mathcal{D}} = \text{soft}_{\mathcal{D}}$  is the soft-thresholder with respect to  $\mathcal{D}$ . For background on convex analysis and optimization, see [3].

## 2 Problem setting and main results

The following assumption will be made in our main results.

**Assumption 1**  $(\mathcal{X}, \mathfrak{A}_{\mathcal{X}})$  is a measurable space,  $\mathcal{Y} \subset \mathbb{R}$  is a nonempty bounded interval, and  $b = \sup_{y \in \mathcal{Y}} |y|$ . Moreover,  $P$  is a probability measure on  $\mathcal{X} \times \mathcal{Y}$  with marginal  $P_{\mathcal{X}}$  on  $\mathcal{X}$ . The risk is

$$R: L^2(P_{\mathcal{X}}) \rightarrow \mathbb{R}_+: f \mapsto \int_{\mathcal{X} \times \mathcal{Y}} |f(x) - y|^2 dP(x, y) \quad (2.1)$$

and  $(\phi_k)_{k \in \mathbb{K}}$  is a family of measurable functions from  $\mathcal{X}$  to  $\mathbb{R}$  such that, for some  $\kappa \in \mathbb{R}_{++}$ ,

$$(\forall x \in \mathcal{X}) \quad \sum_{k \in \mathbb{K}} |\phi_k(x)|^2 \leq \kappa^2. \quad (2.2)$$

The feature map is

$$\Phi: \mathcal{X} \rightarrow \ell^2(\mathbb{K}): x \mapsto (\phi_k(x))_{k \in \mathbb{K}} \quad (2.3)$$

and

$$A: \ell^2(\mathbb{K}) \rightarrow L^2(P_{\mathcal{X}}): u = (\mu_k)_{k \in \mathbb{K}} \mapsto f_u = \sum_{k \in \mathbb{K}} \mu_k \phi_k \text{ (pointwise)}. \quad (2.4)$$

In addition,

- (a)  $(C_k)_{k \in \mathbb{K}}$  is a family of closed intervals in  $\mathbb{R}$  such that  $0 \in \bigcap_{k \in \mathbb{K}} C_k$ .
- (b)  $(h_k)_{k \in \mathbb{K}}$  is a family in  $\Gamma_0^+(\mathbb{R})$  such that  $(\forall k \in \mathbb{K}) h_k(0) = 0$  and  $h_k - \eta|\cdot|^r \in \Gamma_0^+(\mathbb{R})$  for some  $r \in ]1, 2]$  and  $\eta \in \mathbb{R}_{++}$ .
- (c)  $(D_k)_{k \in \mathbb{K}}$  is a family of nonempty closed bounded intervals in  $\mathbb{R}$  such that  $\sum_{k \in \mathbb{K}} |(\inf D_k)_+|^{r^*} < +\infty$  and  $\sum_{k \in \mathbb{K}} |(\inf D_k)_-|^{r^*} < +\infty$ .

We define

$$\begin{cases} (\forall k \in \mathbb{K}) & g_k = \iota_{C_k} + \sigma_{D_k} + h_k \\ F = R \circ A: \ell^2(\mathbb{K}) \rightarrow \mathbb{R} \\ G: \ell^2(\mathbb{K}) \rightarrow ]-\infty, +\infty]: u \mapsto \sum_{k \in \mathbb{K}} g_k(\mu_k) \\ \mathcal{C} = \overline{A(\ell^2(\mathbb{K}) \cap \times_{k \in \mathbb{K}} C_k)} \quad (\text{closure is taken in } L^2(P_{\mathcal{X}})). \end{cases} \quad (2.5)$$

$(X_i, Y_i)_{i \in \mathbb{N}}$  is a sequence of i.i.d. random variables, on an underlying probability space  $(\Omega, \mathfrak{A}, P)$ , taking values in  $\mathcal{X} \times \mathcal{Y}$  and distributed according to  $P$ . For every  $n \in \mathbb{N}^*$ ,  $Z_n = (X_i, Y_i)_{1 \leq i \leq n}$ . The sequence  $(\varepsilon_n)_{n \in \mathbb{N}}$  is in  $[0, 1]$  and  $\varepsilon_n \rightarrow 0$  as  $n \rightarrow +\infty$ . Moreover, for every  $n \in \mathbb{N}^*$ , every  $\lambda \in \mathbb{R}_{++}$ , and every training set  $z_n = (x_i, y_i)_{1 \leq i \leq n} \in (\mathcal{X} \times \mathcal{Y})^n$

$$\hat{u}_{n, \lambda}(z_n) \in \underset{u \in \ell^2(\mathbb{K})}{\text{Argmin}}^{\varepsilon_n} \left( \frac{1}{n} \sum_{i=1}^n |f_u(x_i) - y_i|^2 + \lambda G(u) \right). \quad (2.6)$$

*Remark 1*

- (i) The conditions on the sequences  $((\inf D_k)_+)_{k \in \mathbb{K}}$  and  $((\max D_k)_-)_{k \in \mathbb{K}}$  given in Assumption 1(c) ensure that  $G \in \Gamma_0(\ell^2(\mathbb{K}))$ . Moreover,  $\text{dom } G \subset \ell^r(\mathbb{K})$  and  $G$  is bounded from below and coercive (see Lemma 7).
- (ii) It follows from (2.2) that  $A$  is a bounded linear operator such that  $\|A\| \leq \kappa$  and  $\text{ran } A \subset L^\infty(P_{\mathcal{X}})$ . Moreover, when viewed as an operator from  $\ell^2(\mathbb{K})$  to  $\mathbb{R}^{\mathcal{X}}$ ,  $A$  is continuous with respect to the topology of the pointwise convergence on  $\mathbb{R}^{\mathcal{X}}$ . The feature map  $\Phi$  and  $A$  are connected via the identities

$$(\forall k \in \mathbb{K})(\forall x \in \mathcal{X}) \quad \langle \Phi(x) | e_k \rangle = (Ae_k)(x). \quad (2.7)$$

In [16, Proposition 3] it is shown that  $\text{ran } A$  can be endowed with a reproducing kernel Hilbert space structure for which  $A$  becomes a partial isometry, and the corresponding reproducing kernel is

$$K: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}: (x, x') \mapsto \sum_{k \in \mathbb{K}} \phi_k(x) \phi_k(x'). \quad (2.8)$$

In the above setting, the goal is to minimize the risk  $R$  of (2.1) on the closed convex subset  $\mathcal{C}$  of  $L^2(P_{\mathcal{X}})$  using the  $n$  i.i.d. observations  $Z_n = (X_i, Y_i)_{1 \leq i \leq n}$ . In this respect, recall that the regression function  $f^\dagger$  is the minimizer of the risk on  $L^2(P_{\mathcal{X}})$  [23] and that

$$(\forall f \in L^2(P_{\mathcal{X}})) \quad R(f) - \inf R(L^2(P_{\mathcal{X}})) = \|f - f^\dagger\|_{L^2}^2. \quad (2.9)$$

This means that minimizing  $R$  on  $L^2(P_{\mathcal{X}})$  is equivalent to approximating the regression function  $f^\dagger$ . In our constrained setting, the solution to the regression problem on  $\mathcal{C}$  results in a target function  $f_{\mathcal{C}}$  with the following properties.

**Proposition 1** *Suppose that Assumption 1 is in force. Then there exists a unique  $f_{\mathcal{C}} \in \mathcal{C}$  such that  $R(f_{\mathcal{C}}) = \inf R(\mathcal{C})$ . Moreover, the following hold:*

- (i)  $f_{\mathcal{C}}$  is the projection of  $f^\dagger$  onto  $\mathcal{C}$  in  $L^2(P_{\mathcal{X}})$ .
- (ii)  $(\forall f \in \mathcal{C}) \quad \|f - f_{\mathcal{C}}\|_{L^2}^2 \leq R(f) - \inf R(\mathcal{C})$ .
- (iii)  $(\forall f \in \mathcal{C}) \quad R(f) - \inf R(\mathcal{C}) \leq 2 \left[ (\|f - f_{\mathcal{C}}\|_{L^2} + \sqrt{\inf R(\mathcal{C}) - \inf R(L^2(P_{\mathcal{X}}))})^2 + \inf R(L^2(P_{\mathcal{X}})) \right]^{1/2} \|f - f_{\mathcal{C}}\|_{L^2}$ .

Proposition 1 states that, as in the unconstrained case, minimizing the risk over  $\mathcal{C}$  is still equivalent to approaching  $f_{\mathcal{C}}$  in  $L^2(P_{\mathcal{X}})$ . It is worth noting that we do not assume that  $f_{\mathcal{C}} = f_u$  for some  $u \in \text{dom} G$ , since the infimum of  $R$  on  $A(\text{dom} G)$  may not be attained. A *consistent learning scheme* generates a random variable  $\hat{u}_{n, \lambda_n}(Z_n)$ , taking values in  $\ell^2(\mathbb{K})$ , from  $n$  i.i.d. observations  $Z_n = (X_i, Y_i)_{1 \leq i \leq n}$ , so that the resulting sequence of random functions  $(\hat{f}_n)_{n \in \mathbb{N}} = (A\hat{u}_{n, \lambda_n}(Z_n))_{n \in \mathbb{N}}$  is *weakly consistent* in the sense that

$$R(\hat{f}_n) \rightarrow \inf R(\mathcal{C}) \text{ in probability, i.e., } \|\hat{f}_n - f_{\mathcal{C}}\|_{L^2} \rightarrow 0 \text{ in probability,} \quad (2.10)$$

or *strongly consistent* in the sense that

$$R(\hat{f}_n) \rightarrow \inf R(\mathcal{C}) \text{ P-a.s., i.e., } \|\hat{f}_n - f_{\mathcal{C}}\|_{L^2} \rightarrow 0 \text{ P-a.s.} \quad (2.11)$$

Next, we provide sufficient conditions on the regularization parameters  $(\lambda_n)_{n \in \mathbb{N}}$  and on the errors  $(\varepsilon_n)_{n \in \mathbb{N}}$ , that guarantee consistency and then present an algorithm to compute the proposed estimators.

**Theorem 2** *Suppose that Assumption 1 is in force and let  $f_{\mathcal{C}}$  be defined as in Proposition 1. Let  $(\lambda_n)_{n \in \mathbb{N}}$  be a sequence in  $]0, +\infty[$  converging to 0 and, for every  $n \in \mathbb{N}$ , let  $\hat{f}_n = A\hat{u}_{n, \lambda_n}(Z_n)$ . Then the following hold:*

- (i) *Suppose that  $\varepsilon_n / \lambda_n^{4/r} \rightarrow 0$  and that  $1 / (\lambda_n^{2/r} n^{1/2}) \rightarrow 0$ . Then  $(\hat{f}_n)_{n \in \mathbb{N}}$  is weakly consistent, i.e.,  $\|\hat{f}_n - f_{\mathcal{C}}\|_{L^2} \rightarrow 0$  in probability.*
- (ii) *Suppose that  $\varepsilon_n = O(1/n)$  and that  $(\log n) / (\lambda_n^{2/r} n^{1/2}) \rightarrow 0$ . Then  $(\hat{f}_n)_{n \in \mathbb{N}}$  is strongly consistent, i.e.,  $\|\hat{f}_n - f_{\mathcal{C}}\|_{L^2} \rightarrow 0$  P-a.s.*
- (iii) *Suppose that  $f_{\mathcal{C}} \in A(\text{dom} G)$  and set  $S = \text{Argmin}_{\text{dom} G} F$ . Then there exists a unique  $u^\dagger \in S$  which minimizes  $G$  over  $S$  and  $Au^\dagger = f_{\mathcal{C}}$ . Moreover, the following hold:*

(a) Suppose that  $\varepsilon_n/\lambda_n^2 \rightarrow 0$  and that  $1/(\lambda_n n^{1/2}) \rightarrow 0$ . Then

$$\|\widehat{u}_{n,\lambda_n}(Z_n) - u^\dagger\|_r \rightarrow 0 \quad \text{in probability.} \quad (2.12)$$

(b) Suppose that  $\varepsilon_n = O(1/n)$  and that  $(\log n)/(\lambda_n n^{1/2}) \rightarrow 0$ . Then

$$\|\widehat{u}_{n,\lambda_n}(Z_n) - u^\dagger\|_r \rightarrow 0 \quad \text{P-a.s.} \quad (2.13)$$

*Remark 2*

- (i) In Theorem 2(i)-(ii) the weakest conditions on the regularization parameters  $(\lambda_n)_{n \in \mathbb{N}}$  occur when  $r = 2$ , whereas, in Theorem 2(iii), the consistency conditions do not depend on the exponent  $r$ .
- (ii) Under the hypotheses of Theorem 2(iii), consistency extends to the sequence of coefficients  $(\widehat{u}_{n,\lambda_n}(Z_n))_{n \in \mathbb{N}}$ . This is relevant when one requires the sparsity pattern of the estimators to approximate that of  $u^\dagger$ . We note that, under further hypotheses on the  $D_k$ 's, both  $\widehat{u}_{n,\lambda_n}(Z_n)$  and  $u^\dagger$  have finite support. See Remark 6(v).

*Remark 3*

- (i) In the special case when, in (1.4), for every  $k \in \mathbb{K}$ ,  $h_k = \eta|\cdot|^2$ ,  $C_k = \mathbb{R}$ ,  $D_k = [-\omega_k, \omega_k]$ , for some  $\omega_k \in \mathbb{R}_+$ , we recover the elastic net framework of [16] and the same consistency conditions as in [16, Theorem 2 and Theorem 3]. This special case yields a strongly convex problem. In our general setting, the exponent  $r$  may take any value in  $]1, 2]$  and the objective function is only totally convex on bounded sets (see Lemma 1).
- (ii) Consistency results have been obtained in [20] for a regularizer composed of the sum of the indicator function of an  $\ell^1$  ball and the squared  $\ell^2$  norm. However, such a regularizer is not separable and hence it is not included in the present study.
- (iii) When  $\mathbb{K}$  is finite and, for every  $k \in \mathbb{K}$ ,  $g_k = |\cdot|^r$ , [25] provides an excess risk bound depending on the cardinality of  $\mathbb{K}$  and the level of sparsity of  $u^\dagger$  (see also [22]). The case  $r = 1$  has been considered in [14].
- (iv) Similar consistency results can also be derived using [11, Corollary 4.6], where a general loss function, a more general penalty function, and an arbitrary Banach space are considered. However, due to the generality of the analysis in [11], the conditions imposed on the sequence of parameters  $(\lambda_n)_{n \in \mathbb{N}}$  in [11] are more restrictive, and imply a slower worst-case convergence rate. Here, exploiting the structure of the square loss and the regularizer, we obtain a sharper result and a simpler statement.

We now address the algorithmic aspects. The objective function in (2.6) consists of a smooth (quadratic) data fitting term and a separable nondifferentiable convex term, penalizing each dictionary coefficient individually. Thus a natural choice is to consider the forward-backward splitting algorithm [12]. We stress that, since  $\varepsilon$ -minimizers are employed in (2.6), algorithms that provide minimizing sequences are necessary. Moreover, due to the term  $h_k$  in the regularizer, the proximity operator of  $G$  may not be computable explicitly. Consequently, convergence results in objective function values for an error-tolerant forward-backward algorithm are in order. Note that, to the best of our knowledge, the only work addressing the above issue is [32].

However, only ergodic convergence is established, which is a result weaker than that obtained in the error-free case [4, 15] and certainly not satisfactory when sparsity is concerned, since averaging blurs the structural properties of the iterates. Nesterov-like [28] variants of the inexact forward-backward splitting algorithm may also be suitable for computing the estimators (2.6) to the extent that they also generate minimizing sequences [32, 34]. However, in practice, they sometimes may be slower than the standard version since they are more sensitive to errors [34].

In the Theorem 3 below, we advance the convergence theory on the standard forward-backward splitting algorithm by establishing an  $o(1/m)$  rate of convergence in objective values with relaxation, variable proximal parameters, and in the presence of the following type of errors in the numerical evaluation of the proximity operator [31, 32, 34].

**Definition 1** Let  $\mathcal{H}$  be a real Hilbert space, let  $J \in \Gamma_0(\mathcal{H})$ , let  $(u, w) \in \mathcal{H}^2$ , and let  $\delta \in \mathbb{R}_+$ . The notation  $u \simeq_\delta \text{prox}_J w$  means that

$$J(u) + \frac{1}{2}\|u - w\|^2 \leq \min_{v \in \mathcal{H}} \left( J(v) + \frac{1}{2}\|v - w\|^2 \right) + \frac{\delta^2}{2}. \quad (2.14)$$

*Remark 4* Note that  $u \in \text{dom} J$  and, since  $J + (1/2)\|\cdot - w\|^2$  is 1-strongly convex, we have  $u = \text{prox}_J w + a$  with  $\|a\| \leq \delta$  [31]. Thus the errors considered in Definition 1 are additive and generate perturbations of the proximal point that are feasible. This feature is critical when convergence of the objective function values is under consideration.

**Theorem 3** Let  $\mathcal{H}$  be a real Hilbert space, let  $J_1: \mathcal{H} \rightarrow \mathbb{R}$  be a convex differentiable function with a  $\beta$ -Lipschitz continuous gradient for some  $\beta \in \mathbb{R}_{++}$ . Let  $J_2 \in \Gamma_0(\mathcal{H})$ , set  $J = J_1 + J_2$ , and suppose that  $\text{Argmin} J \neq \emptyset$ . Let  $(\gamma_m)_{m \in \mathbb{N}}$  be a sequence in  $\mathbb{R}_{++}$  such that  $0 < \inf_{m \in \mathbb{N}} \gamma_m \leq \sup_{m \in \mathbb{N}} \gamma_m < 2/\beta$ , let  $(\tau_m)_{m \in \mathbb{N}}$  be a sequence in  $]0, 1]$ , such that  $\inf_{m \in \mathbb{N}} \tau_m > 0$ . Let  $(\delta_m)_{m \in \mathbb{N}}$  be a summable sequence in  $\mathbb{R}_+$  and let  $(b_m)_{m \in \mathbb{N}}$  be a summable sequence in  $\mathcal{H}$ . Fix  $u_0 \in \mathcal{H}$  and set

$$\begin{aligned} & \text{for } m = 0, 1, \dots \\ & \begin{cases} v_m \simeq_{\delta_m} \text{prox}_{\gamma_m J_2}(u_m - \gamma_m(\nabla J_1(u_m) + b_m)) \\ u_{m+1} = u_m + \tau_m(v_m - u_m). \end{cases} \end{aligned} \quad (2.15)$$

Then the following hold:

- (i)  $(u_m)_{m \in \mathbb{N}}$  converges weakly to a point in  $\text{Argmin} J$ .
- (ii) For every  $u \in \text{Argmin} J$ ,  $\sum_{m \in \mathbb{N}} \|\nabla J_1(u_m) - \nabla J_1(u)\|^2 < +\infty$ .
- (iii)  $\sum_{m \in \mathbb{N}} \|v_m - u_m\|^2 < +\infty$ .
- (iv)  $J(u_m) \rightarrow \inf J(\mathcal{H})$  and  $\sum_{m \in \mathbb{N}} |J(v_m) - \inf J(\mathcal{H})|^2 < +\infty$ .
- (v) Suppose that  $\sum_{m \in \mathbb{N}} (1 - \tau_m) < +\infty$ . Then

$$\sum_{m \in \mathbb{N}} (J(v_m) - \inf J(\mathcal{H})) < +\infty \quad \text{and} \quad \sum_{m \in \mathbb{N}} (J(u_m) - \inf J(\mathcal{H})) < +\infty.$$

- (vi) Suppose that  $\sum_{m \in \mathbb{N}} (1 - \tau_m) < +\infty$ ,  $\sum_{m \in \mathbb{N}} m \delta_m < +\infty$ , and  $\sum_{m \in \mathbb{N}} m \|b_m\| < +\infty$ . Then  $J(u_m) - \inf J(\mathcal{H}) = o(1/m)$ .



*Remark 5*

- (i) In [4], the rate  $O(1/m)$  for objective values is proved in the error-free case and no relaxations ( $\delta_m \equiv 0$  and  $\tau_m \equiv 1$ ), assuming that  $J_1 + J_2$  is coercive. On the other hand, an  $o(1/m)$  rate on the objective values was derived in [15] in the special case of a fixed proximal parameter  $\gamma \in ]0, 2/\beta[$ , no relaxation, and no errors.
- (ii) In [32] no relaxation is considered and the proximal parameters  $(\gamma_m)_{m \in \mathbb{N}}$  are fixed to a constant value and limited to  $1/\beta$ . Moreover, only ergodic convergence is proved.

In general the criterion considered in Definition 1 is not explicitly verifiable since the minimum is not known. This is the reason why another type of errors is considered in [34]. However, the following result shows that when computing proximity operators of separable functions of the type considered in (1.3)–(1.4), errors of type of Definition 1 arise, and they can be explicitly checked and implemented in practice.

**Proposition 2** *Let  $\mathcal{H}$  be a separable real Hilbert space and let  $(e_k)_{k \in \mathbb{K}}$  be an orthonormal basis of  $\mathcal{H}$ , where  $\mathbb{K}$  is an at most countable set. Let  $(h_k)_{k \in \mathbb{K}}$  be a family of convex functions from  $\mathbb{R}$  to  $\mathbb{R}$  such that, for every  $k \in \mathbb{K}$ ,  $h_k \geq h_k(0) = 0$ . Let  $(C_k)_{k \in \mathbb{K}}$  be a family of closed intervals in  $\mathbb{R}$  such that  $0 \in \bigcap_{k \in \mathbb{K}} C_k$ , let  $(D_k)_{k \in \mathbb{K}}$  be a family of nonempty closed bounded intervals in  $\mathbb{R}$ . Suppose that  $(h_k^*(-\min D_k)_+)$  and  $(h_k^*((\max D_k)_-))_{k \in \mathbb{K}}$  are summable, and set*

$$G: \mathcal{H} \rightarrow ]-\infty, +\infty]: u \mapsto \sum_{k \in \mathbb{K}} (\iota_{C_k} + \sigma_{D_k} + h_k)(\langle u | e_k \rangle). \quad (2.16)$$

Let  $\gamma \in \mathbb{R}_{++}$ , let  $w \in \mathcal{H}$ , let  $(\alpha_k)_{k \in \mathbb{K}} \in \mathbb{R}^{\mathbb{K}}$ , let  $(\xi_k)_{k \in \mathbb{K}} \in \ell^1(\mathbb{K})$ , set  $\delta = \sqrt{\sum_{k \in \mathbb{K}} \xi_k}$ , and let

$$\begin{array}{l} \text{for every } k \in \mathbb{K} \\ \left[ \begin{array}{l} \chi_k = \langle w | e_k \rangle \\ |\alpha_k| \leq \frac{\xi_k}{4\gamma \max\{h_k(|\chi_k| + 2), h_k(-|\chi_k| - 2)\} + 2|\chi_k| + 1} \\ \pi_k = \text{prox}_{\gamma h_k}(\text{soft}_{\gamma D_k} \chi_k) + \alpha_k \\ v_k = \text{proj}_{C_k}(\text{sign}(\chi_k) \max\{0, \text{sign}(\chi_k) \pi_k\}). \end{array} \right. \quad (2.17) \end{array}$$

Then  $v = (v_k)_{k \in \mathbb{K}} \in \ell^2(\mathbb{K})$  and  $v \simeq_{\delta} \text{prox}_{\gamma G} w$ .

*Remark 6*

- (i) The soft-thresholding operator with respect to a bounded interval  $D_k = [\underline{\omega}_k, \overline{\omega}_k] \subset \mathbb{R}$  is

$$(\forall \mu \in D_k) \quad \text{soft}_{D_k} \mu = \begin{cases} \mu - \overline{\omega}_k & \text{if } \mu > \overline{\omega}_k \\ 0 & \text{if } \mu \in D_k \\ \mu - \underline{\omega}_k & \text{if } \mu < \underline{\omega}_k. \end{cases} \quad (2.18)$$

- (ii) The function  $G$  in (2.16) is well-defined and lies in  $\Gamma_0(\mathcal{H})$  as it is the composition of the continuous linear operator  $\mathcal{H} \rightarrow \ell^2(\mathbb{K}): u \mapsto (\langle u | e_k \rangle)_{k \in \mathbb{K}}$  and the function

$$\ell^2(\mathbb{K}) \rightarrow ]-\infty, +\infty]: (\mu_k)_{k \in \mathbb{K}} \mapsto \sum_{k \in \mathbb{K}} g_k(\mu_k), \quad \text{with } g_k = \iota_{C_k} + \sigma_{D_k} + h_k, \quad (2.19)$$

which is a well-defined function in  $\Gamma_0(\ell^2(\mathbb{K}))$  by Lemma 7.

- (iii) The conditions on  $(\min D_k)_{k \in \mathbb{N}}$  and  $(\max D_k)_{k \in \mathbb{N}}$  appearing in Proposition 2 are weaker than those considered in Assumption 1(c). See Lemma 7.
- (iv) In algorithm (2.17), the computation of  $\text{prox}_{\gamma h_k}$  tolerates an error  $\alpha_k$ . This is necessary since, in general, the proximity operator may not be computable explicitly. In such instances,  $\text{prox}_{\gamma h_k}$  must be computed iteratively (e.g., by bisection) and the bound on  $|\alpha_k|$  in (2.17) gives an explicit stopping rule for the iterations. In Appendix B, the case  $h_k = \eta_k |\cdot|^r$ , with  $r > 1$ , is further analyzed.
- (v) In algorithm (2.17), for every  $k \in \mathbb{K}$ , we have  $\text{sign } v_k = \text{sign } \chi_k$ .
- (vi) As in [16, Corollary 3] one proves that if  $\sup_{k \in \mathbb{K}} \min D_k < 0 < \inf_{k \in \mathbb{K}} \max D_k$ , then  $\{k \in \mathbb{K} \mid v_k \neq 0\}$  is finite. Similarly, if  $\sup_{k \in \mathbb{K}} \min D_k \leq 0 < \inf_{k \in \mathbb{K}} \max D_k$ , then  $\{k \in \mathbb{K} \mid v_k > 0\}$  is finite, so sparsity is enforced only on the positive coefficients.

We now present an inexact forward-backward algorithm to solve problem (1.3) which combines algorithms (2.15) and (2.17).

**Algorithm 4** Let  $(\gamma_m)_{m \in \mathbb{N}}$  be a sequence in  $\mathbb{R}_{++}$  such that  $0 < \inf_{m \in \mathbb{N}} \gamma_m \leq \sup_{m \in \mathbb{N}} \gamma_m < \lambda / \kappa^2$ , let  $(\tau_m)_{m \in \mathbb{N}}$  be a sequence in  $]0, 1]$  such that  $\inf_{m \in \mathbb{N}} \tau_m > 0$ . Let  $(b_m)_{m \in \mathbb{N}} = ((\beta_{m,k})_{k \in \mathbb{K}})_{m \in \mathbb{N}} \in (\ell^2(\mathbb{K}))^{\mathbb{N}}$  be such that  $\sum_{m \in \mathbb{N}} \|b_m\| < +\infty$ , let  $\zeta \in \mathbb{R}_{++}$ , let  $p \in ]1, +\infty[$ , and let  $(\xi_k)_{k \in \mathbb{K}} \in \ell^1(\mathbb{K})$ . Fix  $(\mu_{0,k})_{k \in \mathbb{K}} \in \ell^2(\mathbb{K})$  and iterate

$$\begin{array}{l}
 \text{for } m = 0, 1, \dots \\
 \quad \text{for every } k \in \mathbb{K} \\
 \quad \left| \begin{array}{l}
 \chi_{m,k} = \mu_{m,k} - \frac{\gamma_m}{\lambda} \left( \frac{2}{n} \sum_{i=1}^n \left( \sum_{j \in \mathbb{K}} \mu_{m,j} \phi_j(x_i) - y_i \right) \phi_k(x_i) + \beta_{m,k} \right) \\
 |\alpha_{m,k}| \leq \frac{\zeta m^{-2p} \xi_k}{4\gamma_m \max\{h_k(|\chi_{m,k}| + 2), h_k(-|\chi_{m,k}| - 2)\}} + 2|\chi_{m,k}| + 1 \\
 \pi_{m,k} = \text{prox}_{\gamma_m h_k}(\text{soft}_{\gamma_m D_k} \chi_{m,k}) + \alpha_{m,k} \\
 v_{m,k} = \text{proj}_{C_k}(\text{sign}(\chi_{m,k}) \max\{0, \text{sign}(\chi_{m,k}) \pi_{m,k}\}) \\
 \mu_{m+1,k} = \mu_{m,k} + \tau_m (v_{m,k} - \mu_{m,k}).
 \end{array} \right. \quad (2.20)
 \end{array}$$

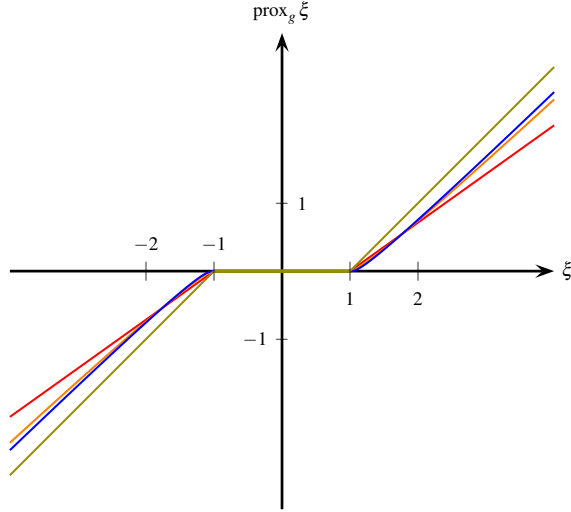
*Remark 7*

- (i) An attractive feature of Algorithm 4 is that, at each iteration, each component of the functions in (2.5) is activated componentwise and individually.
- (ii) The freedom in the choice of the intervals  $(D_k)_{k \in \mathbb{K}}$ ,  $(C_k)_{k \in \mathbb{K}}$ , and of the exponent  $r$  provides flexibility in setting the type of thresholding operation. It is in particular possible to promote selective sparsity. Figures 1 and 2 show a few examples.

**Theorem 5** *Suppose that Assumption 1 is in force. Call*

$$J: \ell^2(\mathbb{K}) \rightarrow ]-\infty, +\infty]: u = (\mu_k)_{k \in \mathbb{K}} \mapsto \frac{1}{n} \sum_{i=1}^n |f_u(x_i) - y_i|^2 + \lambda \sum_{k \in \mathbb{K}} g_k(\mu_k) \quad (2.21)$$

*the objective function in (2.6), and let  $(u_m)_{m \in \mathbb{N}} = ((\mu_{m,k})_{k \in \mathbb{K}})_{m \in \mathbb{N}}$  and  $(v_m)_{m \in \mathbb{N}} = ((v_{m,k})_{k \in \mathbb{K}})_{m \in \mathbb{N}}$  be the sequences generated by Algorithm 4. Then the following hold:*



**Fig. 1** Soft thresholding (green) and  $\text{prox}_g$  for  $g = |\cdot| + 0.2|\cdot|^r$ , with  $r = 2$  (red),  $r = 3/2$  (orange),  $r = 4/3$  (blue).

- (i)  $J$  has a unique minimizer  $\hat{u}$ , and  $\hat{u} \in \ell^r(\mathbb{K})$ .  
(ii)  $\sum_{m \in \mathbb{N}} |J(v_m) - \inf J(\ell^2(\mathbb{K}))|^2 < +\infty$ ,  $J(u_m) \rightarrow \inf J(\ell^2(\mathbb{K}))$ ,  $\|v_m - \hat{u}\|_r \rightarrow 0$ , and  $\|u_m - \hat{u}\|_r \rightarrow 0$  as  $m \rightarrow +\infty$ . Moreover

$$\|v_m - \hat{u}\|_r = O\left(\sqrt{J(v_m) - \inf J(\ell^2(\mathbb{K}))}\right) \quad (2.22)$$

and

$$\|u_m - \hat{u}\|_r = O\left(\sqrt{J(u_m) - \inf J(\ell^2(\mathbb{K}))}\right). \quad (2.23)$$

- (iii) Suppose that  $\sum_{m \in \mathbb{N}} (1 - \tau_m) < +\infty$ . Then

$$\sum_{m \in \mathbb{N}} (J(v_m) - \inf J(\ell^2(\mathbb{K}))) < +\infty \quad \text{and} \quad \sum_{m \in \mathbb{N}} (J(u_m) - \inf J(\ell^2(\mathbb{K}))) < +\infty.$$

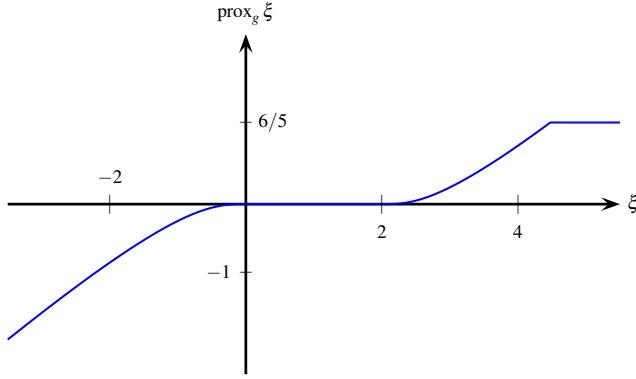
- (iv) Suppose that  $p > 2$ , that  $\sum_{m \in \mathbb{N}} (1 - \tau_m) < +\infty$  and  $\sum_{m \in \mathbb{N}} m \|b_m\| < +\infty$ . Then

$$J(u_m) - \inf J(\ell^2(\mathbb{K})) = o(1/m) \quad \text{and} \quad \|u_m - \hat{u}\|_r = o(1/\sqrt{m}). \quad (2.24)$$

### 3 Statistical analysis

Throughout the section Assumption 1 is made. Our main objective is to prove Theorem 2. To this end, we first observe that, setting  $(\hat{f}_n)_{n \in \mathbb{N}} = (A\hat{u}_{n,\lambda_n}(Z_n))_{n \in \mathbb{N}}$  and using Proposition 1(ii), we have

$$\begin{aligned} (\forall n \in \mathbb{N}) \quad \|\hat{f}_n - f_{\mathcal{E}}\|_{L^2} &\leq \|A\hat{u}_{n,\lambda_n}(Z_n) - Au_{\lambda_n}\|_{L^2} + \|Au_{\lambda_n} - f_{\mathcal{E}}\|_{L^2} \\ &\leq \|A\| \|\hat{u}_{n,\lambda_n}(Z_n) - u_{\lambda_n}\|_r + \sqrt{F(u_{\lambda_n}) - \inf R(\mathcal{E})}. \end{aligned} \quad (3.1)$$



**Fig. 2**  $\text{prox}_g$  for  $g = \iota_{]-\infty, 6/5]} + \sigma_{[0, 2]} + 0.9|\cdot|^{4/3}$ .

This suggests to study separately the convergence of the two terms on the right-hand side of (3.1).

*Proof (of Proposition 1)* (i): For every  $f \in \mathcal{C}$ ,  $R(f) = \|f - f^\dagger\|_{L^2}^2 + \inf R(L^2(P_{\mathcal{X}}))$ . Therefore, minimizing  $R$  over  $\mathcal{C}$  turns to find the element of  $\mathcal{C}$  which is nearest to  $f^\dagger$  in  $L^2(P_{\mathcal{X}})$ .

(ii): It follows from (i), that  $\inf R(\mathcal{C}) = \|f_{\mathcal{C}} - f^\dagger\|_{L^2}^2 + \inf R(L^2(P_{\mathcal{X}}))$ . Therefore, since for every  $f \in \mathcal{C}$ ,  $\langle f - f_{\mathcal{C}} | f^\dagger - f_{\mathcal{C}} \rangle \leq 0$ , we have  $R(f) - \inf R(\mathcal{C}) = \|f - f^\dagger\|_{L^2}^2 - \|f_{\mathcal{C}} - f^\dagger\|_{L^2}^2 = \|f - f_{\mathcal{C}}\|_{L^2}^2 + 2\langle f - f_{\mathcal{C}} | f_{\mathcal{C}} - f^\dagger \rangle \geq \|f - f_{\mathcal{C}}\|_{L^2}^2$ .

(iii): Let  $f \in \mathcal{C}$ . Using the fact that, for every  $(a, b, c) \in \mathbb{R}_+^3$  with  $a \geq b$ ,  $\sqrt{a+c} - \sqrt{b+c} \leq \sqrt{a} - \sqrt{b}$ , we derive that

$$\sqrt{R(f)} - \sqrt{\inf R(\mathcal{C})} \leq \|f - f^\dagger\|_{L^2} - \|f_{\mathcal{C}} - f^\dagger\|_{L^2} \leq \|f - f_{\mathcal{C}}\|_{L^2}. \quad (3.2)$$

Therefore, using the inequality  $a^2 - b^2 \leq 2a(a - b)$ , we obtain

$$\begin{aligned} R(f) - \inf R(\mathcal{C}) &\leq 2\sqrt{R(f)} \|f - f_{\mathcal{C}}\|_{L^2} \\ &= 2\sqrt{\|f - f^\dagger\|_{L^2}^2 + \inf R(L^2(\mathcal{X}))} \|f - f_{\mathcal{C}}\|_{L^2} \\ &\leq 2\left(\|f - f_{\mathcal{C}}\|_{L^2} + \|f_{\mathcal{C}} - f^\dagger\|_{L^2}\right)^2 + \inf R(L^2(\mathcal{X})) \Big)^{1/2} \|f - f_{\mathcal{C}}\|_{L^2} \\ &= 2\left(\|f - f_{\mathcal{C}}\|_{L^2} + \sqrt{\inf_{\mathcal{C}} R - \inf_{L^2(P_{\mathcal{X}})} R} + \inf R(L^2(\mathcal{X}))\right)^{1/2} \|f - f_{\mathcal{C}}\|_{L^2} \end{aligned} \quad (3.3)$$

The following result establishes that  $G$  is totally convex [6] on bounded subsets of  $\ell^r(\mathbb{K})$  and gives an explicit lower bound for the relative modulus of total convexity.

**Lemma 1** *Suppose that Assumption 1 is in force. Let  $\rho \in \mathbb{R}_{++}$ , let  $u_0 \in \ell^r(\mathbb{K})$  be such that  $\|u_0\|_r \leq \rho$ , let  $u_0^* \in \partial G(u_0)$ , and set  $M = (7/32)r(r-1)(1 - (2/3)^{r-1})$ . Then*

$$(\forall u \in \ell^r(\mathbb{K})) \quad G(u) - G(u_0) \geq \langle u - u_0 \mid u_0^* \rangle + \frac{\eta M \|u - u_0\|_r^2}{(\rho + \|u - u_0\|_r)^{2-r}}. \quad (3.4)$$

*Proof* Let  $G_\downarrow$  be the restriction of  $G$  to  $\ell^r(\mathbb{K})$ , endowed with the norm  $\|\cdot\|_r$ . Since  $u_0 \in \ell^r(\mathbb{K})$  and  $u_0^* \in \ell^{r^*}(\mathbb{K})$ , we have that  $u_0^* \in \partial G_\downarrow(u_0)$ . Let  $\psi$  be the modulus of total convexity of  $G_\downarrow$  and let  $\varphi$  be the modulus of total convexity of  $\|\cdot\|_r^r$  in  $\ell^r(\mathbb{K})$ . Then, for every  $u \in \ell^r(\mathbb{K})$ ,  $G(u) - G(u_0) \geq \langle u - u_0, u_0^* \rangle + \psi(u_0; \|u - u_0\|_r)$ . Moreover, since  $G_\downarrow = H + \eta \|\cdot\|_r^r$ , with  $H \in \Gamma_0(\ell^r(\mathbb{K}))$  (see Lemma 7), we have  $\psi \geq \eta \varphi$ . The statement follows from [11, Proposition A.9-Remark A.10].

The next proposition concerns the second term in the right-hand side of (3.1). It revisits some results of [2] about Tikhonov-like regularization specialized to our setting.

**Proposition 3** *Suppose that Assumption 1 is in force. For every  $(\lambda, \varepsilon) \in \mathbb{R}_{++} \times \mathbb{R}_+$ , let  $u_{\lambda, \varepsilon}$  be an  $\varepsilon$ -minimizer of  $F + \lambda G$  and let  $u_G$  be the minimizer of  $G$ . Let  $M \in \mathbb{R}_{++}$  be defined as in Lemma 1. Then the following hold:*

- (i)  $\inf R(\mathcal{C}) = \inf F(\text{dom } G)$ .
- (ii)  $(\forall (\lambda, \varepsilon) \in \mathbb{R}_{++} \times \mathbb{R}_+) \|u_{\lambda, \varepsilon} - u_G\|_r \leq \max \{ \|u_G\|_r, (2(F(u_G) + \varepsilon)/(\eta M \lambda))^{1/r} \}$ .
- (iii)  $F(u_{\lambda, \varepsilon}) \rightarrow \inf F(\text{dom } G)$  as  $(\lambda, \varepsilon) \rightarrow (0^+, 0^+)$ .
- (iv) *Suppose that  $S = \text{Argmin}_{\text{dom } G} F \neq \emptyset$ . Then there exists  $u^\dagger \in \ell^r(\mathbb{K})$  such that  $\text{Argmin}_S G = \{u^\dagger\}$  and  $u_{\lambda, 0} \rightarrow u^\dagger$  as  $\lambda \rightarrow 0^+$ .*

*Proof* We first note that it follows from Remark 1(i) that  $G$  has a minimizer.

(i): We prove that  $\mathcal{C} = \overline{A(\text{dom } G)}$  and then the statement will follow. Let  $u = (\mu_k)_{k \in \mathbb{K}} \in \ell^2(\mathbb{K}) \cap \times_{k \in \mathbb{K}} C_k$  and take  $\delta \in \mathbb{R}_{++}$ . Then there exists a finite set  $\mathbb{K}_1 \subset \mathbb{K}$  such that  $\sum_{k \in \mathbb{K} \setminus \mathbb{K}_1} |\mu_k|^2 < \delta^2$ . Now let  $v = (v_k)_{k \in \mathbb{K}}$  be such that, for every  $k \in \mathbb{K}_1$ ,  $v_k = \mu_k$  and, for every  $k \in \mathbb{K} \setminus \mathbb{K}_1$ ,  $v_k = 0$ . We have  $v \in \text{dom } G$  and  $\|Au - Av\|_{L^2} < \|A\| \delta$ .

(ii): Let  $(\lambda, \varepsilon) \in \mathbb{R}_{++}^2$ . We derive from the definition of  $u_{\lambda, \varepsilon}$ , that  $F(u_{\lambda, \varepsilon}) + \lambda G(u_{\lambda, \varepsilon}) \leq F(u_G) + \lambda G(u_G) + \varepsilon$  hence, since  $0 \in \partial G(u_G)$ , it follows from Lemma 1 and the fact that  $F(u_{\lambda, \varepsilon}) \geq 0$ , that

$$\frac{\eta M \|u_{\lambda, \varepsilon} - u_G\|_r^2}{(\|u_G\|_r + \|u_{\lambda, \varepsilon} - u_G\|_r)^{2-r}} \leq G(u_{\lambda, \varepsilon}) - G(u_G) \leq \frac{F(u_G) + \varepsilon}{\lambda}. \quad (3.5)$$

If  $\|u_{\lambda, \varepsilon} - u_G\|_r \geq \|u_G\|_r$ , then

$$\frac{\eta M \|u_{\lambda, \varepsilon} - u_G\|_r^2}{(\|u_G\|_r + \|u_{\lambda, \varepsilon} - u_G\|_r)^{2-r}} \geq \frac{\eta M \|u_{\lambda, \varepsilon} - u_G\|_r^2}{(2\|u_{\lambda, \varepsilon} - u_G\|_r)^{2-r}} \geq \frac{\eta M}{2} \|u_{\lambda, \varepsilon} - u_G\|_r^r \quad (3.6)$$

and hence  $\|u_{\lambda, \varepsilon} - u_G\|_r^r \leq 2(F(u_G) + \varepsilon)/(\eta M \lambda)$ .

(iii): Let  $u \in \text{dom } G$ . Then

$$\begin{aligned}
\inf F(\text{dom } G) &\leq \underline{\lim}_{(\lambda, \varepsilon) \rightarrow (0,0)} F(u_{\lambda, \varepsilon}) \\
&\leq \overline{\lim}_{(\lambda, \varepsilon) \rightarrow (0,0)} F(u_{\lambda, \varepsilon}) \\
&\leq \overline{\lim}_{(\lambda, \varepsilon) \rightarrow (0,0)} (F(u) + \lambda(G(u) - G(u_G)) + \varepsilon) \\
&\leq F(u).
\end{aligned} \tag{3.7}$$

Since  $u$  is arbitrary, the statement follows.

(iv): Since  $S$  is convex and  $G \in \Gamma_0(\ell^2(\mathbb{K}))$  is strictly convex, coercive, and  $\text{dom } G \subset \ell^r(\mathbb{K})$ , it follows from [3, Corollary 11.16(ii)] that there exists a unique minimizer  $u^\dagger \in \ell^r(\mathbb{K})$  of  $G$  over  $S$ . Moreover, for every  $\lambda \in \mathbb{R}_{++}$ ,

$$G(u_{\lambda,0}) \leq (F(u^\dagger) - F(u_{\lambda,0}))/\lambda + G(u^\dagger) \leq G(u^\dagger) < +\infty, \tag{3.8}$$

which implies that  $(G(u_{\lambda,0}))_{\lambda \in \mathbb{R}_{++}}$  is bounded. Since  $G$  is coercive, the family  $(u_{\lambda,0})_{\lambda \in \mathbb{R}_{++}}$  is bounded as well, and it therefore has weak sequential cluster points. Next, we show that any such cluster point is necessarily equal to  $u^\dagger$ , which implies that  $u_\lambda \rightharpoonup u^\dagger$ . Indeed let  $(\lambda_n)_{n \in \mathbb{N}}$  be a vanishing sequence in  $\mathbb{R}_{++}$  and suppose that  $u_{\lambda_n} \rightharpoonup v$ , for some  $v \in \ell^2(\mathbb{K})$ . Then it follows from (iii) and (3.8) that

$$F(v) \leq \underline{\lim} F(u_{\lambda_n,0}) = \inf F(\text{dom } G) \quad \text{and} \quad G(v) \leq \underline{\lim} G(u_{\lambda_n,0}) \leq G(u^\dagger), \tag{3.9}$$

which implies that  $v \in S$  and  $v \in \text{Argmin}_S G = \{u^\dagger\}$ . However, thanks to Lemma 1,  $G$  is totally convex on bounded sets in  $\ell^r(\mathbb{K})$ . Therefore, by [37, Proposition 3.6.5],  $G$  is uniformly convex on bounded sets too. So, since  $(u_{\lambda,0})_{\lambda \in \mathbb{R}_{++}}$  is bounded, there exists an increasing function  $\phi: \mathbb{R}_+ \rightarrow \mathbb{R}_+$  such that

$$\begin{cases} \phi(0) = 0 \\ (\forall t \in \mathbb{R}_{++}) \quad \phi(t) > 0 \\ (\forall \lambda \in \mathbb{R}_{++}) \quad \phi\left(\frac{\|u_{\lambda,0} - u^\dagger\|}{2}\right) \leq \frac{G(u^\dagger) + G(u_{\lambda,0})}{2} - G\left(\frac{u_{\lambda,0} + u^\dagger}{2}\right). \end{cases} \tag{3.10}$$

In turn, we obtain

$$\overline{\lim}_{\lambda \rightarrow 0^+} \phi\left(\frac{\|u_{\lambda,0} - u^\dagger\|}{2}\right) \leq \overline{\lim}_{\lambda \rightarrow 0^+} \frac{G(u^\dagger) + G(u_{\lambda,0})}{2} + \overline{\lim}_{\lambda \rightarrow 0^+} (-G)\left(\frac{u_{\lambda,0} + u^\dagger}{2}\right) \tag{3.11}$$

and, arguing as in [8, Proof of Proposition 3.1(vi)], we get  $u_{\lambda,0} \rightarrow u^\dagger$  as  $\lambda \rightarrow 0^+$ .

We now address the convergence of the term  $\|\widehat{u}_{n,\lambda_n}(Z_n) - u_{\lambda_n}\|_r$  in (3.1) and hence the proof of Theorem 2. We first give a representer and stability theorem which generalizes existing results [18, 33] to our class of regularization functions.

**Theorem 6** *Suppose that Assumption 1 is in force. Set  $M = (7/32)r(r-1)(1 - (2/3)^{r-1})$ , let  $\lambda \in \mathbb{R}_{++}$ , and let  $u_\lambda \in \ell^r(\mathbb{K})$  be the minimizer of  $F + \lambda G$ . Then the following hold:*

(i) *The function*

$$\Psi_\lambda : \mathcal{X} \times \mathcal{Y} \rightarrow \ell^2(\mathbb{K}) : (x, y) \mapsto 2(f_{u_\lambda}(x) - y)\Phi(x) \quad (3.12)$$

is bounded and  $\|\Psi_\lambda\|_\infty \leq 2\kappa(\kappa\|u_\lambda\|_2 + b)$ . Moreover  $\|\Psi_\lambda\|_2 \leq 2\kappa\sqrt{R(f_{u_\lambda})}$  and  $-\mathbb{E}_P(\Psi_\lambda) \in \lambda\partial G(u_\lambda)$ .

(ii) Let  $n \in \mathbb{N}^*$ . Then there exists  $\widehat{v} \in \ell^r(\mathbb{K})$  such that  $\|\widehat{v} - \widehat{u}_{n,\lambda}(z_n)\|_r \leq \sqrt{\varepsilon_n}$

$$\frac{\eta M \|\widehat{v} - u_\lambda\|_r}{(\|u_\lambda\|_r + \|\widehat{v} - u_\lambda\|_r)^{2-r}} \leq \frac{1}{\lambda} \left( \left\| \frac{1}{n} \sum_{i=1}^n \Psi_\lambda(x_i, y_i) - \mathbb{E}_P(\Psi_\lambda) \right\|_2 + \sqrt{\varepsilon_n} \right). \quad (3.13)$$

*Proof* (i): First note that, by [3, Corollary 11.16(ii)] and Remark 1(i),  $u_\lambda$  is well defined since  $G$  is proper, lower semicontinuous, strictly convex, and coercive. Furthermore [3, Corollary 27.3(vi)] implies that  $-\nabla F(u_\lambda) \in \lambda\partial G(u_\lambda)$ . We derive from (2.3) that  $A^* : L^2(\mathcal{P}_{\mathcal{X}}) \rightarrow \ell^2(\mathbb{K}) : f \mapsto \mathbb{E}_{P_{\mathcal{X}}}(f\Phi)$ , and hence, since  $F = R \circ A$ ,

$$(\forall u \in \ell^2(\mathbb{K})) \quad \nabla F(u) = A^* \nabla R(f_u) = \mathbb{E}_P(\varphi), \quad (3.14)$$

where  $\varphi : (x, y) \mapsto 2(f_u(x) - y)\Phi(x)$ . Let  $(x, y) \in \mathcal{X} \times \mathcal{Y}$ . Then

$$|f_{u_\lambda}(x) - y| \leq |f_{u_\lambda}(x)| + |y| \leq \sum_{k \in \mathbb{K}} |\langle u_\lambda | e_k \rangle| |\phi_k(x)| + b \leq \kappa \|u_\lambda\|_2 + b \quad (3.15)$$

and hence  $\|\Psi_\lambda(x, y)\|_2 \leq 2|f_{u_\lambda}(x) - y| \|\Phi(x)\|_2 \leq 2(\kappa\|u_\lambda\|_2 + b)\kappa$ . Moreover,

$$\begin{aligned} \int_{\mathcal{X} \times \mathcal{Y}} \|\Psi_\lambda(x, y)\|_2^2 dP(x, y) &\leq \int_{\mathcal{X} \times \mathcal{Y}} (2\kappa|f_{u_\lambda}(x) - y|)^2 dP(x, y) \\ &= 4\kappa^2 R(f_{u_\lambda}). \end{aligned} \quad (3.16)$$

(ii): Let  $\widehat{F}_n : \ell^2(\mathbb{K}) \rightarrow \mathbb{R}_+ : u \mapsto (1/n) \sum_{i=1}^n |f_u(x_i) - y_i|^2$ . Since the restriction of  $G$  to  $\ell^r(\mathbb{K})$  is in  $\Gamma_0(\ell^r(\mathbb{K}))$  by Lemma 7, equation (2.6) and Ekeland's variational principle [26, Corollary 4.2.12] imply that there exists  $\widehat{v} \in \ell^r(\mathbb{K})$  and  $\widehat{e}^* \in \partial(\widehat{F}_n + \lambda G)(\widehat{v})$  such that  $\|\widehat{u}_{n,\lambda}(z_n) - \widehat{v}\|_r \leq \sqrt{\varepsilon_n}$  and  $\|\widehat{e}^*\|_{r^*} \leq \sqrt{\varepsilon_n}$ . Using the inequality  $a^2 - b^2 \geq 2(a - b)b$ , we derive from definitions (3.12) and (2.4) that, for every  $i \in \{1, \dots, n\}$ ,

$$\begin{aligned} \sum_{k \in \mathbb{K}} \langle \widehat{v} - u_\lambda | e_k \rangle \langle \Psi_\lambda(x_i, y_i) | e_k \rangle &= \sum_{k \in \mathbb{K}} \langle \widehat{v} - u_\lambda | e_k \rangle 2(f_{u_\lambda}(x_i) - y_i) \phi_k(x_i) \\ &= 2(f_{\widehat{v}}(x_i) - f_{u_\lambda}(x_i))(f_{u_\lambda}(x_i) - y_i) \\ &\leq (y_i - f_{\widehat{v}}(x_i))^2 - (y_i - f_{u_\lambda}(x_i))^2 \end{aligned} \quad (3.17)$$

and, summing over  $i$  and dividing by  $n$ , we obtain

$$\widehat{F}_n(\widehat{v}) - \widehat{F}_n(u_\lambda) \geq \sum_{k \in \mathbb{K}} \langle \widehat{v} - u_\lambda | e_k \rangle \left\langle \frac{1}{n} \sum_{i=1}^n \Psi_\lambda(x_i, y_i) \middle| e_k \right\rangle. \quad (3.18)$$

Lemma 1 and (i) yield

$$\lambda G(\widehat{v}) - \lambda G(u_\lambda) \geq \langle \widehat{v} - u_\lambda | -\mathbb{E}_P(\Psi_\lambda) \rangle + \lambda \eta M \frac{\|\widehat{v} - u_\lambda\|_r^2}{(\|u_\lambda\|_r + \|\widehat{v} - u_\lambda\|_r)^{2-r}}. \quad (3.19)$$

Next, since  $\widehat{e}^* \in \partial(\widehat{F}_n + \lambda G)(\widehat{v})$ , we have  $\langle u_\lambda - \widehat{v} | \widehat{e}^* \rangle \leq (\widehat{F}_n + \lambda G)(u_\lambda) - (\widehat{F}_n + \lambda G)(\widehat{v})$ . Summing inequalities (3.18) and (3.19), we have

$$\begin{aligned} \sqrt{\varepsilon_n} \|\widehat{v} - u_\lambda\|_r &\geq (\widehat{F}_n + \lambda G)(\widehat{v}) - (\widehat{F}_n + \lambda G)(u_\lambda) \\ &\geq \sum_{k \in \mathbb{K}} \langle \widehat{v} - u_\lambda | e_k \rangle \left\langle \frac{1}{n} \sum_{i=1}^n \Psi_\lambda(x_i, y_i) - \mathbb{E}_P(\Psi_\lambda) \middle| e_k \right\rangle \\ &\quad + \frac{\lambda \eta M \|\widehat{v} - u_\lambda\|_r^2}{(\|u_\lambda\|_r + \|\widehat{v} - u_\lambda\|_r)^{2-r}}. \end{aligned} \quad (3.20)$$

Hence, using Hölder's inequality,

$$\frac{\lambda \eta M \|\widehat{v} - u_\lambda\|_r^2}{(\|u_\lambda\|_r + \|\widehat{v} - u_\lambda\|_r)^{2-r}} \leq \|\widehat{v} - u_\lambda\|_r \left( \left\| \frac{1}{n} \sum_{i=1}^n \Psi_\lambda(x_i, y_i) - \mathbb{E}_P(\Psi_\lambda) \right\|_{r^*} + \sqrt{\varepsilon_n} \right) \quad (3.21)$$

and the statement follows from the fact that  $\|\cdot\|_{r^*} \leq \|\cdot\|_2$ .

We recall the following concentration inequality in Hilbert spaces [36] and give the proof of the main result of this section.

**Lemma 2 (Bernstein's inequality)** *Let  $(U_i)_{1 \leq i \leq n}$  be a finite sequence of i.i.d. random variables on a probability space  $(\Omega, \mathfrak{A}, \mathbb{P})$  and taking values in a real separable Hilbert space  $\mathcal{H}$ . Let  $\beta > 0$ , let  $\sigma > 0$  and suppose that  $\max_{1 \leq i \leq n} \|U_i\| \leq \beta$  and that  $\mathbb{E}_P \|U_i\|^2 \leq \sigma^2$ . Then for every  $\tau > 0$  and every integer  $n \geq 1$*

$$P \left[ \left\| \frac{1}{n} \sum_{i=1}^n (U_i - \mathbb{E}_P U_i) \right\| \geq \frac{2\sigma}{\sqrt{n}} + 4\sigma \sqrt{\frac{\tau}{n} + \frac{4\beta\tau}{3n}} \right] \leq e^{-\tau}. \quad (3.22)$$

*Proof (of Theorem 2)* (i): Let  $n \in \mathbb{N}^*$ , let  $z_n = (x_i, y_i)_{1 \leq i \leq n} \in (\mathcal{X} \times \mathcal{Y})^n$  and let  $\widehat{F}_n: u \in \ell^2(\mathbb{K}) \rightarrow \mathbb{R}_+ : u \mapsto (1/n) \sum_{i=1}^n |f_u(x_i) - y_i|^2$ . Let  $u_G \in \text{Argmin} G$ , let  $\lambda \in \mathbb{R}_{++}$ , and let  $\rho_\lambda = \max \{1, \|u_G\|_r, (2(b + \kappa \|u_G\| + 1)^2 / (\eta M \lambda))^{1/r}\}$ . Since  $F(u_G) \leq (b + \kappa \|u_G\|)^2$  and  $\widehat{F}_n(u_G) \leq (b + \kappa \|u_G\|)^2$ , from the definition of  $\rho_\lambda$  and Proposition 3(ii) we derive that  $\|u_\lambda - u_G\|_r \leq \rho_\lambda$  and  $\|\widehat{u}_{n,\lambda}(z_n) - u_G\|_r \leq \rho_\lambda$ . It follows from Theorem 6 that there exist  $\Psi_\lambda: \mathcal{X} \times \mathcal{Y} \rightarrow \ell^2(\mathbb{K})$  and  $\widehat{v} \in \ell^r(\mathbb{K})$  such that  $\|\widehat{v} - \widehat{u}_{n,\lambda}(z_n)\| \leq \sqrt{\varepsilon_n}$  and

$$\frac{M\eta \|\widehat{v} - u_\lambda\|_r}{(\|u_\lambda\|_r + \|\widehat{v} - u_\lambda\|_r)^{2-r}} \leq \frac{1}{\lambda} \left( \left\| \frac{1}{n} \sum_{i=1}^n \Psi_\lambda(x_i, y_i) - \mathbb{E}_P(\Psi_\lambda) \right\|_2 + \sqrt{\varepsilon_n} \right). \quad (3.23)$$

Therefore, since  $\|u_\lambda\|_r \leq 2\rho_\lambda$  and  $\|\widehat{v} - u_\lambda\|_r \leq \|\widehat{v} - \widehat{u}_{n,\lambda}(z_n)\|_r + \|\widehat{u}_{n,\lambda}(z_n) - u_\lambda\|_r \leq \sqrt{\varepsilon_n} + 2\rho_\lambda \leq 3\rho_\lambda$ , we have

$$\|\widehat{v} - u_\lambda\|_r \leq \frac{(5\rho_\lambda)^{2-r}}{M\eta\lambda} \left( \left\| \mathbb{E}_P(\Psi_\lambda) - \frac{1}{n} \sum_{i=1}^n \Psi_\lambda(x_i, y_i) \right\|_2 + \sqrt{\varepsilon_n} \right). \quad (3.24)$$



Thus,

$$\|\widehat{u}_{n,\lambda}(z_n) - u_\lambda\|_r \leq \sqrt{\varepsilon_n} + \frac{(5\rho_\lambda)^{2-r}}{M\eta\lambda} \left( \|\mathbb{E}_P(\Psi_\lambda) - \frac{1}{n} \sum_{i=1}^n \Psi_\lambda(x_i, y_i)\|_2 + \sqrt{\varepsilon_n} \right). \quad (3.25)$$

Now, consider the i.i.d. random vectors  $\Psi_\lambda(X_i, Y_i): \Omega \rightarrow \ell^2(\mathbb{K})$ , for  $1 \leq i \leq n$ . It follows from Theorem 6(i) that  $\max_{1 \leq i \leq n} \|\Psi_\lambda(X_i, Y_i)\| \leq 2\kappa(\kappa\rho_\lambda + b)$  and that  $\max_{1 \leq i \leq n} \mathbb{E}_P \|\Psi_\lambda(X_i, Y_i)\|^2 \leq 4\kappa^2 R(f_{u_\lambda})$ . Now set  $\beta_\lambda = 2\kappa(\kappa\rho_\lambda + b)$  and  $\sigma_\lambda^2 = \kappa^2 R(f_{u_\lambda})$ . Then Bernstein's inequality in Hilbert spaces (Lemma 2) gives

$$(\forall \tau \in \mathbb{R}_{++}) \quad \mathbb{P} \left[ \left\| \mathbb{E}(\Psi_\lambda(X, Y)) - \frac{1}{n} \sum_{i=1}^n \Psi_\lambda(X_i, Y_i) \right\|_2 \leq \delta(n, \lambda, \tau) \right] \geq 1 - e^{-\tau}, \quad (3.26)$$

where  $\delta(n, \lambda, \tau) = 2\sigma_\lambda/\sqrt{n} + 4\sigma_\lambda\sqrt{\tau/n} + 4\beta_\lambda\tau/(3n)$ . Thus, recalling (3.25) we have

$$\mathbb{P} \left[ \|\widehat{u}_{n,\lambda}(Z_n) - u_\lambda\|_r > \sqrt{\varepsilon_n} + \frac{(5\rho_\lambda)^{2-r}}{M\eta\lambda} (\delta(n, \lambda, \tau) + \sqrt{\varepsilon_n}) \right] \leq e^{-\tau}. \quad (3.27)$$

Set  $\gamma_0 = 2(b + \kappa\|u_G\| + 1)^2$  and  $\gamma_1 = 5^{2-r}\gamma_0^{2/r-1}/(\eta M)^{2/r}$ . We note that, since  $\sigma_\lambda$  is bounded, say by  $\gamma_2$ , for  $\lambda < 1$  sufficiently small, we have

$$\begin{aligned} & \frac{(5\rho_\lambda)^{2-r}}{M\eta\lambda} (\delta(n, \lambda, \tau) + \sqrt{\varepsilon_n}) \\ &= \frac{5^{2-r}}{M\eta} \left( \frac{\gamma_0}{\eta M} \right)^{\frac{2}{r}-1} \frac{1}{\lambda^{2/r}} \left( \frac{2\sigma_\lambda}{\sqrt{n}} + 4\sigma_\lambda \sqrt{\frac{\tau}{n}} + \frac{4\beta_\lambda\tau}{3n} + \sqrt{\varepsilon_n} \right) \\ &\leq \gamma_1 \left( \frac{2\gamma_2}{\lambda^{2/r}n^{1/2}} + 4\gamma_2 \frac{\sqrt{\tau}}{\lambda^{2/r}n^{1/2}} + \frac{8\tau\kappa^2\gamma_0/(\eta M)^{1/r}}{3n\lambda^{3/r}} + \frac{8\tau\kappa b}{3n\lambda^{2/r}} + \frac{\sqrt{\varepsilon_n}}{\lambda^{2/r}} \right). \end{aligned} \quad (3.28)$$

Therefore, since  $1/(\lambda_n^{2/r}n^{1/2}) \rightarrow 0$  and  $\sqrt{\varepsilon_n}/\lambda_n^{2/r} \rightarrow 0$  it follows that

$$\frac{(5\rho_{\lambda_n})^{2-r}}{\lambda_n} (\delta(n, \lambda_n, \tau) + \sqrt{\varepsilon_n}) \rightarrow 0 \quad (3.29)$$

and hence, in view of (3.27), we get  $\|\widehat{u}_{n,\lambda_n}(Z_n) - u_{\lambda_n}\|_r \rightarrow 0$  in probability. Now recalling (3.1), since  $F(u_{\lambda_n}) - \inf F(\text{dom } G) \rightarrow 0$  by Proposition 3(iii), and  $\|\widehat{u}_{n,\lambda}(Z_n) - u_{\lambda_n}\|_r \rightarrow 0$  in probability, we derive that  $\|\widehat{f}_n - f_{\mathcal{E}}\|_{L^2} \rightarrow 0$  in probability.

(ii): Let  $n \in \mathbb{N}^*$ , let  $\eta \in \mathbb{R}_{++}$ , and set

$$\Omega_{n,\eta} = \left\{ \|\widehat{f}_n - f_{\mathcal{E}}\|_{L^2} > \|A\|\eta + \sqrt{F(u_{\lambda_n}) - \inf F(\text{dom } G)} \right\}. \quad (3.30)$$

Since  $\varepsilon_n = O(1/n)$ , it follows from (3.28) that there exists  $\gamma_3 \in \mathbb{R}_{++}$  such that, for every  $\tau \in ]1, +\infty[$ , and every  $n \in \mathbb{N}^*$ ,

$$\frac{(5\rho_{\lambda_n})^{2-r}}{\eta M \lambda_n} (\delta(n, \lambda_n, \tau) + \sqrt{\varepsilon_n}) \leq \frac{\gamma_3 \tau}{\lambda_n^{2/r} n^{1/2}}. \quad (3.31)$$

Let  $\xi \in ]1, +\infty[$ . There exists  $\bar{n} \in \mathbb{N}^*$ , such that, for every integer  $n \geq \bar{n}$ ,

$$\frac{\gamma_3}{\lambda_n^{2/r} n^{1/2}} \leq \frac{\gamma_3 \xi \log n}{\lambda_n^{2/r} n^{1/2}} \leq \eta. \quad (3.32)$$

Therefore, it follows from (3.27), (3.1), (3.31), and (3.32) that, for  $n$  large enough,

$$\mathbb{P}\Omega_{n,\eta} \leq \exp\left(-\frac{\eta \lambda_n^{2/r} n^{1/2}}{\gamma_3}\right) \leq \exp(-\xi \log n) = n^{-\xi}. \quad (3.33)$$

Thus,  $\sum_{n=\bar{n}}^{+\infty} \mathbb{P}\Omega_{n,\eta} < +\infty$  and we derive from the Borel-Cantelli lemma that  $\mathbb{P}\left(\bigcap_{k \geq \bar{n}} \bigcup_{n \geq k} \Omega_{n,\eta}\right) = 0$ . Recalling Proposition 3(iii), we conclude that the sequence  $\|\widehat{f}_n - f_{\mathcal{C}}\|_{L^2} \rightarrow 0$  P-a.s.

(iii): First note that Proposition 3(iii) implies that  $u^\dagger$  is well defined and that  $\rho = \sup_{\lambda \in \mathbb{R}_{++}} \|u_\lambda\| < +\infty$ . Now, let  $\lambda \in \mathbb{R}_{++}$  and let  $n \in \mathbb{N}^*$ . Since  $\|u_\lambda\| \leq \rho$ , arguing as in the proof of (i), we obtain

$$(\forall \tau \in \mathbb{R}_{++}) \mathbb{P}\left[\|\widehat{u}_{n,\lambda}(Z_n) - u_\lambda\|_r > \sqrt{\varepsilon_n} + \frac{(5\rho)^{2-r}}{M\lambda}(\delta(n,\tau) + \sqrt{\varepsilon_n})\right] \leq e^{-\tau}, \quad (3.34)$$

where  $\sigma = 2\kappa(\kappa\rho + b)$  and  $\delta(n,\tau) = 4\sigma/\sqrt{n} + 4\sigma\sqrt{\tau/n} + 4\sigma\tau/(3n)$ .

(iii)(a): Since  $1/(\lambda_n n^{1/2}) \rightarrow 0$ , we have  $(1/\lambda_n)\delta(n,\tau) \rightarrow 0$  and hence in view of (3.34),  $\|\widehat{u}_{n,\lambda_n}(Z_n) - u_{\lambda_n}\|_r \rightarrow 0$  in probability. Moreover, since  $\|u_{n,\lambda_n}(Z_n) - u^\dagger\| \leq \|u_{n,\lambda_n}(Z_n) - u_{\lambda_n}\| + \|u_{\lambda_n} - u^\dagger\|$ , the statement follows by Proposition 3(iv).

(iii)(b): The proof follows the same line as that of (ii).

#### 4 Algorithmic analysis

The goal of this section is to prove Theorem 3 and Theorem 5. The proof of Theorem 3 is based on the following fact.

**Lemma 3** [34, Lemma 4.1] *Let  $\mathcal{H}$  be a real Hilbert space, let  $\beta \in \mathbb{R}_{++}$ , and let  $\delta \in \mathbb{R}_+$ . Let  $J_1: \mathcal{H} \rightarrow \mathbb{R}$  be a convex differentiable function with a  $\beta$ -Lipschitz continuous gradient, and let  $J_2 \in \Gamma_0(\mathcal{H})$ . Then, for every  $(u, v, w) \in \mathcal{H}^3$  and every  $v^* \in \partial_\delta J_2(v)$ ,*

$$(J_1 + J_2)(v) \leq (J_1 + J_2)(u) + \langle v - u | \nabla J_1(w) + v^* \rangle + \frac{\beta}{2} \|v - w\|^2 + \delta. \quad (4.1)$$

*Proof (of Theorem 3)* Let  $m \in \mathbb{N}$  and set

$$\check{v}_m = \text{prox}_{\gamma_m J_2}(u_m - \gamma_m(\nabla J_1(u_m) + b_m)). \quad (4.2)$$

Since

$$v_m \in \text{Argmin}_{w \in \mathcal{H}}^{\delta_m^2/(2\gamma_m)} \left\{ J_2(w) + \frac{1}{2\gamma_m} \|w - (u_m - \gamma_m(\nabla J_1(u_m) + b_m))\|^2 \right\}, \quad (4.3)$$

using the strong convexity of the objective function in (4.3), we get

$$\|v_m - \tilde{v}_m\| \leq \delta_m. \quad (4.4)$$

Therefore, setting  $a_m = v_m - \tilde{v}_m$ , we have

$$u_{m+1} = u_m + \tau_m \left( \text{prox}_{\gamma_m J_2}(u_m - \gamma_m (\nabla J_1(u_m) + b_m)) + a_m - u_m \right). \quad (4.5)$$

Hence (2.15) is an instance of the inexact forward-backward algorithm studied in [12] and we can therefore use the results of [12, Theorem 3.4].

(i)–(ii): The statements follow from [12, Theorem 3.4(i)–(ii)].

(iii): We have

$$\begin{aligned} \|u_m - v_m\|^2 &\leq 2\|u_m - \tilde{v}_m\|^2 + 2\|a_m\|^2 \\ &\leq 4\|u_m - \text{prox}_{\gamma_m J_2}(u_m - \gamma_m \nabla J_1(u_m))\|^2 + 4\|b_m\|^2 + 2\|a_m\|^2. \end{aligned} \quad (4.6)$$

Therefore, the statement follows from [12, Theorem 3.4(iii)].

(iv): By (4.3) and [31, Lemma 1], there exist  $\delta_{1,m} \in [0, +\infty[$ ,  $\delta_{2,m} \in [0, +\infty[$ , and  $e_m \in \mathcal{H}$  with  $\delta_{1,m}^2 + \delta_{2,m}^2 \leq \delta_m^2$  and  $\|e_m\| \leq \delta_{2,m}$  such that

$$v_m^* = \frac{u_m - v_m}{\gamma_m} - (\nabla J_1(u_m) + b_m) + \frac{e_m}{\gamma_m} \in \partial_{\delta_{1,m}/(2\gamma_m)} J_2(v_m). \quad (4.7)$$

Since  $J = J_1 + J_2$ , it follows from Lemma 3 that, for every  $u \in \mathcal{H}$ ,

$$\begin{aligned} J(v_m) - J(u) &\leq \langle v_m - u \mid \nabla F(u_m) + v_m^* \rangle + \frac{L}{2} \|v_m - u_m\|^2 + \frac{\delta_{1,m}^2}{2\gamma_m} \\ &= \frac{1}{\gamma_m} \langle v_m - u \mid u_m - v_m \rangle + \frac{L}{2} \|v_m - u_m\|^2 + \frac{1}{\gamma_m} \langle v_m - u \mid e_m - \gamma_m b_m \rangle + \frac{\delta_{1,m}^2}{2\gamma_m} \\ &= \frac{1}{2\gamma_m} (\|u_m - u\|^2 - \|v_m - u\|^2) + \frac{1}{2} \left( \beta - \frac{1}{\gamma_m} \right) \|v_m - u_m\|^2 \\ &\quad + \frac{1}{\gamma_m} \langle v_m - u \mid e_m - \gamma_m b_m \rangle + \frac{\delta_{1,m}^2}{2\gamma_m}. \end{aligned} \quad (4.8)$$

We derive from (i) and (iii) that  $(\langle v_m - u \mid u_m - v_m \rangle)_{m \in \mathbb{N}}$  is square summable. Therefore, if we let  $u \in \text{Argmin} J$ , it follows from (4.8) that  $(J(v_m) - \inf J(\mathcal{H}))_{m \in \mathbb{N}}$  is square summable. Now, if we let  $u = u_m$  in (4.8) we have

$$\begin{aligned} J(v_m) - J(u_m) &\leq \left( \frac{\beta}{2} - \frac{1}{\gamma_m} \right) \|u_m - v_m\|^2 + \frac{1}{\gamma_m} \left( \|u_m - v_m\| \|e_m - \gamma_m b_m\| + \frac{\delta_{1,m}^2}{2} \right) \\ &\leq \frac{1}{\gamma_m} \left( \|u_m - v_m\| \|e_m - \gamma_m b_m\| + \frac{\delta_{1,m}^2}{2} \right). \end{aligned} \quad (4.9)$$

Set  $\underline{\gamma} = \inf_{m \in \mathbb{N}} \gamma_m$ . Since  $u_{m+1} = u_m + \tau_m(v_m - u_m)$ , using the convexity of  $J$  and (4.9), we get

$$\begin{aligned} J(u_{m+1}) - \inf J(\mathcal{H}) &\leq J(u_m) - \inf J(\mathcal{H}) + \tau_m(J(v_m) - J(u_m)) \\ &\leq J(u_m) - \inf J(\mathcal{H}) + \underline{\gamma}^{-1}(\|u_m - v_m\| \|e_m - \gamma_m b_m\| + \delta_{1,m}^2/2). \end{aligned} \quad (4.10)$$

Thus, since  $(\|u_m - v_m\| \|e_m - \gamma_m b_m\| + \delta_{1,m}^2/2)_{m \in \mathbb{N}}$  is summable, [29, Lemma 2.2.2], ensures that  $(J(u_m) - \inf J(\mathcal{H}))_{m \in \mathbb{N}}$  converges and, in view of the inequalities in (4.10), its limit must be 0.

(v): Let  $u \in \text{Argmin} J$ . Since,  $u_m - u = (1 - \tau_{m-1})(u_{m-1} - u) + \tau_{m-1}(v_{m-1} - u)$ , it follows from the convexity of  $\|\cdot\|^2$  that

$$\|u_m - u\| - \|v_m - u\|^2 \leq (1 - \tau_{m-1})\|u_{m-1} - u\|^2 + \|v_{m-1} - u\|^2 - \|v_m - u\|^2. \quad (4.11)$$

Therefore, it follows from (4.8) that

$$\begin{aligned} 0 &\leq J(v_m) - J(u) \\ &\leq \frac{1 - \tau_{m-1}}{2\underline{\gamma}} \|u_{m-1} - u\|^2 + \frac{1}{2\underline{\gamma}} (\|v_{m-1} - u\|^2 - \|v_m - u\|^2) \\ &\quad + \frac{1}{2} \left( \beta - \frac{1}{\gamma_m} \right) \|v_m - u_m\|^2 + \frac{1}{\underline{\gamma}} \left( \|v_m - u\| \|e_m - \gamma_m b_m\| + \frac{\delta_{1,m}^2}{2} \right). \end{aligned} \quad (4.12)$$

Hence,  $(J(v_m) - \inf J(\mathcal{H}))_{m \in \mathbb{N}}$  is summable, for each term on the right hand side of (4.12) is summable. Since  $u_{m+1} = (1 - \tau_m)u_m + \tau_m v_m$ , convexity of  $J$  yields

$$0 \leq J(u_{m+1}) - \inf J(\mathcal{H}) \leq (1 - \tau_m)(J(u_m) - \inf J(\mathcal{H})) + \tau_m(J(v_m) - \inf J(\mathcal{H})). \quad (4.13)$$

The summability of  $(1 - \tau_m)_{m \in \mathbb{N}}$  and  $(J(v_m) - \inf J(\mathcal{H}))_{m \in \mathbb{N}}$  implies that of  $(J(u_m) - \inf J(\mathcal{H}))_{m \in \mathbb{N}}$ .

(vi): Since  $(J(u_m) - \inf J(\mathcal{H}))_{m \in \mathbb{N}}$  is summable, it follows from (4.10) and [15, Lemma 3] that  $(J(u_m) - \inf J(\mathcal{H})) = o(1/m)$ .

The purpose of the rest of the section is to prove Proposition 2 and Theorem 5.

**Lemma 4** *Let  $h: \mathbb{R} \rightarrow \mathbb{R}$  be convex and such that  $0 \in \text{Argmin}_{\mathbb{R}} h$ , let  $(s, \mu) \in \mathbb{R}^2$ , and let  $\alpha \in [-1, 1]$ . Let  $\beta \in \mathbb{R}_+$  be the Lipschitz constant of  $h$  in  $[-1 - |\mu|, 1 + |\mu|]$  and set*

$$\delta = \sqrt{(2\beta + 2|\mu| + 1)|\alpha|} \quad \text{and} \quad s = \text{prox}_h \mu + \alpha. \quad (4.14)$$

*Then  $s \simeq_{\delta} \text{prox}_h \mu$ . Moreover,  $\hat{s} = \text{sign}(\mu) \max\{0, \text{sign}(\mu)s\}$  satisfies  $\hat{s} \simeq_{\delta} \text{prox}_h \mu$  and  $\mu \hat{s} \geq 0$ .*

*Proof* Let  $t = \text{prox}_h \mu$ . Since  $0 \in \text{Argmin} h$ ,  $\text{prox}_h 0 = 0$ . Hence, since  $\text{prox}_h$  is non-expansive and increasing [12, Lemma 2.4],  $|t| \leq |\mu|$  and  $\text{sign}(t) = \text{sign}(\mu)$ . We note that  $|s| \leq |s-t| + |t| \leq 1 + |\mu|$ . Thus,

$$\begin{aligned} h(s) + \frac{1}{2}|s-\mu|^2 - h(t) - \frac{1}{2}|t-\mu|^2 &\leq \beta|s-t| + \frac{1}{2}|s-t||s-\mu+t-\mu| \\ &\leq \frac{1}{2}(2\beta + 1 + 2|\mu|)|\alpha|. \end{aligned} \quad (4.15)$$

To conclude, it is enough to note that  $|\hat{s} - \text{prox}_h(\mu)| \leq |\alpha|$ .

**Lemma 5** *Let  $h \in \Gamma_0(\mathbb{R})$ , let  $\sigma \in \Gamma_0(\mathbb{R})$  be a support function, and set  $\phi = h + \sigma$ . Let  $(s, x) \in \mathbb{R}^2$  be such that  $s \text{prox}_\sigma(x) \geq 0$ , and let  $\delta \in \mathbb{R}_+$ . Then*

$$s \simeq_\delta \text{prox}_h(\text{prox}_\sigma x) \quad \Rightarrow \quad s \simeq_\delta \text{prox}_\phi x. \quad (4.16)$$

*Proof* Let  $\mu = \text{prox}_\sigma x$  and  $s \simeq_\delta \text{prox}_h(\text{prox}_\sigma x)$ . By [31, Lemma 2.4] there exist  $(\delta_1, \delta_2) \in \mathbb{R}_+^2$  and  $e \in \mathbb{R}$ , such that

$$\mu - s + e \in \partial_{\delta_1^2/2} h(s), \quad |e| \leq \delta_2, \quad \text{and} \quad \delta_1^2 + \delta_2^2 \leq \delta^2. \quad (4.17)$$

Hence

$$x - s + e = x - \mu + \mu - s + e \in \partial\sigma(\mu) + \partial_{\delta_1^2/2} h(s). \quad (4.18)$$

Since  $s\mu \geq 0$ , there exists  $t \in \mathbb{R}_+$  such that  $\mu = ts$ . Moreover, since  $\sigma$  is positively homogeneous,  $\partial\sigma(ts) \subset \partial\sigma(s)$ . Therefore  $x - s + e \in \partial\sigma(s) + \partial_{\delta_1^2/2} h(s) \subset \partial_{\delta_1^2/2} \phi(s)$ , which implies that  $s \simeq_\delta \text{prox}_\phi x$  by [31, Lemma 2.4].

*Remark 8* Let  $h \in \Gamma_0(\mathbb{R})$ , let  $(s, \mu) \in \mathbb{R}^2$ , and let  $\delta \in \mathbb{R}_{++}$ . Suppose that  $0 \in \text{Argmin}_{\mathbb{R}} h$  and that  $s \simeq_\delta \text{prox}_h \mu$  with  $\delta \leq |s|$ . Then  $s\mu \geq 0$ . Indeed, since  $h(0) = \inf h(\mathbb{R})$ , we have

$$h(s) + \frac{1}{2}|s-\mu|^2 \leq h(0) + \frac{1}{2}\mu^2 + \frac{1}{2}\delta^2 \leq h(s) + \frac{1}{2}\mu^2 + \frac{1}{2}\delta^2 \quad (4.19)$$

and hence  $0 \leq (1/2)(s^2 - \delta^2) \leq s\mu$ . This shows that Lemma 5, when  $\delta = 0$ , gives  $\text{prox}_\phi = \text{prox}_h \circ \text{prox}_\sigma$  and consequently generalizes [9, Proposition 3.6], relaxing also the condition on the differentiability of  $h$  at 0. With the help of this result one can compute general thresholding operators as the proximity operator of  $|\cdot| + \eta|\cdot|^r$ . Figure 1 depicts some instances of these thresholders (see also [9]).

The following lemma is an error-tolerant version of [10, Proposition 12].

**Lemma 6** *Let  $\phi \in \Gamma_0(\mathbb{R})$ , let  $(s, x, p) \in \mathbb{R}^3$ , let  $\delta \in \mathbb{R}_+$ , and let  $C \subset \mathbb{R}$  be a nonempty closed interval. Then*

$$s \simeq_\delta \text{prox}_\phi x, \quad \text{and} \quad p = \text{proj}_C s \quad \Rightarrow \quad p \simeq_\delta \text{prox}_{\phi + \iota_C} x. \quad (4.20)$$

*Proof* Let  $g = \phi + (1/2)|\cdot - x|^2$  and let  $\varepsilon = (\delta^2/2)$ . Since  $g$  is convex and  $\bar{s} = \text{prox}_\phi x$  is its minimum,  $g$  is decreasing on  $]-\infty, \bar{s}]$  and increasing on  $[\bar{s}, +\infty[$ . By definition  $s$  is an  $\varepsilon$ -minimizer of  $g$ . The statement is equivalent to the fact that  $p$  is an  $\varepsilon$ -minimizer of  $g + \iota_C$ . If  $s \in C$ , then  $p$  is a fortiori an  $\varepsilon$ -minimizer of  $g + \iota_C$ . We now consider two cases. First suppose that  $s < \inf C$ . If  $s < \inf C \leq \bar{s}$ , then  $\inf C$  is still an  $\varepsilon$ -minimizer of  $g$  and  $\inf C \in C$ . Thus  $p = \inf C$  is an  $\varepsilon$ -minimizer of  $g + \iota_C$ . If either  $s \leq \bar{s} \leq \inf C$  or  $\bar{s} \leq s < \inf C$ , we have  $p = \text{proj}_C \bar{s} = \inf C$ , which is the minimum of  $g + \iota_C$ , since  $g$  is increasing on  $[\bar{s}, +\infty[$ . The second case  $\sup C < s$  is treated likewise.

*Proof (of Proposition 2)* Set

$$(\forall k \in \mathbb{K}) \quad \begin{cases} \mu_k = \text{soft}_{\gamma D_k} \chi_k \\ s_k = \text{sign}(\mu_k) \max \{0, \text{sign}(\mu_k)(\text{prox}_{\gamma h_k} \mu_k + \alpha_k)\} \\ v_k = \text{proj}_{C_k} s_k. \end{cases} \quad (4.21)$$

Let  $k \in \mathbb{K}$ . Since  $\text{soft}_{\gamma D_k}$  is nonexpansive and  $2\gamma \max\{h_k(|\chi_k| + 2), h_k(-|\chi_k| - 2)\}$  is a Lipschitz constant for  $\gamma h_k$  on the interval  $[-|\chi_k| - 1, |\chi_k| + 1]$ , it follows from (4.21) and Lemma 4 that

$$\begin{cases} \delta_k^2 = (4\gamma \max\{h_k(|\chi_k| + 2), h_k(-|\chi_k| - 2)\} + 2|\chi_k| + 1)|\alpha_k| \\ s_k \simeq_{\delta_k} \text{prox}_{\gamma h_k}(\text{prox}_{\gamma \sigma_{D_k}} \chi_k) \\ s_k \text{prox}_{\gamma \sigma_{D_k}} \chi_k \geq 0. \end{cases} \quad (4.22)$$

Thus, Lemma 5 yields

$$s_k \simeq_{\delta_k} \text{prox}_{\gamma(h_k + \sigma_{D_k})} \chi_k, \quad (4.23)$$

and, using Lemma 6, we obtain  $v_k \simeq_{\delta_k} \text{prox}_{\gamma g_k} \chi_k$ . Hence, by Definition 1,

$$\gamma g_k(v_k) + \frac{1}{2}|v_k - \chi_k|^2 \leq \gamma g_k(\text{prox}_{\gamma g_k} \chi_k) + \frac{1}{2}|\text{prox}_{\gamma g_k} \chi_k - \chi_k|^2 + \frac{\delta_k^2}{2}. \quad (4.24)$$

On the other hand, we derive from [12, Example 2.19] and [9, Proposition 3.6] that

$$\langle \text{prox}_{\gamma G} w \mid e_k \rangle = \text{prox}_{\gamma g_k} \chi_k. \quad (4.25)$$

Thus, summing the inequalities (4.24) over  $k$ , we obtain

$$\begin{aligned} & \gamma \sum_{k \in \mathbb{K}} g_k(v_k) + \frac{1}{2} \sum_{k \in \mathbb{K}} |v_k - \chi_k|^2 \\ & \leq \gamma \sum_{k \in \mathbb{K}} \left( g_k(\langle \text{prox}_{\gamma g} w \mid e_k \rangle) + \frac{1}{2} |\langle \text{prox}_{\gamma G} w \mid e_k \rangle - \langle w \mid e_k \rangle|^2 \right) + \frac{1}{2} \sum_{k \in \mathbb{K}} \delta_k^2 \\ & \leq \gamma G(\text{prox}_{\gamma G} w) + \frac{1}{2} \|\text{prox}_{\gamma G} w - w\|^2 + \frac{1}{2} \sum_{k \in \mathbb{K}} \xi_k \\ & < +\infty. \end{aligned} \quad (4.26)$$

Thus, since, by Lemma 7,  $\sum_{k \in \mathbb{K}} g_k(v_k)$  is either convergent or divergent to  $+\infty$ , (4.26) yields  $(v_k)_{k \in \mathbb{N}} \in \ell^2(\mathbb{K})$  and one can find  $v \in \mathcal{H}$  such that, for every  $k \in \mathbb{K}$ ,  $\langle v | e_k \rangle = v_k$ . Hence,

$$\gamma G(v) + \frac{1}{2} \|v - w\|^2 \leq \gamma G(\text{prox}_{\gamma G} w) + \frac{1}{2} \|\text{prox}_{\gamma G} w - w\|^2 + \frac{1}{2} \sum_{k \in \mathbb{K}} \xi_k \quad (4.27)$$

and finally  $v \simeq_{\delta} \text{prox}_{\gamma G} w$ , where  $\delta = \sqrt{\sum_{k \in \mathbb{K}} \xi_k}$ .

*Proof (of Theorem 5)* (i): Lemma 7 guarantees that  $G \in \Gamma_0(\ell^2(\mathbb{K}))$ , that  $G$  is coercive, and that  $\text{dom} G \subset \ell^r(\mathbb{K})$ . The statement therefore follows from [3, Corollary 11.16(ii)].

(ii)–(iv): Let  $\widehat{F}_n: \ell^2(\mathbb{K}) \rightarrow \mathbb{R}: u \rightarrow (1/n) \sum_{i=1}^n (\langle u | \Phi(x_i) \rangle - y_i)^2$ . Then, for every  $u \in \ell^2(\mathbb{K})$ ,  $\nabla \widehat{F}_n(u) = (2/n) \sum_{i=1}^n (\langle u | \Phi(x_i) \rangle - y_i) \Phi(x_i)$ . Hence, since  $\|\Phi(x_i)\|_2 \leq \kappa$ ,  $\nabla(1/\lambda) \widehat{F}_n$  is Lipschitz continuous with constant  $2\kappa^2/\lambda$ . Therefore, the statement follows from Theorem 3, with  $J_1 = (1/\lambda) \widehat{F}_n$  and  $J_2 = G$ , by noting that, in view of Proposition 2, for every  $m \in \mathbb{N}$ ,  $v_m = (v_{m,k})_{k \in \mathbb{K}} \in \ell^2(\mathbb{K})$  and  $v_m \simeq_{\delta_m} \text{prox}_{\gamma_m G} w_m$ , where  $\delta_m \leq \zeta \sqrt{\sum_{k \in \mathbb{K}} \xi_k / m^p}$ . It remains to show the convergence properties of  $(\|u_m - \widehat{u}\|_r)_{m \in \mathbb{N}}$  and  $(\|v_m - \widehat{u}\|_r)_{m \in \mathbb{N}}$ . We focus on the sequence  $(\|u_m - \widehat{u}\|_r)_{m \in \mathbb{N}}$ , since  $(\|v_m - \widehat{u}\|_r)_{m \in \mathbb{N}}$  can be treated analogously. It follows from Lemma 1 and the convexity of  $\widehat{F}_n$  that

$$(\forall m \in \mathbb{N}) \quad ((1/\lambda) \widehat{F}_n + G)(u_m) - ((1/\lambda) \widehat{F}_n + G)(\widehat{u}) \geq \frac{\eta M \|u_m - \widehat{u}\|_r^2}{(\|\widehat{u}\|_r + \|u_m - \widehat{u}\|_r)^{2-r}}. \quad (4.28)$$

Therefore, since  $((1/\lambda) \widehat{F}_n + G)(u_m) - ((1/\lambda) \widehat{F}_n + \lambda G)(\widehat{u}) \rightarrow 0$  as  $m \rightarrow +\infty$  and  $\psi: \mathbb{R}_+ \rightarrow \mathbb{R}: t \mapsto t^2 / (\|\widehat{u}\|_r + t)^{2-r}$  is strictly increasing with  $\psi(0) = 0$ , we obtain  $\|u_m - \widehat{u}\|_r \rightarrow 0$ . Moreover, taking  $\rho \in \mathbb{R}_{++}$  such that  $\sup_{m \in \mathbb{N}} (\|\widehat{u}\|_r + \|u_m - \widehat{u}\|_r)^{2-r} \leq \rho$ , (2.22) follows from (4.28).

## References

1. A. Antoniadis, D. Leporini, and J.-C. Pesquet. Wavelet thresholding for some classes of non-Gaussian noise. *Statistica Neerlandica*, 56:434–453, 2002.
2. H. Attouch. Viscosity solutions of minimization problems. *SIAM Journal on Optimization*, 6:769–805, 1996.
3. H. H. Bauschke and P. L. Combettes. *Convex Analysis and Monotone Operator Theory in Hilbert Spaces*, second edition. Springer, New York, 2017.
4. K. Bredies. A forward-backward splitting algorithm for the minimization of non-smooth convex functionals in Banach space. *Inverse Problems*, 25: art. 015005, 2009.
5. P. Bühlmann and S. van de Geer. *Statistics for High-Dimensional Data*. Springer, Heidelberg, 2011.
6. D. Butnariu, A. N. Iusem, and C. Zălinescu. On uniform convexity, total convexity and convergence of the proximal point and outer Bregman projection algorithms in Banach spaces. *Journal of Convex Analysis*, 10:35–61, 2003.
7. C. Chau, P. L. Combettes, J.-C. Pesquet, and V. Wajs. A variational formulation for frame-based inverse problems. *Inverse Problems*, 23:1495–1518, 2007.
8. P. L. Combettes. Strong convergence of block-iterative outer approximation methods for convex optimization. *SIAM Journal on Control and Optimization*, 38:538–565, 2000.

9. P. L. Combettes and J.-C. Pesquet. Proximal thresholding algorithm for minimization over orthonormal bases. *SIAM Journal on Optimization*, 18:1351–1376, 2007.
10. P. L. Combettes and J.-C. Pesquet. A Douglas-Rachford splitting approach to nonsmooth convex variational signal recovery. *IEEE Journal of Selected Topics in Signal Processing*, 1:564–574, 2007.
11. P. L. Combettes, S. Salzo, and S. Villa. Regularized learning schemes in feature Banach spaces. *Analysis and Applications*, published online 2016-12-07.
12. P. L. Combettes and V. R. Wajs. Signal recovery by proximal forward-backward splitting. *Multiscale Modeling and Simulation*, 4:1168–1200, 2005.
13. F. Cucker and S. Smale. On the mathematical foundations of learning. *Bulletin of the American Mathematical Society (N.S.)*, 39:1–49, 2002.
14. I. Daubechies, M. Defrise, and C. De Mol. An iterative thresholding algorithm for linear inverse problems with a sparsity constraint. *Communications on Pure and Applied Mathematics*, 57:1413–1457, 2004.
15. D. Davis and Y. Yin. Convergence rate analysis of several splitting schemes. In: *Splitting Methods in Communication, Imaging, Science, and Engineering*, pp. 115–163. Springer, New York, 2016.
16. C. De Mol, E. De Vito, and L. Rosasco. Elastic-net regularization in learning theory. *Journal of Complexity*, 25:201–230, 2009.
17. C. De Mol, S. Mosci, M. Traskine, and A. Verri. A regularized method for selecting nested groups of relevant genes from microarray data. *Journal of Computational Biology*, 16:677–690, 2009.
18. E. De Vito, L. Rosasco, A. Caponnetto, M. Piana, and A. Verri. Some properties of regularized kernel methods. *Journal of Machine Learning Research*, 5:1363–1390, 2004.
19. E. De Vito, L. Rosasco, A. Caponnetto, U. De Giovannini, and F. Odone. Learning from examples as an inverse problem. *Journal of Machine Learning Research*, 6:883–904, 2005.
20. E. De Vito, V. Umanità, and S. Villa. A consistent algorithm to solve Lasso, elastic-net and Tikhonov regularization. *Journal of Complexity*, 27:188–200, 2011.
21. T. Evgeniou, M. Pontil, and T. Poggio. Regularization networks and support vector machines. *Advances in Computational Mathematics*, 13:1–50, 2000.
22. W. J. Fu. Penalized regressions: the bridge versus the lasso. *Journal of Computational and Graphical Statistics*, 7:397–416, 1998.
23. L. Györfi, M. Kohler, A. Krzyżak, and H. Walk. *A Distribution-Free Theory of Nonparametric Regression*. Springer-Verlag, New York, 2002.
24. A. E. Hoerl and R. W. Kennard. Ridge regression: Biased estimation for nonorthogonal problems. *Technometrics*, 12:55–67, 1970.
25. V. Koltchinskii. Sparsity in penalized empirical risk minimization. *Annales de l’Institut Henri Poincaré Probabilités et Statistiques*, 45:7–57, 2009.
26. R. Lucchetti. *Convexity and Well-Posed Problems*. Springer, New York, 2006.
27. J. J. Moreau. Fonctions convexes duales et points proximaux dans un espace hilbertien. *Comptes Rendus de l’Académie des Sciences de Paris*, A255:2897–2899, 1962.
28. Yu. Nesterov. Gradient methods for minimizing composite functions. *Mathematical Programming, Series B*, 140:125–161, 2013.
29. B. T. Polyak. *Introduction to Optimization*. Optimization Software Inc., New York, 1987.
30. S. Salzo, S. Masecchia, A. Verri, and A. Barla. Alternating proximal regularized dictionary learning. *Neural Computation*, 26:2855–2895, 2014.
31. S. Salzo and S. Villa. Inexact and accelerated proximal point algorithms. *Journal of Convex Analysis*, 19:1167–1192, 2012.
32. M. Schmidt, N. Le Roux, and F. Bach. Convergence rates of inexact proximal-gradient methods for convex optimization. In *Advances in Neural Information Processing Systems*, 24:1458–1466, 2011.
33. B. Schölkopf, R. Herbrich, and A. J. Smola. A generalized representer theorem. In D. Helmbold and B. Williamson, editors, *Computational Learning Theory*, volume 2111 of *Lecture Notes in Computer Science*, pp. 416–426. Springer, Berlin, 2001.
34. S. Villa, S. Salzo, L. Baldassarre, and A. Verri. Accelerated and inexact forward-backward algorithms. *SIAM Journal on Optimization*, 23:1607–1633, 2013.
35. V. Wajs. *Décompositions et Algorithmes Proximaux pour l’Analyse et le Traitement Itératif des Signaux*. 2007. Thèse de doctorat, Université Pierre et Marie Curie, Paris.
36. V. Yurinsky. *Sums and Gaussian Vectors*, volume 1617 of *Lecture Notes in Mathematics*. Springer-Verlag, Berlin, 1995.
37. C. Zălinescu. *Convex Analysis in General Vector Spaces*. World Scientific, River Edge, NJ, 2002.
38. Z. Zou and T. Hastie. Regularization and variable selection via the elastic net. *Journal of the Royal Statistical Society: Series B*, 67:301–320, 2005.



## A An auxiliary result

The following result is a generalization of [12, Proposition 5.14].

**Lemma 7** *Let  $\mathbb{K}$  be an at most countable set. For every  $k \in \mathbb{K}$ , let  $C_k$  be a closed interval in  $\mathbb{R}$  such that  $0 \in C_k$ , let  $D_k$  be a nonempty closed bounded interval in  $\mathbb{R}$ , and let  $h_k \in \Gamma_0^+(\mathbb{R})$  be such that  $h_k(0) = 0$ . Set*

$$G: \ell^2(\mathbb{K}) \rightarrow ]-\infty, +\infty]: (\xi_k)_{k \in \mathbb{K}} \mapsto \sum_{k \in \mathbb{K}} g_k(\xi_k), \quad \text{where } g_k = \mathbf{1}_{C_k} + \sigma_{D_k} + h_k. \quad (\text{A.1})$$

Let  $r \in ]1, 2]$  and consider the following statements:

- (a)  $\sum_{k \in \mathbb{K}} |(\min D_k)_+|^2 < +\infty$  and  $\sum_{k \in \mathbb{K}} |(\max D_k)_-|^2 < +\infty$ .
- (b)  $\sum_{k \in \mathbb{K}} h_k^*(-\min D_k)_+ < +\infty$  and  $\sum_{k \in \mathbb{K}} h_k^*((\max D_k)_-) < +\infty$ .
- (c)  $\sum_{k \in \mathbb{K}} |(\min D_k)_+|^{r^*} < +\infty$  and  $\sum_{k \in \mathbb{K}} |(\max D_k)_-|^{r^*} < +\infty$ .

Then the following hold:

- (i) Suppose that (a) or (b) is satisfied. Then  $G \in \Gamma_0(\ell^2(\mathbb{K}))$ .
- (ii) Suppose that (b) is satisfied. Then  $\inf G(\ell^2(\mathbb{K})) > -\infty$ .
- (iii) Suppose that, for every  $k \in \mathbb{K}$ ,  $h_k \geq \eta|\cdot|^r$  for some  $\eta \in \mathbb{R}_{++}$ . Then (a) $\Rightarrow$ (c) $\Rightarrow$ (b).
- (iv) Suppose that, for every  $k \in \mathbb{K}$ ,  $h_k - \eta|\cdot|^r \in \Gamma_0^+(\mathbb{R})$  for some  $\eta \in \mathbb{R}_{++}$  and that (c) holds. Then, for every  $\eta' \in ]0, \eta[$ , there exists  $H \in \Gamma_0(\ell^2(\mathbb{K}))$  such that  $G: u \mapsto H(u) + \eta' \sum_{k \in \mathbb{K}} |\mu_k|^r$ ,  $\text{dom } G \subset \ell^r(\mathbb{K})$ , and  $G$  is coercive in  $\ell^2(\mathbb{K})$ .

*Proof* We first observe that, if there exist  $(\chi_k)_{k \in \mathbb{K}} \in \ell_+^1(\mathbb{K})$  and  $b \in \mathbb{R}_+$  such that

$$(\forall k \in \mathbb{K}) \quad -g_k \leq \chi_k + b|\cdot|^2, \quad (\text{A.2})$$

then  $G \in \Gamma_0(\ell^2(\mathbb{K}))$ .

(i): Let  $k \in \mathbb{K}$ . Since

$$(\forall \mu \in \mathbb{R}) \quad \sigma_{D_k}(\mu) = \begin{cases} \mu \max D_k & \text{if } \mu \geq 0 \\ \mu \min D_k & \text{if } \mu < 0, \end{cases} \quad (\text{A.3})$$

we have

$$(\forall \mu \in \mathbb{R}) \quad -g_k(\mu) \leq -\sigma_{D_k}(\mu) - h_k(\mu) \leq \max\{\mu_-(\min D_k)_+, \mu_+(\max D_k)_-\} - h_k(\mu). \quad (\text{A.4})$$

Hence, in order to guarantee condition (A.2) for some  $(\chi_k)_{k \in \mathbb{K}} \in \ell_+^1(\mathbb{K})$  and  $b \in \mathbb{R}_{++}$ , it is sufficient to require condition (a) or (b) (note that  $h_k^* \geq 0$ , since  $h_k(0) = 0$ ). Therefore in this case  $G \in \Gamma_0(\ell^2(\mathbb{K}))$ .

(ii): It follows from (A.4) that

$$(\forall k \in \mathbb{K}) \quad -g_k \leq \max\{h_k^*(-(\min D_k)_+), h_k^*((\max D_k)_-)\}. \quad (\text{A.5})$$

Hence, for every  $u \in \ell^2(\mathbb{K})$ ,  $-G(u) \leq \sum_{k \in \mathbb{K}} \max\{h_k^*(-(\min D_k)_+), h_k^*((\max D_k)_-)\} < +\infty$ .

(iii): For every  $k \in \mathbb{K}$ ,  $h_k^* \leq (r\eta)^{1-r^*} (r^*)^{-1} |\cdot|^{r^*}$ . The statement therefore follows by observing that, since  $2 \leq r^*$ ,  $\ell^2(\mathbb{K}) \subset \ell^{r^*}(\mathbb{K})$ .

(iv): Setting, for every  $k \in \mathbb{K}$ ,  $\tilde{h}_k = h_k - \eta'|\cdot|^r$ , we have  $g_k = \mathbf{1}_{C_k} + \sigma_{D_k} + \tilde{h}_k + \eta'|\cdot|^r$ , with  $(\eta - \eta')|\cdot|^r \leq \tilde{h}_k \in \Gamma_0^+(\mathbb{R})$ . It follows from (i) and (iii) that, for every  $u = (\mu_k)_{k \in \mathbb{K}} \in \ell^2(\mathbb{K})$ ,  $G(u) = H(u) + \eta' \sum_{k \in \mathbb{K}} |\mu_k|^r$ , for some  $H \in \Gamma_0(\ell^2(\mathbb{K}))$ .

## B Proximity operators of power functions

It follows from [7, Example 4.4] that, for every  $\gamma \in \mathbb{R}_{++}$  and every  $r \in [1, 2]$ ,

$$(\forall \mu \in \mathbb{R}) \quad \text{prox}_{\gamma|\cdot|^r} \mu = \xi \text{sign}(\mu), \quad \text{where } \xi \geq 0 \quad \text{and} \quad \xi + r\gamma\xi^{r-1} = |\mu|. \quad (\text{B.1})$$

Equation (B.1) can be solved explicitly for  $r \in \{3/2, 4/3, 5/4\}$  [7, 35]. However, in general, it must be solved iteratively.

**Proposition 4** *Let  $\mu \in \mathbb{R}$ , let  $\gamma \in \mathbb{R}_{++}$ , let  $r \in [1, 2]$ , and let  $(r_1, r_2) \in [1, 2]^2$ , be such that  $r_1 < r_2$ . Then the following hold:*

- (i)  $\text{prox}_{\gamma|\cdot|^r} : \mathbb{R} \rightarrow \mathbb{R}$  is strictly increasing, nonexpansive, odd, and differentiable, and  $\text{prox}_{\gamma|\cdot|^r} + \text{id}_{\mathbb{R}_+}$  is convex.
- (ii) We have

$$\min \left\{ \frac{|\mu|}{1+r\gamma}, \left( \frac{|\mu|}{1+r\gamma} \right)^{\frac{1}{r-1}} \right\} \leq |\text{prox}_{\gamma|\cdot|^r} \mu| \leq \max \left\{ \frac{|\mu|}{1+r\gamma}, \left( \frac{|\mu|}{1+r\gamma} \right)^{\frac{1}{r-1}} \right\}. \quad (\text{B.2})$$

(iii) Suppose that  $|\mu| > 1 + r_2\gamma$ . Then  $|\text{prox}_{\gamma|\cdot|^{r_2}} \mu| < |\text{prox}_{\gamma|\cdot|^{r_1}} \mu|$ .

(iv) Suppose that  $r > 1$  and that  $|\mu| > 1 + r\gamma$ . Then  $\frac{|\mu|}{1+r\gamma} \leq |\text{prox}_{\gamma|\cdot|^r} \mu| < |\mu| - \gamma$ .

*Proof* (i): It follows from [9, Lemma 2.2(iv) and Proposition 2.4] that  $\text{prox}_{\tau|\cdot|}$  is nonexpansive, increasing, and odd. Now set  $\psi : \mathbb{R}_+ \rightarrow \mathbb{R}_+ : \xi \mapsto \xi + r\tau\xi^{r-1}$ . Clearly  $\psi$  is strictly increasing and concave. Moreover it is differentiable on  $\mathbb{R}_{++}$  and, for every  $\xi \in \mathbb{R}_{++}$ ,  $\psi'(\xi) = 1 + r(r-1)\tau\xi^{r-2}$ . Hence, from (B.1), for every  $\mu \in \mathbb{R}_+$ ,  $\text{prox}_{\tau|\cdot|} \mu = \psi^{-1}(\mu)$ . This shows that  $\text{prox}_{\tau|\cdot|}$  is strictly increasing, convex, differentiable on  $\mathbb{R}_{++}$  with, for every  $\mu \in \mathbb{R}_{++}$ ,  $(\text{prox}_{\tau|\cdot|})' \mu = 1/\psi'(\psi^{-1}(\mu))$ , that is

$$(\text{prox}_{\gamma|\cdot|^r})' \mu = \left( 1 + \frac{r(r-1)\gamma}{(\text{prox}_{\gamma|\cdot|^r} \mu)^{2-r}} \right)^{-1}. \quad (\text{B.3})$$

(ii): According to (B.1), there exists  $\xi \in \mathbb{R}_+$  such that  $\text{prox}_{\tau|\cdot|^r} \mu = \text{sign}(\mu)\xi$  and  $\xi + r\tau\xi^{r-1} = |\mu|$ . If  $\xi \geq 1$ , then  $|\mu| = \xi + r\tau\xi^{r-1} \leq (1+r\tau)\xi$ , hence  $|\mu|/(1+r\tau) \leq \xi = |\text{prox}_{\tau|\cdot|^r} \mu|$ . If  $\xi < 1$ , then  $|\mu| = \xi + r\tau\xi^{r-1} \leq (1+r\tau)\xi^{r-1}$ , hence  $(|\mu|/(1+r\tau))^{1/(r-1)} \leq \xi = |\text{prox}_{\tau|\cdot|^r} \mu|$ . The first inequality in (B.2) follows and the second is proved analogously.

(iii): In view of (B.1) there exist  $\xi_1 \in \mathbb{R}_+$  and  $\xi_2 \in \mathbb{R}_+$  such that

$$\begin{cases} \text{prox}_{\tau|\cdot|^{r_1}} \mu = \text{sign}(\mu)\xi_1 & \text{and} & \xi_1 + r_1\tau\xi_1^{r_1-1} = |\mu| \\ \text{prox}_{\tau|\cdot|^{r_2}} \mu = \text{sign}(\mu)\xi_2 & \text{and} & \xi_2 + r_2\tau\xi_2^{r_2-1} = |\mu|. \end{cases} \quad (\text{B.4})$$

If  $|\mu| > 1 + \tau r_2 > 1 + \tau r_1$ , it follows from (B.2) that

$$1 < \frac{|\mu|}{1+r_1\tau} \leq |\xi_1| \quad \text{and} \quad 1 < \frac{|\mu|}{1+r_2\tau} \leq |\xi_2|. \quad (\text{B.5})$$

Therefore, since  $r_1 < r_2$  and  $\xi_1 > 1$ ,

$$\xi_2 + r_2\tau\xi_2^{r_2-1} = |\mu| = \xi_1 + r_1\tau\xi_1^{r_1-1} < \xi_1 + r_2\tau\xi_1^{r_2-1}. \quad (\text{B.6})$$

Hence, since  $\xi \mapsto \xi + r_2\tau\xi^{r_2-1}$  is strictly increasing on  $\mathbb{R}_+$ , we conclude that  $\xi_2 < \xi_1$ .

(iv): Since (B.1) implies that  $\text{prox}_{\tau|\cdot|^r} \mu = \text{sign}(\mu)(|\mu| - \tau)$ , we derive from (iii) that

$$|\mu| > 1 + r\tau \quad \Rightarrow \quad |\text{prox}_{\tau|\cdot|^r} \mu| < |\mu| - \tau, \quad (\text{B.7})$$

The first inequality in (iv) follows directly from (B.2).

*Remark 9*

- (i) The bounds given in (B.2) can be useful to initialize the bisection method to solve (B.1).
- (ii)  $(\text{prox}_{\gamma|\cdot|^r})' 0 = 0$ ,  $(\text{prox}_{\gamma|\cdot|^r})' \mu \leq 1$  and  $(\text{prox}_{\gamma|\cdot|^r})' \mu \rightarrow 1$  as  $\mu \rightarrow +\infty$ .
- (iii)  $\text{prox}_{\gamma|\cdot|^r}$  has no asymptote as  $\mu \rightarrow +\infty$ , since (B.1) yields  $\text{prox}_{\gamma|\cdot|^r} \mu - \mu = -r\gamma(\text{prox}_{\gamma|\cdot|^r} \mu)^{r-1} \rightarrow -\infty$  as  $\mu \rightarrow +\infty$ .