

## VARIATIONAL PROPERTIES OF VALUE FUNCTIONS\*

ALEKSANDR Y. ARAVKIN<sup>†</sup>, JAMES V. BURKE<sup>‡</sup>, AND MICHAEL P. FRIEDLANDER<sup>§</sup>

**Abstract.** Regularization plays a key role in a variety of optimization formulations of inverse problems. A recurring theme in regularization approaches is the selection of regularization parameters and their effect on the solution and on the optimal value of the optimization problem. The sensitivity of the value function to the regularization parameter can be linked directly to the Lagrange multipliers. This paper characterizes the variational properties of the value functions for a broad class of convex formulations, which are not all covered by standard Lagrange multiplier theory. An inverse function theorem is given that links the value functions of different regularization formulations (not necessarily convex). These results have implications for the selection of regularization parameters, and the development of specialized algorithms. Numerical examples illustrate the theoretical results.

**Key words.** convex optimization, optimal value functions, Lagrange multipliers, inverse problems

**AMS subject classifications.** 65K05, 65K10, 90C25

**DOI.** 10.1137/120899157

**1. Introduction.** It is well known that there is a close connection between the sensitivity of the optimal value of a parametric optimization problem and its Lagrange multipliers. Consider the family of feasible convex optimization problems

$$\mathcal{P}(b, \tau) \quad \underset{r, x}{\text{minimize}} \quad \rho(r) \quad \text{subject to} \quad Ax + r = b, \quad \phi(x) \leq \tau,$$

where  $b \in \mathbb{R}^m$ ,  $A \in \mathbb{R}^{m \times n}$ , and the functions  $\phi : \mathbb{R}^n \rightarrow \overline{\mathbb{R}} := (-\infty, \infty]$  and  $\rho : \mathbb{R}^m \rightarrow \overline{\mathbb{R}}$  are closed, proper, and convex and continuous relative to their domains. The value function

$$v(b, r) := \inf_{r, x} \{ \rho(r) \mid Ax + r = b, \phi(x) \leq \tau \}$$

gives the optimal objective value of problem  $\mathcal{P}(b, \tau)$  for fixed parameters  $b$  and  $\tau$ . If  $\mathcal{P}(b, \tau)$  is a feasible *ordinary convex program* [34, section 28], then under standard hypotheses the subdifferential of  $v$  is the set of pairs  $(u, \mu)$ , where  $u \in \mathbb{R}^m$  and  $\mu \in \mathbb{R}$  are the Lagrange multipliers of  $\mathcal{P}(b, \tau)$  corresponding to the equality and inequality constraints, respectively. This connection is extensively explored in Rockafeller's 1993 survey paper [35].

If we allow  $\phi$  to take on infinite values on the domain of the objective—which can occur, for example, if  $\phi$  is an arbitrary gauge—then  $\mathcal{P}(b, \tau)$  is no longer an ordinary convex program, and so the standard Lagrange multiplier theory does not apply.

---

\*Received by the editors November 16, 2012; accepted for publication (in revised form) May 24, 2013; published electronically August 22, 2013.

<http://www.siam.org/journals/siopt/23-3/89915.html>

<sup>†</sup>IBM T.J. Watson Research Center, Yorktown Heights, NY 10598 (saravkin@us.ibm.com). The work of this author was partially supported by an NSERC CRD grant through the University of British Columbia.

<sup>‡</sup>Department of Mathematics, University of Washington, Seattle, WA 98195 (burke@math.washington.edu).

<sup>§</sup>Department of Computer Science, University of British Columbia, Vancouver B.C., V6T 1Z4, Canada (mpf@cs.ubc.ca). The work of this author was partially supported by NSERC Discovery grant 312104.

Multiplier theories that do apply to more general contexts can be found in [21, 16, 45, 8]. Remarkably, even in this general setting, it is possible to obtain explicit formulas of the subdifferential of the value function  $v$  useful in many applications.

**1.1. Examples.** We give two simple examples that illustrate the need for the extended Lagrange multiplier theory. Both are of the form

$$(1.1) \quad \underset{x}{\text{minimize}} \quad \frac{1}{2} \|Ax - b\|_2^2 \quad \text{subject to} \quad \gamma(x | U) \leq 1,$$

where

$$\gamma(x | U) := \inf \{ \lambda \geq 0 \mid x \in \lambda U \}$$

is the *gauge function* for the nonempty closed convex set  $U \subset \mathbb{R}^n$ , which contains 0. Let  $A = I$  and  $b = (0, -1)^T$ . Then the solution to (1.1) is just the 2-norm projection onto the set  $\{x \mid \gamma(x | U) \leq 1\} = U$ .

For our first example, we consider the set

$$U = \{x \in \mathbb{R}^2 \mid \frac{1}{2}x_1^2 \leq x_2\},$$

defined in [34, section 10]. The gauge for this set is an example of a closed, proper, and convex function that is not locally bounded and therefore not continuous at a point in its effective domain. It is straightforward to show that

$$\gamma(x | U) = \begin{cases} \frac{x_2^2}{2x_1}, & x_2 > 0, \\ 0, & x_1 = 0 = x_2, \\ +\infty, & \text{otherwise.} \end{cases}$$

The constraint region for (1.1) is the set  $U$ , and the unique global solution is the point  $x = 0$ . However, since  $0 = \gamma(0 | U) < 1$ , the classical Lagrange multiplier theory fails: the solution is on the boundary of the feasible region, and yet no classical Lagrange multiplier exists. The problem is that the constraint is active at the solution but not active in the functional sense, i.e.,  $\gamma(0 | U) < 1$ . In contrast, the extended multiplier theory of [45, Theorem 2.9.3] succeeds with the multiplier choice of 0.

For the second example, take  $U = \mathbb{B}_2 \cap K$ , where  $\mathbb{B}_2$  is the unit ball associated with the Euclidean norm on  $\mathbb{R}^2$ . Then  $\gamma(x | \mathbb{B}_2 \cap K) = \|x\|_2 + \delta(x | K)$ , and the constraint region for (1.1) is the set  $\mathbb{B}_2 \cap K$ . Set  $K = \{(x_1, x_2) \mid x_2 \geq 0\}$ . Again, the origin is the unique global solution to this optimization problem, and no classical Lagrange multiplier for this problem exists.

In both of these examples, the multiplier theory in [45] can be applied to obtain a Lagrange multiplier theorem. In Theorem 5.2, we extend this theory and provide a characterization of these Lagrange multipliers that is useful in computation.

**1.2. Formulations.** Appropriate definitions of the functions  $\rho$  and  $\phi$  can be used to represent a range of practical problems. Choosing  $\rho$  to be the 2-norm and  $\phi$  to be any norm yields the canonical regularized least-squares problem

$$(1.2) \quad \underset{x, r}{\text{minimize}} \quad \|r\|_2 \quad \text{subject to} \quad Ax + r = b, \quad \|x\| \leq \tau,$$

which optimizes the misfit between the data  $b$  and the forward model  $Ax$ , subject to keeping  $x$  appropriately bounded in some norm. The 2-norm constraint on  $x$  yields a Tikhonov regularization, popular in many inversion applications. A 1-norm constraint on  $x$  yields the Lasso problem [41], often used in sparse recovery and model-selection applications. Interestingly, when the optimal residual  $r$  is nonzero, the value function for this family of problems is always differentiable in both  $b$  and  $\tau$  with

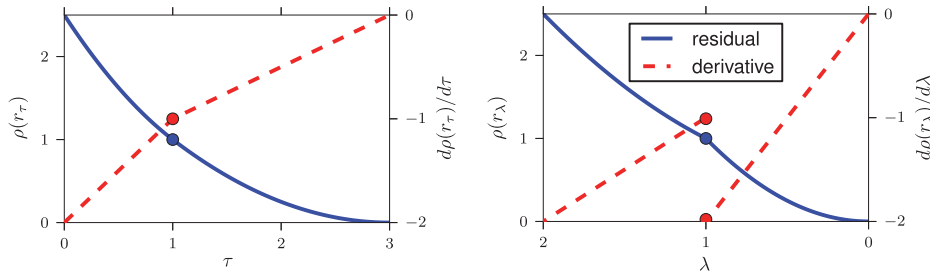


FIG. 1.1. The misfit  $\rho(r)$  (solid line) and its derivative (dashed line) as a function of the regularization parameter for a 1-norm regularized example. The left panel shows the constrained formulation  $\mathcal{P}(b, \tau)$ , which varies smoothly with  $\tau$ ; the right panel shows that the penalized formulation does not vary smoothly with  $\lambda$  (note the reversed axis).

$$\nabla v(b, \tau) = \left( \frac{r}{\|r\|_2}, \frac{\|A^T r\|_*}{\|r\|_2} \right),$$

where  $\|\cdot\|_*$  is the norm dual to  $\|\cdot\|$ . This gradient is derived by van den Berg and Friedlander [10, Theorem 2.2]. The analysis of the sensitivity in  $\tau$  of the value function for the Lasso problem led to the development of the SPGL1 solver [9], currently used in a variety of sparse inverse problems, with particular success in large-scale sparse inverse problems [27]. A subsequent analysis [12] that allows  $\phi(x)$  to be a gauge paved the way for other applications, such as group-sparsity promotion [11].

An alternative to  $\mathcal{P}(b, \tau)$  is the class of penalized formulations

$$\mathcal{P}_L(b, \lambda) \quad \underset{x}{\text{minimize}} \quad \rho(b - Ax) + \lambda\phi(x).$$

(The subscript “L” in the label reminds us that it can be interpreted as a Lagrangian of the original problem.) The nonnegative regularization parameter  $\lambda$  is used to control the tradeoff between the data misfit  $\rho$  and the regularization term  $\phi$ . For example, taking  $\rho(r) = \|r\|_2$  and  $\phi(x) = \|x\|$  yields a formulation analogous to (1.2). This penalized formulation is commonly used in applications of Bayesian parametric regression [31, 37, 30, 42, 44], inference problems on dynamic linear systems [1, 15], feature selection, selective shrinkage, and compressed sensing [25, 20, 19], robust formulations [29, 24, 2, 23], support-vector regression [43, 26], classification [22, 33, 39], and functional reconstruction [6, 38, 17].

From an algorithmic point of view, the unconstrained formulation  $\mathcal{P}_L(b, \lambda)$  may be preferable. However, the constrained formulation  $\mathcal{P}(b, \tau)$  has the distinction that its value function  $v(b, \tau)$  is jointly convex in its parameters; see section 1.3. In contrast, the optimal value of the penalized formulation  $\mathcal{P}_L(b, \lambda)$  is not in general a convex function of its parameters. The following simple example

$$\rho(r) = \frac{1}{2}\|r\|_2^2, \quad \phi(x) = \|x\|_1, \quad A = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \quad b = \begin{bmatrix} 2 \\ 1 \end{bmatrix}$$

illustrates this situation. The optimal values of  $\rho$  in the formulations  $\mathcal{P}(b, \tau)$  and  $\mathcal{P}_L(b, \lambda)$ , as functions of  $\tau$  and  $\lambda$ , respectively, are given by

$$\rho_\tau = \begin{cases} \frac{1}{2} + \frac{1}{2}(\tau - 2)^2 & \text{for } \tau \in [0, 1), \\ \frac{1}{4}(\tau - 3)^2 & \text{for } \tau \in [1, 3), \\ 0 & \text{otherwise;} \end{cases} \quad \rho_\lambda = \begin{cases} \lambda^2 & \text{for } \lambda \in [0, 1), \\ \frac{1}{2} + \frac{1}{2}\lambda^2 & \text{for } \lambda \in [1, 2), \\ 5/2 & \text{otherwise.} \end{cases}$$

The optimal values and their derivatives are shown in Figure 1.1, where it is clear that  $\rho_\tau$  is convex (and in this case also smooth) in  $\tau$ , but  $\rho_\lambda$  is not convex in  $\lambda$ .

The admissibility of variational analysis and convexity of the value function may convince some practitioners to explore formulations of type  $\mathcal{P}(b, \tau)$  rather than  $\mathcal{P}_L(b, \lambda)$ . In fact, we give an example (in section 7) of how this variational information can be used for algorithm design in the context of large-scale inverse problems.

**1.3. Approach.** For many practical inverse problems, the formulation of primary interest is the residual-constrained formulation

$$\mathcal{P}_R(b, \sigma) \quad \underset{x}{\text{minimize}} \quad \phi(x) \quad \text{subject to} \quad \rho(b - Ax) \leq \sigma,$$

(the subscript “R” reminds us that this formulation reverses the objective and constraint functions from that of  $\mathcal{P}(b, \tau)$ ) in part because estimates of a tolerance level  $\sigma$  on fitting the error  $\rho(b - Ax)$  are more easily available than estimates of a bound on the penalty parameter on the regularization  $\phi$ ; cf.  $\mathcal{P}_L(b, \lambda)$ . However, the formulation  $\mathcal{P}(b, \tau)$  can sometimes be easier to solve. The underlying numerical theme is to develop methods for solving  $\mathcal{P}_R(b, \sigma)$  that use a sequence of solutions to the possibly easier problem  $\mathcal{P}(b, \tau)$ .

In section 2, we present an inverse function theorem for value functions that characterizes the relationship between  $\mathcal{P}(b, \tau)$  and  $\mathcal{P}_R(b, \sigma)$  and applies more generally to nonconvex problems. Pairs of problems of this type are classical, though typically paired in a max-min fashion. For example, the isoperimetric inequality and Queen Dido’s problem are of this type; the greatest area surrounded by a curve of given length is related to the problem of finding the curve of least arc length surrounding a given area. (See [40] for a modern survey.) The Markowitz mean-variance portfolio theory is also based on such a pairing; minimizing volatility subject to a lower bound on expected return is related to maximizing expected return subject to an upper bound on volatility [32].

The application motivating our investigation is establishing conditions under which it is possible to implement a root-finding approach for the nonlinear equation

$$(1.3) \quad \text{find } \tau \text{ such that } v(b, \tau) = \sigma,$$

where  $\mathcal{P}_R(b, \sigma)$  can be solved via a sequence of approximate solutions of  $\mathcal{P}(b, \tau)$ . This generalizes the approach used by van den Berg and Friedlander [10, 12] for large-scale sparse optimization applications. The convex case is especially convenient, because both value functions are decreasing and convex. When the value function is differentiable, Newton’s method is globally monotonic and locally quadratic. In section 5 we establish the variational properties (including conditions necessary for differentiability) of  $\mathcal{P}(b, \tau)$ .

In section 4 we derive dual representations of  $\mathcal{P}(b, \tau)$  and their optimality conditions. These are used in section 5 to characterize the variational properties of the value function  $v$ . The conjugate, horizon, and perspective functions arise naturally as part of the analysis, and we present a calculus (section 3) for these functions that allows explicit computation of the subdifferential of  $v$  for large classes of misfit functions  $\rho$  and regularization functions  $\phi$  (see section 6).

One of the motivating problems for the general analysis and methods we present is the treatment of a robust misfit function  $\rho$  (such as the popular Huber penalty) in the context of sparsity promotion, which typically involves a nonsmooth regularizer  $\phi$ . In

section 7 we demonstrate that the sensitivity analysis can be applied to solve a sparse nonnegative denoising problem with convex and nonconvex robust misfit measures.

The proofs of all of the results are relegated to the appendix (section 8).

**1.4. Notation.** For a matrix  $A \in \mathbb{R}^{m \times n}$ , the image and inverse image of the sets  $E$  and  $F$ , respectively, are given by the sets

$$AE = \{y \mid y = Ax, x \in E\} \quad \text{and} \quad A^{-1}F = \{x \mid Ax \in F\}.$$

For a function  $p : \mathbb{R}^n \rightarrow \overline{\mathbb{R}}$ , its epigraph is denoted  $\text{epi } p = \{(x, \mu) \mid p(x) \leq \mu\}$ , and its level set is denoted  $\text{lev}_p(\tau) = \{x \mid p(x) \leq \tau\}$ . The function  $p$  is said to be proper if  $\text{dom } p \neq \emptyset$  and closed if  $\text{epi } p$  is a closed set. The function  $\delta(x \mid X)$  is the indicator to a set  $X$ , i.e.,  $\delta(x \mid X) = 0$  if  $x \in X$  and  $\delta(x \mid X) = +\infty$  if  $x \notin X$ .

**2. An inverse function theorem for optimal value functions.** Let  $\psi_i : X \subseteq \mathbb{R}^n \rightarrow \overline{\mathbb{R}}$ ,  $i \in \{1, 2\}$ , be arbitrary scalar-valued functions, and consider the following pair of related problems and their associated value functions:

$$\begin{aligned} \mathcal{P}_{1,2}(\tau) \quad v_1(\tau) &:= \inf_{x \in X} \psi_1(x) + \delta((x, \tau) \mid \text{epi } \psi_2), \\ \mathcal{P}_{2,1}(\sigma) \quad v_2(\sigma) &:= \inf_{x \in X} \psi_2(x) + \delta((x, \sigma) \mid \text{epi } \psi_1). \end{aligned}$$

This pair corresponds to the problems  $\mathcal{P}(b, \tau)$  and  $\mathcal{P}_R(b, \sigma)$ , defined in section 1, with the identifications

$$\psi_1(x) = \rho(b - Ax) \quad \text{and} \quad \psi_2(x) = \phi(x).$$

Our goal in this section is to establish general conditions under which the value functions  $v_1$  and  $v_2$  satisfy the inverse-function relationship

$$v_1 \circ v_2 = \text{id},$$

and for which the the pair of problems  $\mathcal{P}_{1,2}(\tau)$  and  $\mathcal{P}_{2,1}(\sigma)$  have the same solution sets. The pair of problems  $\mathcal{P}(b, \tau)$  and  $\mathcal{P}_R(b, \sigma)$  always satisfy the conditions of the next theorem, which applies to functions that are not necessarily convex.

**THEOREM 2.1.** *Let  $\psi_i : X \subseteq \mathbb{R}^n \rightarrow \overline{\mathbb{R}}$ ,  $i \in \{1, 2\}$ , be as defined in  $\mathcal{P}_{1,2}(\tau)$ , and define*

$$\mathcal{S}_{1,2} := \{\tau \in \overline{\mathbb{R}} \mid \emptyset \neq \arg \min \mathcal{P}_{1,2}(\tau) \subset \{x \in X \mid \psi_2(x) = \tau\}\}.$$

*Let  $\mathcal{S}_{2,1}$  be defined symmetrically to  $\mathcal{S}_{1,2}$  by interchanging the roles of the indices. Then, for every  $\tau \in \mathcal{S}_{1,2}$ ,*

- (a)  $v_2(v_1(\tau)) = \tau$  and
- (b)  $\arg \min \mathcal{P}_{1,2}(\tau) = \arg \min \mathcal{P}_{2,1}(v_1(\tau)) \subset \{x \in X \mid \psi_1(x) = v_1(\tau)\}$ .

*Moreover,  $\mathcal{S}_{2,1} = \{v_1(\tau) \mid \tau \in \mathcal{S}_{1,2}\}$ , and so*

$$\{(\tau, v_1(\tau)) \mid \tau \in \mathcal{S}_{1,2}\} = \{(v_2(\sigma), \sigma) \mid \sigma \in \mathcal{S}_{2,1}\}.$$

**3. Convex analysis.** In order to present the duality results of section 4, we require a few basic tools from convex analysis. There are many excellent references for the necessary background material, with several appearing within the past 10 years. In this study we make use of Rockafellar [34] and Rockafellar and Wets [36], although similar results can be found elsewhere [8, 13, 14, 21, 28, 45]. We review the necessary results here.

**3.1. Functional operations.** The proper convex function  $h : \mathbb{R}^n \rightarrow \overline{\mathbb{R}}$  generates the following convex functions:

1. *Legendre–Fenchel conjugate* of  $h$ :

$$h^*(y) := \sup_x [\langle y, x \rangle - h(x)].$$

2. *Horizon function* of  $h$ :

$$h^\infty(z) := \sup_{x \in \text{dom } h} [h(x+z) - h(x)].$$

3. *Perspective function* of  $h$ :

$$\tilde{h}(z, \lambda) := \begin{cases} \lambda h(\lambda^{-1}z) & \text{if } \lambda > 0, \\ \delta(z \mid 0) & \text{if } \lambda = 0, \\ +\infty & \text{if } \lambda < 0. \end{cases}$$

4. *Closure of the perspective function* of  $h$ :

$$h^\pi(z, \lambda) := \begin{cases} \lambda h(\lambda^{-1}z) & \text{if } \lambda > 0, \\ h^\infty(z) & \text{if } \lambda = 0, \\ +\infty & \text{if } \lambda < 0. \end{cases}$$

Each of these functions can also be defined by considering the epigraphical perspective and properties of convex sets. Indeed, the horizon function  $h^\infty$  is usually defined to be the function whose epigraph is the horizon cone of the epigraph of  $h$  (see section 3.2 below). The definition given above is a consequence of [34, Theorem 8.5].

The perspective function of  $h$  is the positively homogeneous function generated by the convex function  $\hat{h}(x, \lambda) := h(x) + \delta(\lambda \mid \{1\})$  [34, pp. 35 and 67]. If  $h$  is additionally closed and proper, then so are  $h^*$  (Theorem 12.2),  $h^\infty$  (Theorem 8.5), and  $h^\pi$  (Corollary 8.5.2), where these results are from Rockafellar [34].

Note that for every closed, proper, and convex function  $h$ , the associated horizon and perspective function,  $h^\infty$  and  $h^\pi$ , are positively homogeneous and so can be represented as the support functional for some convex set [34, Theorem 13.2]. Moreover, if  $h$  is a support function, then  $h^\infty = h^\pi = h$ .

**3.2. Cones.** We associate the following cones with a convex set  $C$  and a convex function  $h$ :

1. *Polar cone*: The polar cone of  $C$  is denoted by

$$C^\circ := \{x^* \mid \langle x^*, x \rangle \leq 0 \ \forall x \in C\}.$$

2. *Recession cone*: The recession cone of  $C$  is denoted by

$$C^\infty := \{x \mid C + x \subset C\} = \{x \mid y + \lambda x \in C \ \forall \lambda \geq 0, \forall y \in C\}.$$

3. *Barrier cone*: The barrier cone of  $C$  is denoted by

$$\text{bar}(C) := \{x^* \mid \text{for some } \beta \in \mathbb{R}, \langle x, x^* \rangle \leq \beta \ \forall x \in C\}.$$

4. *Horizon cone of  $h$* : The horizon cone [34, Theorem 8.7] of  $h$  is denoted by

$$\text{hzn}(h) := \{y \mid h^\infty(y) \leq 0\} = [\text{lev}_h(\tau)]^\infty \ \forall \tau > \inf h.$$

A further excellent reference for horizon cones and functions is [7], where they are referred to as *asymptotic* cones and functions.

**3.3. Calculus rules.** The conjugate, horizon, and perspective transformations defined in section 3.1 possess a rich calculus. We use this calculus to obtain explicit expressions for the functions  $\rho^*$ ,  $\phi^*$ ,  $(\phi^*)^\infty$ , and  $(\phi^*)^\pi$ , which play a crucial role in the applications of section 6. The calculus for conjugates and horizons is developed in many references (e.g., [8, 13, 14, 21, 28, 45]); specific citations from [34] are provided. In order to establish the perspective calculus rules for affine composition and the inverse linear image, we note that addition is a special case of affine composition, and that infimal convolution is a special case of inverse linear image. Hence, we need only establish the perspective calculus formulas for affine composition and the inverse linear image: the formula for affine-composition follows from [34, Theorem 9.5] and the definition of the perspective transformation; the formula for inverse linear image is established in section 8.

**Affine composition.** Let  $p : \mathbb{R}^m \rightarrow \overline{\mathbb{R}}$  be a closed proper convex function,  $A \in \mathbb{R}^{m \times n}$ , and  $b \in \mathbb{R}^m$ , such that  $(\text{Ran}(A) - b) \cap \text{ri}(\text{dom } p) \neq \emptyset$ . Let

$$h(x) := p(Ax - b).$$

Then

$$h^*(y) = \inf_{A^T u = y} [\langle b, u \rangle + p^*(u)] \quad [34, \text{Theorem 16.3}],$$

$$h^\infty(z) = p^\infty(Az) \quad [34, \text{Theorem 9.5}],$$

$$h^\pi(x, \lambda) = p^\pi(Ax - \lambda b, \lambda),$$

where, for  $\lambda = 0$ ,

$$h^\pi(x, 0) = p^\pi(Ax, 0) = p^\infty(Ax).$$

All three functions are closed, proper, and convex. The derivation of  $h^*$  also makes use of the observation that

$$(3.1) \quad \text{if } g(x) := h(x - b), \quad \text{then } g^*(v) = h^*(v) + \langle v, b \rangle.$$

**Inverse linear image.** Let  $p : \mathbb{R}^n \rightarrow \overline{\mathbb{R}}$  be closed, proper, and convex, and let  $A \in \mathbb{R}^{m \times n}$ . Let

$$h(w) := \inf_{Ax=w} p(x)$$

be the inverse linear image of  $p$  under  $A$ . Then

$$h^*(y) = p^*(A^T y) \quad [34, \text{Theorem 16.3}].$$

If  $(A^T)^{-1} \text{ri}(\text{dom } p^*) \neq \emptyset$ , then

$$h^\infty(z) = \inf_{Ax=z} p^\infty(x) \quad [34, \text{Theorem 9.2}],$$

$$h^\pi(w, \lambda) = \inf_{Ax=w} p^\pi(x, \lambda) \quad (\text{Proof in section 8}),$$

where all of the functions  $h$ ,  $h^*$ ,  $h^\infty$ , and  $h^\pi$  are closed, proper, and convex.

**Addition.** Let  $h_i : \mathbb{R}^n \rightarrow \overline{\mathbb{R}}$ , for  $i = 1, \dots, m$ , be closed proper convex functions. If  $h := h_1 + \dots + h_m$  is not identically  $+\infty$ , then

$$\begin{aligned} h^\infty &= h_1^\infty + \dots + h_m^\infty & [34, \text{Theorem 9.2}], \\ h^\pi &= h_1^\pi + \dots + h_m^\pi, \end{aligned}$$

where both are closed, proper, and convex. Moreover, if  $\bigcap_{i=1}^m \text{ri}(\text{dom } h_i) \neq \emptyset$ , then

$$h^* = h_1^* \nabla \dots \nabla h_m^* \quad [34, \text{Theorem 16.4}]$$

is closed, proper, and convex, where  $\nabla$  denotes infimal convolution.

**Infimal convolution.** Let  $h_i : \mathbb{R}^n \rightarrow \overline{\mathbb{R}}$ , for  $i = 1, \dots, m$ , be closed, proper, and convex functions. Let  $h := h_1 \nabla \dots \nabla h_m$ . Then  $h^* = h_1^* + \dots + h_m^*$ , and

$$\text{if } \bigcap_{i=1}^m \text{ri}(\text{dom } h_i^*) \neq \emptyset, \quad \text{then } h^\infty = h_1^\infty \nabla \dots \nabla h_m^\infty \quad [34, \text{Corollary 9.2.1}],$$

and

$$h^\pi(x, \lambda) = \inf_{\sum_{i=1}^m x_i = x} [h_1^\pi(x_1, \lambda) + \dots + h_m^\pi(x_m, \lambda)].$$

All three functions are closed, proper, and convex.

**4. The dual problem.** For our analysis, it is convenient to consider the (equivalent) epigraphical formulation

$$(\mathcal{P}) \quad v(b, \tau) = \underset{x}{\text{minimize}} \quad f(x, b, \tau)$$

of  $\mathcal{P}(b, \tau)$ , where

$$f(x, b, \tau) := \rho(b - Ax) + \delta((x, \tau) \mid \text{epi } \phi).$$

Because the functions  $\rho$  and  $\phi$  are convex, it immediately follows that  $f$  is also convex. This fact gives the convexity of the value function  $v$ , since it is the inf-projection of the objective function in  $x$  [34, Theorem 5.3].

We use a duality framework derived from the one described in Rockafellar and Wets [36, Chapter 11, section H] and associate with  $\mathcal{P}$  its dual problem and corresponding dual value function:

$$(\mathcal{D}) \quad \hat{v}(b, \tau) := \underset{u, \mu}{\text{maximize}} \quad \langle b, u \rangle + \tau\mu - f^*(0, u, \mu).$$

To derive this dual from [36, Theorem 11.39], define

$$f_{(b, \tau)}(x, \Delta b, \Delta \tau) := f(x, b + \Delta b, \tau + \Delta \tau).$$

Then, by (3.1),  $f_{(b, \tau)}^*(v, u, \mu) = f^*(v, u, \mu) - \langle b, u \rangle - \tau\mu$ . Substituting this expression into [36, Theorem 11.39] gives  $\mathcal{D}$ .

The dual  $\mathcal{D}$  is the key to understanding the variational behavior of the value function. To access these results we must compute the conjugate of  $f$ . For this it is useful to have an alternative representation for the support function of the epigraph, which is the conjugate of the indicator function appearing in  $f$ .



**4.1. Reduced dual problem.** In Theorem 4.2, we derive an equivalent representation of the dual problem  $\mathcal{D}$  in terms of  $u$  alone. This is the *reduced* dual problem for  $\mathcal{P}$ . We first present a result about conjugates for epigraphs and lower level sets.

LEMMA 4.1 (Conjugates for epigraphs and lower level sets). *Let  $h : \mathbb{R}^n \rightarrow \overline{\mathbb{R}}$  be closed, proper, and convex. Then*

$$(4.1a) \quad \delta^*((y, \mu) \mid \text{epi } h) = (h^*)^\pi(y, -\mu),$$

$$(4.1b) \quad \delta^*(y \mid \text{lev}_h(\tau)) = \text{cl} \left( \inf_{\mu \geq 0} [\tau\mu + (h^*)^\pi(y, \mu)] \right).$$

Expressions (4.1b) and (4.1a) are easily derived from the case where  $\tau = 0$  which is established in [34, Theorem 13.5] and [34, Corollary 13.5.1], respectively. In [34], it is shown that (4.1a) is a consequence of (4.1b). In section 8 we provide a different proof of Lemma 4.1, where it is shown that (4.1b) follows from (4.1a). The arguments provided in the proof are instructive for later computations.

The conjugate  $f^*(y, u, \mu)$  of the perturbation function  $f(x, b, \tau)$  defined in  $\mathcal{P}$  is now easily computed:

$$\begin{aligned} f^*(y, u, \mu) &= \sup_{x, b, \tau} [\langle y, x \rangle + \langle u, b \rangle + \mu\tau - \rho(b - Ax) - \delta((x, \tau) \mid \text{epi } \phi)] \\ &= \sup_{x, w, \tau} [\langle y, x \rangle + \langle u, w + Ax \rangle + \mu\tau - \rho(w) - \delta((x, \tau) \mid \text{epi } \phi)] \\ &= \sup_{x, \tau} [\langle y + A^T u, x \rangle + \mu\tau - \delta((x, \tau) \mid \text{epi } \phi)] + \sup_w [\langle u, w \rangle - \rho(w)] \\ (4.2) \quad &= (\phi^*)^\pi(y + A^T u, -\mu) + \rho^*(u), \end{aligned}$$

where the final equality follows from (4.1a). With this representation of the conjugate of  $f$ , we obtain the following equivalent representations for the dual problem  $\mathcal{D}$ . The representation labeled  $\mathcal{D}_r$  is of particular importance to our discussion. We refer to  $\mathcal{D}_r$  as the *reduced dual*.

THEOREM 4.2 (Dual representations). *For problem  $\mathcal{P}$  define the functions*

$$\begin{aligned} g_\tau(u) &:= \rho^*(u) + \delta^*(A^T u \mid \text{lev}_\phi(\tau)), \\ p_\tau(s, \mu) &:= \tau\mu + (\phi^*)^\pi(s, \mu). \end{aligned}$$

*Then the value function for  $\mathcal{D}$  has the following equivalent characterizations:*

$$\begin{aligned} \hat{v}(b, \tau) &= \sup_u \left[ \langle b, u \rangle - \rho^*(u) - \inf_{\mu \geq 0} p_\tau(A^T u, \mu) \right] \\ (\mathcal{D}_r) \quad &= \sup_u [\langle b, u \rangle - \rho^*(u) - \delta^*(A^T u \mid \text{lev}_\phi(\tau))] \end{aligned}$$

$$(4.3a) \quad = g_\tau^*(b)$$

$$(4.3b) \quad = \text{cl}(v(\cdot, \tau))(b),$$

where the closure operation in the (4.3b) refers to the lower semicontinuous hull of the convex function  $b \mapsto v(b, \tau)$ . In particular, this implies the weak duality inequality  $\hat{v}(b, \tau) \leq v(b, \tau)$ . Moreover, if the function  $\rho$  is differentiable, the solution  $u$  to  $\mathcal{D}_r$  is unique.

In the large-scale setting, the primal problem  $\mathcal{P}(b, \tau)$  is usually solved using a primal method that does not give direct access to the multiplier  $\bar{\mu}$  for the inequality constraint  $\phi(x) \leq \tau$ . For example,  $\mathcal{P}(b, \tau)$  may be solved using a variant of the

gradient projection algorithm. However, one can still obtain an approximation to the optimal dual variable  $\bar{u}$  in  $\mathcal{D}_r$ , typically through the residual corresponding to the current iterate. For this reason, one needs a way to obtain an approximation to  $\bar{\mu}$  from an approximation to  $\bar{u}$  (i.e., given  $\bar{u}$ , compute  $\bar{\mu}$ ). Lemma 4.1 and Theorem 4.2 show that this can be done by solving the problem  $\inf_{\mu \geq 0} p_\tau(A^T \bar{u}, \mu)$  for  $\bar{\mu}$ . Indeed, in the sequel we show that in many important cases there is a closed form expression for the solution  $\bar{\mu}$ . The following lemma serves to establish a precise relationship between the solution  $\bar{u}$  to the reduced dual  $\mathcal{D}_r$  and the solution pair  $(\bar{u}, \bar{\mu})$  to the dual  $\mathcal{D}$ .

LEMMA 4.3. *Let  $\phi$  be as in  $\mathcal{P}$  with  $\tau > \inf \phi$  and  $\bar{x} \in \text{lev}_\phi(\tau)$ .*

1. *For every  $s$ , we have*

$$(4.4) \quad \delta^*(s \mid \text{lev}_\phi(\tau)) \leq \inf_{\mu \geq 0} p_\tau(s, \mu).$$

2. *Let  $(\bar{x}, \bar{s})$  satisfy  $\bar{s} \in N(\bar{x} \mid \text{lev}_\phi(\tau))$  and define*

$$S_1 = \arg \min_{\mu \geq 0} p_\tau(\bar{s}, \mu) \quad \text{and} \quad S_2 = \left\{ \bar{\mu} \geq 0 \left| \begin{array}{l} \bar{s} \in \bar{\mu}^+ \partial \phi(\bar{x}) \\ 0 = \bar{\mu}(\phi(\bar{x}) - \tau) \end{array} \right. \right\},$$

where, for  $x \in \text{dom } \phi$ ,

$$\mu^+ \partial \phi(x) := \begin{cases} \{\mu z \mid z \in \partial \phi(x)\} & \text{if } \mu > 0 \text{ and } x \in \text{dom } \partial \phi, \\ N(x \mid \text{dom } \phi) & \text{if } \mu = 0 \text{ or } \partial \phi(x) = \emptyset. \end{cases}$$

*If either  $S_1$  or  $S_2$  is nonempty, then  $S_1 = S_2$  and equality holds in (4.4).*

In Zălinescu the object  $\mu^+ \partial \phi(x)$  is denoted as  $\partial(\mu\phi)(x)$  [45, p. 141], where

$$(\mu\phi)(x) := \begin{cases} \mu\phi(x) & \text{if } \lambda > 0 \text{ and} \\ \delta(x \mid \text{dom } \phi) & \text{if } \lambda = 0. \end{cases}$$

We choose the notation  $\mu^+ \partial \phi(x)$  to emphasize that there is an underlying limiting operation at play, e.g., see [36, Definition 8.3 and Proposition 8.12].

The final lemma of this section concerns conditions under which solutions to  $\mathcal{P}$  and  $\mathcal{D}_r$  exist. This is closely tied to the horizon behavior of these problems and the notion of coercivity.

DEFINITION 4.4. *A function  $h : \mathbb{R}^n \rightarrow \overline{\mathbb{R}}$  is said to be  $\alpha$ -coercive if*

$$\lim_{\|x\| \rightarrow \infty} \frac{f(x)}{\|x\|^\alpha} = +\infty.$$

*In particular,  $h$  is said to be 0-coercive, or simply coercive, if  $\lim_{\|x\| \rightarrow \infty} f(x) = \infty$ .*

LEMMA 4.5 (Coercivity of primal and dual objectives).

1. *The objective function  $f(\cdot, b, \tau)$  of  $\mathcal{P}$  is coercive if and only if*

$$(4.5a) \quad \text{hzn}(\phi) \cap [-A^{-1} \text{hzn}(\rho)] = \{0\}.$$

2. *The objective function of the reduced dual  $\mathcal{D}_r$  is coercive if and only if*

$$(4.5b) \quad b \in \text{int}(\text{dom } \rho + \text{Alev}_\phi(\tau)).$$

**5. Variational properties of the value function.** Using  $\mathcal{D}$  and representation of the conjugate of the objective of  $\mathcal{P}$  (cf. (4.2)), we can specialize [36, Theorem 11.39] to obtain a characterization of the subdifferential of the value function, as well as sufficient conditions for strong duality.

**THEOREM 5.1** (Strong duality and subgradient of the value function). *Let  $v$  and  $\hat{v}$  be as in  $\mathcal{P}$  and  $\mathcal{D}$ , respectively. It is always the case that*

$$v(b, \tau) \geq \hat{v}(b, \tau) \quad (\text{weak duality}).$$

If  $(b, \tau) \in \text{int}(\text{dom } v)$ , then

$$v(b, \tau) = \hat{v}(b, \tau) \quad (\text{strong duality})$$

and  $\partial v(b, \tau) \neq \emptyset$  with

$$\partial v(b, \tau) := \arg \max_{u, \mu \geq 0} [\langle b, u \rangle - \rho^*(u) - p_\tau(A^T u, -\mu)].$$

Furthermore, for fixed  $(b, \tau) \in \mathbb{R}^m \times \mathbb{R}$ ,

$$\text{dom } f(\cdot, b, \tau) \neq \emptyset \iff b \in \text{dom } \rho + A(\text{lev}_\phi(\tau)).$$

In particular, this implies that

$$(b, \tau) \in \text{int}(\text{dom } v) \iff b \in \text{int}(\text{dom } \rho + A(\text{lev}_\phi(\tau))).$$

We now derive a characterization of the subdifferential  $\partial v(b, \tau)$  based on the solutions of the reduced dual  $\mathcal{D}_\tau$ .

**THEOREM 5.2** (Value function subdifferential). *Suppose that*

$$(5.1a) \quad b \in \text{ri}(\text{dom } \rho) + A \text{ri}(\text{lev}_\phi(\tau)) \quad \text{and}$$

$$(5.1b) \quad \text{ri}(\text{dom } \rho^*) \cap [A^{-T} \text{ri}(\text{bar}(\text{lev}_\phi(\tau)))] \neq \emptyset.$$

1. If the pair  $(\bar{x}, \bar{u})$  satisfies

$$(5.1c) \quad \bar{x} \in \text{lev}_\phi(\tau), \quad \bar{u} \in \partial \rho(b - A\bar{x}), \quad \text{and} \quad A^T \bar{u} \in N(\bar{x} | \text{lev}_\phi(\tau)),$$

then  $\bar{x}$  solves  $\mathcal{P}$  and  $\bar{u}$  solves  $\mathcal{D}_\tau$ .

2. If  $\bar{x}$  solves  $\mathcal{P}$  and (5.1a) holds, there exists  $\bar{u}$  such that  $(\bar{x}, \bar{u})$  satisfies (5.1c).

3. If  $\bar{u}$  solves  $\mathcal{D}_\tau$  and (5.1b) holds, there exists  $\bar{x}$  such that  $(\bar{x}, \bar{u})$  satisfies (5.1c).

4. If either (4.5a) and (5.1a) holds, or (4.5b) and (5.1b) holds, then  $\partial v(b, \tau) \neq \emptyset$  and  $\arg \min_{\mu \geq 0} p_\tau(A^T \bar{u}, \mu) \neq \emptyset$  for all  $(\bar{x}, \bar{u}) \in \mathbb{R}^n \times \mathbb{R}^m$  satisfying (5.1c) with

$$(5.1d) \quad \partial v(b, \tau) = \left\{ \begin{pmatrix} \bar{u} \\ -\bar{\mu} \end{pmatrix} \mid \begin{array}{l} (\bar{x}, \bar{u}) \in \mathbb{R}^n \times \mathbb{R}^m \text{ satisfy (5.1c) and} \\ \bar{\mu} \in \arg \min_{\mu \geq 0} p_\tau(A^T \bar{u}, \mu) \end{array} \right\}$$

$$(5.1e) \quad = \left\{ \begin{pmatrix} \bar{u} \\ -\bar{\mu} \end{pmatrix} \mid \begin{array}{l} \exists \bar{x} \in \text{lev}_\phi(\tau) \text{ s.t. } 0 \in -A^T \bar{u} + \bar{\mu}^+ \partial \phi(\bar{x}), \\ \text{where } \bar{u} \in \partial \rho(b - A\bar{x}), \\ \bar{\mu} \geq 0, \text{ and } \bar{\mu}(\phi(\bar{x}) - \tau) = 0 \end{array} \right\}.$$

The representation (5.1e) expresses the elements of  $\partial v(b, \tau)$  in terms of classical Lagrange multipliers when  $\bar{\mu} > 0$ , and extends the classical theory when  $\bar{\mu} = 0$ . (See Lemma 4.3 for the definition of  $\mu^+ \partial \phi(x)$ .) Because  $v$  is convex, it is subdifferentially regular, and so for fixed  $b$ , we can obtain the subdifferential of  $v$  with respect to  $\tau$  alone [36, Corollary 10.11], i.e.,

$$\partial_\tau v(b, \tau) = \left\{ \omega \left| \begin{pmatrix} u \\ \omega \end{pmatrix} \in \partial v(b, \tau) \right. \right\}.$$

**6. Applications.** In this section we apply the calculus rules of section 3.3 in conjunction with Theorem 5.2 to evaluate the subdifferential of the value function in three important special cases: where  $\phi$  is a gauge-plus-indicator (section 6.1), a quadratic support function (section 6.2), and an affine composition with a quadratic support function (section 6.3). In all cases we allow  $\rho$  to be an arbitrary convex function.

**6.1. Gauge-plus-indicator.** The case where  $\rho$  is a linear least-squares objective and  $\phi$  is a gauge function is studied in [12]. We generalize this case by allowing the convex function  $\rho$  to be possibly nonsmooth and non-finite-valued, and take

$$(6.1) \quad \phi(x) := \gamma(x | U) + \delta(x | X),$$

where  $U$  is a nonempty closed convex set containing the origin. Here,  $\gamma(x | U)$  is the gauge function defined in (1.1). It is evident from the definition of a gauge that  $\phi$  is also a gauge if and only if  $X$  is a convex cone. Since  $0 \in U$ , it follows from [34, Theorem 14.5] that  $\gamma(\cdot | U) = \delta^*(\cdot | U^\circ)$ , where

$$U^\circ = \{v | \langle v, u \rangle \leq 1 \ \forall u \in U\}$$

is the polar of the set  $U$ .

Observe that the requirement  $x \in X$  is unaffected by varying  $\tau$  in the constraint  $\phi(x) \leq \tau$ . Indeed, the problem  $\mathcal{P}$  is unchanged if we replace  $\rho$  and  $\phi$  by

$$(6.2) \quad \hat{\rho}(y, x) := \rho(y) + \delta(x | X) \quad \text{and} \quad \hat{\phi}(x) := \gamma(x | U)$$

with  $A$  and  $b$  replaced by

$$\hat{b} := \begin{pmatrix} b \\ 0 \end{pmatrix} \quad \text{and} \quad \hat{A} := \begin{bmatrix} A \\ -I \end{bmatrix}.$$

Hence, the generalization of [12] discussed here only concerns the application to more general convex functions  $\rho$ .

There are two ways one can proceed with this application. One can use  $\phi$  as given in (6.1) or use  $\hat{\rho}$  and  $\hat{\phi}$  as defined in (6.2). We choose the former in order to highlight the presence of the abstract constraint  $x \in X$ . But we emphasize—regardless of the formulation chosen—that the end result is the same.

LEMMA 6.1. *Let  $\phi$  be as given in (6.1). The following formulas hold:*

$$(6.3a) \quad \gamma(\cdot | U) = \delta^*(\cdot | U^\circ),$$

$$(6.3b) \quad \text{dom } \gamma(\cdot | U) = \text{cone}(U) = \text{bar}(U^\circ),$$

$$(6.3c) \quad \text{dom } \phi = \text{cone}(U) \cap X,$$

$$(6.3d) \quad \text{lev}_\phi(\tau) = (\tau U) \cap X,$$

$$(6.3e) \quad \text{hzn}(\phi) = U^\infty \cap X^\infty, \text{ and}$$

$$(6.3f) \quad \text{cl}(\text{bar}(\text{lev}_\phi(\tau))) = \text{cl}(\text{bar}(U) + \text{bar}(X)).$$

If it is further assumed that

$$(6.4) \quad \text{ri}(\tau U) \cap \text{ri}(X) \neq \emptyset,$$

and then we also have

$$(6.5a) \quad \phi^*(z) = \min_s [\delta^*(z - s | X) + \delta(s | U^\circ)],$$

$$(6.5b) \quad (\phi^*)^\pi(z, \mu) = \min_s [\delta^*(z - s | X) + \delta(s | \mu U^\circ)],$$

$$(6.5c) \quad \delta^*(z | \text{lev}_\phi(\tau)) = \min_{\mu \geq 0} [\tau\mu + (\phi^*)^\pi(z, \mu)]$$

$$(6.5d) \quad = \min_s [\delta^*(z - s | X) + \tau\gamma(s | U^\circ)], \text{ and}$$

$$(6.5e) \quad N(x | \text{lev}_\phi(\tau)) = N(x | \tau U) + N(x | X).$$

If  $\bar{s}$  minimizes (6.5d), then  $\bar{\mu} := \gamma(\bar{s} | U^\circ)$  minimizes (6.5c).

By Theorem 5.1, the subdifferential of  $v(b, \tau)$  is obtained by solving the dual problem (8.4) or the reduced dual  $\mathcal{D}_r$ . When  $\phi$  is given by (6.1), the results of Lemma 6.1 show that the dual and the reduced dual take the form

$$(6.6) \quad \begin{aligned} \sup_{u, \mu} [\langle b, u \rangle + \tau\mu - (\phi^*)^\pi(A^T u, -\mu) - \rho^*(u)] \\ = \sup_u [\langle b, u \rangle - \rho^*(u) - \delta^*(A^T u | \text{lev}_\phi(\tau))] \\ = \sup_u [\langle b, u \rangle - \rho^*(u) - \min_s [\delta^*(A^T u - s | X) + \tau\gamma(s | U^\circ)]] \\ (6.7) \quad = \sup_{u, s} [\langle b, u \rangle - \rho^*(u) - \delta^*(A^T u - s | X) - \delta^*(s | \tau U)]. \end{aligned}$$

Moreover, if  $(\bar{u}, \bar{s})$  solves (6.7), then  $(\bar{u}, \bar{\mu})$  solves (6.6) with  $\bar{\mu} = -\gamma(\bar{s} | U^\circ)$ , and

$$(\bar{u}, -\gamma(\bar{s} | U^\circ)) \in \partial v(b, \tau).$$

We have the following version of Theorem 5.2 when  $\phi$  is given by (6.1).

**THEOREM 6.2.** *Let  $\phi$  be given by (6.1) under the assumption that (6.4) holds, and consider the following two conditions:*

$$(6.8) \quad b \in \text{ri}(\text{dom } \rho + A[\tau U \cap X]) = \text{ri}(\text{dom } \rho) + A[\text{ri}(\tau U) \cap \text{ri}(X)]$$

and

$$(6.9) \quad \exists \hat{u} \in \text{ri}(\text{dom } \rho^*) \text{ such that } A^T \hat{u} \in \text{ri}(\text{bar}(U)) + \text{ri}(\text{bar}(X)).$$

1. *If the triple  $(\bar{x}, \bar{u}, \bar{s})$  satisfies*

$$(6.10a) \quad \bar{u} \in \partial \rho(b - A\bar{x}), \quad \bar{x} \in X \cap (\tau U),$$

$$(6.10b) \quad \bar{s} \in N(\bar{x} | \tau U), \quad \text{and} \quad A^T \bar{u} - \bar{s} \in N(\bar{x} | X),$$

*then  $\bar{x}$  solves  $\mathcal{P}(b, \tau)$  and  $(\bar{u}, \bar{s})$  solves (6.7).*

2. *If  $\bar{x}$  solves  $\mathcal{P}(b, \tau)$  and (6.8) holds, then there exists a pair  $(\bar{u}, \bar{s})$  such that  $(\bar{x}, \bar{u}, \bar{s})$  satisfies (6.10).*

3. *If  $(\bar{u}, \bar{s})$  solves (6.7) and (6.9) holds, then there exists  $\bar{x}$  such that  $(\bar{x}, \bar{u}, \bar{s})$  satisfies (6.10).*

4. If either

$$(6.11) \quad U^\infty \cap X^\infty \cap [-A^{-1}\text{hzn}(\rho)] = \{0\} \quad \text{and (6.8) holds}$$

or

$$(6.12) \quad b \in \text{int}(\text{dom } \rho + A[\tau U \cap X]) \text{ and (6.9) holds,}$$

then  $\partial v(b, \tau) \neq \emptyset$  and is given by

$$(6.13) \quad \begin{aligned} \partial v(b, \tau) &= \left\{ \left( \begin{array}{c} \bar{u} \\ -\gamma(\bar{s} \mid U^\circ) \end{array} \right) \middle| (\bar{x}, \bar{u}, \bar{s}) \in \mathbb{R}^n \times \mathbb{R}^m \times \mathbb{R}^n \text{ satisfy (6.10)} \right\} \\ &= \left\{ \left( \begin{array}{c} \bar{u} \\ -\bar{\mu} \end{array} \right) \middle| \begin{array}{l} \exists \bar{x} \in X \cap (\tau U) \text{ s.t.} \\ 0 \in -A^T \bar{u} + N(\bar{x} \mid X) + \bar{\mu}^+ \partial \gamma(\bar{x} \mid U), \text{ where} \\ \bar{u} \in \partial \rho(b - A\bar{x}), 0 \leq \bar{\mu} \text{ and } \bar{\mu}(\gamma(\bar{x} \mid U) - \tau) = 0 \end{array} \right\}. \end{aligned}$$

**6.1.1. Gauge penalties.** In [12], the authors study the case where  $\rho$  is a linear least-squares objective,  $\phi$  is a gauge functional, and  $X = \mathbb{R}^n$ . In this case, [12, Lemma 2.1] and [12, Theorem 2.2(b)] can be deduced from (6.7) and (6.13), respectively. Another application is to the case where  $\rho$  is finite-valued and smooth,  $\phi$  is a norm, and  $X$  is a generalized box. In this case, all of the conditions of Theorems 5.1 and 6.2 are satisfied, solutions to both  $\mathcal{P}(b, \tau)$  and (6.7) exist, and  $v$  is differentiable. In particular, consider the nonnegative 1-norm-constrained inversion, where

$$\phi(x) = \|x\|_1 + \delta(x \mid \mathbb{R}_+^n),$$

and  $\rho$  is any differentiable convex function. The subdifferential characterization given in Theorem 5.1 can be explicitly computed via Theorem 6.2. In the notation of (6.1),

$$U = \{x \mid \|x\|_1 \leq 1\} : \mathbb{B}_1,$$

and  $X$  in (6.1) is  $\mathbb{R}_+^n$ . Since the function  $\rho$  is differentiable, the solution  $\bar{u}$  to the dual (6.7) is unique [34, Theorem 26.3]. Therefore, Theorem 6.2 gives the existence of a *unique* gradient

$$\nabla_b v(b, \tau) = -A^T \nabla \rho(b - A\bar{x}),$$

where  $\bar{x}$  is any solution that achieves the optimal value. The derivative with respect to  $\tau$  is immediately given by Theorem 6.2 as

$$(6.14) \quad \nabla_\tau v(b, \tau) = -\gamma(A^T \nabla \rho(b - A\bar{x}) \mid U^\circ) = -\|A^T \nabla \rho(b - A\bar{x})\|_\infty.$$

Note that (6.14) has the same algebraic form when  $x$  is unconstrained. The nonnegativity constraint on  $x$  is reflected in the derivative only through its effect on the optimal point  $\bar{x}$ .

**6.2. Quadratic support functions.** We now consider the case

$$(6.15) \quad \phi(x) := \sup_{w \in U} [\langle x, w \rangle - \frac{1}{2} \langle w, Bw \rangle],$$

where  $U \subset \mathbb{R}^n$  is nonempty, closed, and convex with  $0 \in U$ , and  $B \in \mathbb{R}^{n \times n}$  is positive semidefinite. We call this class of functions quadratic support (QS) functions. This surprisingly rich and useful class is found in many applications. A deeper study of its properties and uses can be found in [4]. Note that the conjugate of  $\phi$  is given by

$$(6.16) \quad \phi^*(w) = \frac{1}{2} \langle w, Bw \rangle + \delta(w | U).$$

If the set  $U$  is polyhedral convex, then the function  $\phi$  is called a piecewise linear-quadratic (PLQ) penalty function [36, Example 11.18]. Since  $B$  is positive semidefinite there is a matrix  $L \in \mathbb{R}^{n \times k}$  such that  $B = LL^T$ , where  $k$  is the rank of  $B$ . Using  $L$ , the calculus rules in section 3.3 give the following alternative representation for  $\phi$ :

$$(6.17) \quad \begin{aligned} \phi(x) &= \sup_{w \in U} [\langle w, x \rangle - \frac{1}{2} \|L^T w\|_2^2 - \delta(w | U)] \\ &= \inf_{x_1 + x_2 = x} \left[ \delta^*(x_1 | U) + \inf_{Ls = x_2} \frac{1}{2} \|s\|_2^2 \right] \\ &= \inf_s \left[ \frac{1}{2} \|s\|_2^2 + \delta^*(x - Ls | U) \right] \\ &= \inf_s \left[ \frac{1}{2} \|s\|_2^2 + \gamma(x - Ls | U^\circ) \right], \end{aligned}$$

where the final equality follows from [34, Theorem 14.5] since  $0 \in U$ . Note that the function class (6.15) includes all gauge functionals for sets containing the origin. By (6.16), it easily follows that

$$(\phi^*)^\pi(w, \mu) = \begin{cases} \frac{1}{2\mu} \|w\|_B^2 + \delta(w | \mu U) & \text{if } \mu > 0, \\ \delta(w | U^\infty \cap \text{Nul}(B)) & \text{if } \mu = 0, \\ +\infty & \text{if } \mu < 0, \end{cases}$$

where  $\|\cdot\|_B$  denotes the seminorm induced by  $B$ , i.e.,

$$\|w\|_B := \sqrt{w^T B w}.$$

The next result catalogues important properties of the function  $\phi$  given in (6.15).

LEMMA 6.3. *Let  $\phi$  be given by (6.15) with  $\tau > 0$ . Then*

$$\begin{aligned} \text{dom } \phi &= \text{cone}(U^\circ) + \text{Ran}(B) \quad \text{and} \\ \text{hzn}(\phi) &= \text{cone}(U)^\circ, \end{aligned}$$

and in particular,  $\phi$  is coercive if and only if  $0 \in \text{int}(U)$ . Moreover,

$$(6.18) \quad \delta^*(w | \text{lev}_\phi(\tau)) = \min_{\lambda \geq 0} [\tau\lambda + (\phi^*)^\pi(w, \lambda)]$$

$$(6.19) \quad = \begin{cases} \tau\gamma(w | U) + \frac{\|w\|_B^2}{2\gamma(w | U)} & \text{if } \gamma(w | U) > \|w\|_B / \sqrt{2\tau}, \\ \sqrt{2\tau} \|w\|_B & \text{if } \gamma(w | U) \leq \|w\|_B / \sqrt{2\tau}, \end{cases}$$

where the minimizing  $\lambda$  in (6.18) is given by

$$(6.20) \quad \lambda = \max \left\{ \gamma(w | U), \frac{\|w\|_B}{\sqrt{2\tau}} \right\}.$$

In particular, the formula (6.19) implies that

$$\text{bar}(\text{lev}_\phi(\tau)) = \text{dom}(\delta^*(\cdot | \text{lev}_\phi(\tau))) = \text{dom}(\gamma(\cdot | U)) = \text{cone}(U).$$

We now apply Theorem 5.2 to the case where  $\phi$  is given by (6.15).

THEOREM 6.4. *Let  $\phi$  be given by (6.15) and consider the following two conditions:*

$$(6.21) \quad \exists \hat{x} \in \text{ri}(\text{dom } \phi) \quad \text{such that} \quad \phi(\hat{x}) < \tau \quad \text{and} \quad b - A\hat{x} \in \text{ri}(\text{dom } \rho)$$

and

$$(6.22) \quad \exists \hat{u} \in \text{ri}(\text{dom } \rho^*) \quad \text{such that} \quad A^T \hat{u} \in \text{ri}(\text{cone}(U)).$$

1. *If the pair  $(\bar{x}, \bar{u})$  satisfy*

$$(6.23) \quad \bar{x} \in \text{lev}_\phi(\tau), \quad \bar{u} \in \partial\rho(b - A\bar{x}), \quad \text{and} \quad A^T \bar{u} \in N(\bar{x} | \text{lev}_\phi(\tau)),$$

*then  $\bar{x}$  solves  $\mathcal{P}(b, \tau)$  and  $\bar{u}$  solves  $\mathcal{D}_r$ .*

2. *If  $\bar{x}$  solves  $\mathcal{P}(b, \tau)$  and (6.21) holds, then there exists  $\bar{u}$  such that (6.23) holds.*

3. *If  $\bar{u}$  solves  $\mathcal{D}_r$  and (6.22) holds, then there exists  $\bar{x}$  such that (6.23) holds.*

4. *If either*

$$(6.24) \quad \text{cone}(U)^\circ \cap [-A^{-1}\text{hzn}(\rho)] = \{0\} \quad \text{and} \quad (6.21) \quad \text{holds}$$

*or*

$$(6.25) \quad b \in \text{int}(\text{dom } \rho + \text{Alev}_\phi(\tau)) \quad \text{and} \quad (6.22) \quad \text{holds,}$$

*then  $\partial v(b, \tau) \neq \emptyset$  and is given by*

$$(6.26) \quad \partial v(b, \tau) = \left\{ \left( \begin{array}{l} \bar{u} \\ -\bar{\mu} \end{array} \right) \mid \begin{array}{l} \exists \bar{x} \text{ s.t. } (\bar{x}, \bar{u}) \text{ satisfy (6.23) and} \\ \bar{\mu} = \max \{ \gamma(A^T \bar{u} | U), \|A^T \bar{u}\|_B / \sqrt{2\tau} \} \end{array} \right\} \\ = \left\{ \left( \begin{array}{l} \bar{u} \\ -\bar{\mu} \end{array} \right) \mid \begin{array}{l} \exists \bar{x} \in \text{lev}_\phi(\tau) \text{ s.t. } 0 \in -A^T \bar{u} + \bar{\mu}^+ \partial\phi(\bar{x}), \text{ where} \\ \bar{u} \in \partial\rho(b - A\bar{x}), \quad 0 \leq \bar{\mu}, \text{ and } \bar{\mu}(\phi(\bar{x}) - \tau) = 0 \end{array} \right\}.$$

In the following corollary we exploit the structure of  $\phi$  to refine the multiplier description of the  $\partial v(b, \tau)$  given in (6.26).

COROLLARY 6.5. *Consider the problem  $\mathcal{P}(b, \tau)$  with  $\phi$  given by (6.15). A pair  $(\bar{x}, \bar{u})$  satisfies (6.23) if and only if  $\bar{x} \in \text{lev}_\phi(\tau)$ ,  $\bar{u} \in \partial\rho(b - A\bar{x})$ , and either*

$$(6.27a) \quad A^T \bar{u} \in N(\bar{x} | \text{dom } \phi) \quad \text{or}$$

$$(6.27b) \quad \exists \bar{\mu} > 0, \quad \bar{w} \in U \quad \text{such that} \quad \bar{x} \in B\bar{w} + N(\bar{w} | U) \quad \text{and} \quad A^T \bar{u} = \bar{\mu}\bar{w}.$$

**6.2.1. Huber penalty.** A popular function in the PLQ class is the Huber penalty [29]:

$$\phi(x) = \sup_{w \in [-\kappa, \kappa]^n} \left[ \langle x, w \rangle - \frac{1}{2} \|w\|_2^2 \right] = \sum_{i=1}^n \phi_i(x_i); \quad \phi_i(x_i) := \begin{cases} \frac{1}{2} x_i^2 & \text{if } |x_i| \leq \kappa, \\ \kappa |x_i| - \kappa^2/2 & \text{otherwise.} \end{cases}$$

The Huber function is of form (6.15) with  $B = I$  and  $U = [-\kappa, \kappa]^n$ . In this case,  $U^\infty \cap \text{Nul}(B) = \{0\}$  so that the conditions of Corollary 6.5 hold.

A graph of the scalar component function  $\phi_i$  is shown in Figure 6.1. The Huber penalty is robust to outliers, since it increases linearly rather than quadratically outside the threshold defined by  $\kappa$ . For any misfit function  $\rho$ , Theorem 6.4 can be



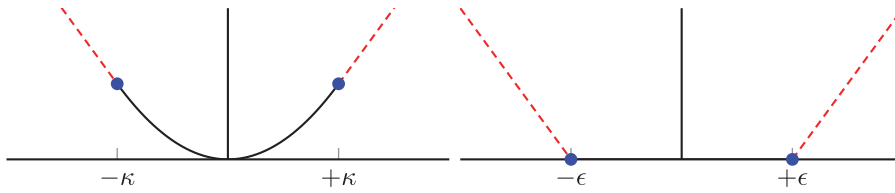


FIG. 6.1. Huber (left) and Vapnik (right) penalties.

used to easily compute the subgradient  $\partial v(b, \tau)$  of the value function. If the regularity condition (6.21) is satisfied (e.g., if  $\rho$  is finite valued), then Theorem 6.4 implies that

$$\partial v(b, \tau) = \left\{ \begin{pmatrix} \bar{u} \\ -\bar{\mu} \end{pmatrix} \mid \begin{matrix} (\bar{x}, \bar{u}) \text{ satisfy (6.23) and} \\ \bar{\mu} = \max \{ \kappa \|A^T \bar{u}\|_\infty, \|A^T \bar{u}\|_2 / \sqrt{2\tau} \} \end{matrix} \right\}.$$

In particular, if  $\rho$  is differentiable finite-valued,  $\bar{u} = \nabla \rho(b - A\bar{x})$  is unique and

$$\nabla v(b, \tau) = \begin{pmatrix} \bar{u} \\ -\bar{\mu} \end{pmatrix}.$$

**6.3. Affine composition with QS functions.** Next consider the case where  $\phi$  takes the form

$$(6.28) \quad \phi(x) := \psi(Hx + c), \quad \text{where} \quad \psi(y) := \sup_{w \in U} [ \langle y, w \rangle - \frac{1}{2} \langle w, Bw \rangle ],$$

$H \in \mathbb{R}^{\nu \times n}$  is injective,  $c \in \mathbb{R}^\nu$ ,  $U \subset \mathbb{R}^\nu$  is nonempty, closed, and convex with  $0 \in U$ , and  $B \in \mathbb{R}^{\nu \times \nu}$  is symmetric and positive semidefinite. We assume that

$$\exists \hat{x} \text{ such that } H\hat{x} + c \in \text{ri}(\text{dom } \psi),$$

where  $\text{dom } \psi = \text{cone}(U^\circ) + \text{Ran}(B)$  (Lemma 6.3). We show that the function  $\phi$  in (6.28) is an instance of the quadratic support functions considered in section 6.2. To see this we make the following definitions:

$$\tilde{x} = \begin{pmatrix} x \\ s \end{pmatrix}, \quad \tilde{y} = \begin{pmatrix} y \\ z \end{pmatrix}, \quad \tilde{w} = \begin{pmatrix} v \\ w \end{pmatrix}, \quad \tilde{U} = \{0\} \times U,$$

$$\tilde{b} = \begin{pmatrix} b \\ c \end{pmatrix}, \quad \tilde{A} = \begin{bmatrix} A & 0 \\ -H & I \end{bmatrix}, \quad \tilde{B} = \begin{bmatrix} 0 & 0 \\ 0 & B \end{bmatrix}, \quad \tilde{\rho} \begin{pmatrix} y \\ z \end{pmatrix} = \rho(y) + \delta(z \mid \{0\}), \quad \text{and}$$

$$\tilde{\phi} \begin{pmatrix} x \\ s \end{pmatrix} = \sup_{\begin{pmatrix} v \\ w \end{pmatrix} \in \tilde{U}} \left[ \left\langle \begin{pmatrix} v \\ w \end{pmatrix}, \begin{pmatrix} x \\ s \end{pmatrix} \right\rangle - \frac{1}{2} \left\langle \begin{pmatrix} v \\ w \end{pmatrix}, \tilde{B} \begin{pmatrix} v \\ w \end{pmatrix} \right\rangle \right] = \delta^*(x \mid \{0\}) + \psi(s).$$

With these definitions, the two problems  $\mathcal{P}(b, \tau)$  and

$$\text{minimize} \quad \tilde{\rho}(\tilde{b} - \tilde{A}\tilde{x}) \quad \text{subject to} \quad \tilde{\phi}(\tilde{x}) \leq \tau$$

are equivalent. In addition, we have the relationships

$$\tilde{\rho}^* \begin{pmatrix} u \\ r \end{pmatrix} = \rho^*(u) + \delta^*(r \mid \{0\}), \quad \tilde{\phi}^* \begin{pmatrix} v \\ w \end{pmatrix} = \delta(v \mid \{0\}) + \psi^*(w),$$

$$\gamma \left( \begin{pmatrix} v \\ w \end{pmatrix} \mid \tilde{U} \right) = \delta(v \mid \{0\}) + \gamma(w \mid U), \quad \text{and} \quad \left\| \begin{pmatrix} v \\ w \end{pmatrix} \right\|_{\tilde{B}} = \delta^*(v \mid \{0\}) + \|w\|_B.$$

Moreover, the reduced dual  $\mathcal{D}_r$  becomes

$$(6.29) \quad \sup_{H^T r = A^T u} [\langle b, u \rangle + \langle c, r \rangle - \rho^*(u) - \delta^*(r | \text{lev}_\psi(\tau))].$$

Using standard methods of convex analysis, we obtain the following result as a direct consequence of Theorem 6.4 and [36, Corollary 10.11].

**THEOREM 6.6.** *Let  $\phi$  be given by (6.28), and consider the following two conditions:*

$$(6.30) \quad \exists \hat{x} \text{ such that } H\hat{x} + c \in \text{ri}(\text{dom } \psi), \psi(H\hat{x} + c) < \tau, \text{ and } b - A\hat{x} \in \text{ri}(\text{dom } \rho)$$

and

$$(6.31) \quad \exists \hat{u} \in \text{ri}(\text{dom } \rho^*) \text{ and } \hat{r} \in \text{ri}(\text{cone}(U)) \text{ such that } \begin{pmatrix} \hat{u} \\ \hat{r} \end{pmatrix} \in \text{Nul} \left( \begin{bmatrix} A \\ -H \end{bmatrix}^T \right).$$

1. If the triple  $(\bar{x}, \bar{u}, \bar{r})$  satisfies

$$(6.32)$$

$$\bar{x} \in \text{lev}_\phi(\tau), \bar{u} \in \partial\rho(b - A\bar{x}), \bar{r} \in N(H\bar{x} + c | \text{lev}_\psi(\tau)), \text{ and } A^T \bar{u} = H^T \bar{r},$$

then  $\bar{x}$  solves  $\mathcal{P}(b, \tau)$  and  $(\bar{u}, \bar{r})$  solves (6.29).

2. If  $\bar{x}$  solves  $\mathcal{P}(b, \tau)$  and (6.30) holds, there exists  $(\bar{u}, \bar{r})$  such that (6.32) holds.
3. If  $(\bar{u}, \bar{r})$  solves (6.29) and (6.31) holds, there exists  $\bar{x}$  such that (6.32) holds.
4. If either

$$H^{-1}[\text{cone}(U)^\circ] \cap [-A^{-1} \text{hzn}(\rho)] = \{0\} \text{ and (6.30) holds}$$

or

$$\begin{pmatrix} b \\ c \end{pmatrix} \in \text{int} \left( \text{dom } \rho \times \text{lev}_\psi(\tau) + \text{Ran} \left( \begin{bmatrix} A \\ -H \end{bmatrix} \right) \right) \text{ and (6.31) holds,}$$

then  $\partial v(b, \tau) \neq \emptyset$  and is given by

$$\begin{aligned} \partial v(b, c, \tau) &= \left\{ \begin{pmatrix} \bar{u} \\ \bar{r} \\ -\bar{\mu} \end{pmatrix} \mid \begin{array}{l} \exists \bar{x} \in \mathbb{R}^n \text{ s.t. } (\bar{x}, \bar{u}, \bar{r}) \text{ satisfy (6.32) and} \\ \bar{\mu} = \max \{ \gamma(\bar{r} | U), \|\bar{r}\|_B / \sqrt{2\tau} \} \end{array} \right\} \\ &= \left\{ \begin{pmatrix} \bar{u} \\ \bar{r} \\ -\bar{\mu} \end{pmatrix} \mid \begin{array}{l} \exists \bar{x} \in \mathbb{R}^n \text{ s.t. } c + H\bar{x} \in \text{lev}_\psi(\tau), \\ \bar{u} \in \partial\rho(b - A\bar{x}), \bar{r} \in \bar{\mu}^+ \partial\psi(c + H\bar{x}), \bar{\mu} \geq 0, \\ \bar{\mu}(\psi(c + H\bar{x}) - \tau) = 0, \text{ and } A^T \bar{u} = H^T \bar{r} \end{array} \right\}. \end{aligned}$$

**COROLLARY 6.7.** *Consider the problem  $\mathcal{P}(b, \tau)$  with  $\phi$  given by (6.28). Then  $(\bar{x}, \bar{u}, \bar{r})$  satisfies (6.32) if and only if*

$$H\bar{x} + c \in \text{lev}_\psi(\tau), \bar{u} \in \partial\rho(b - A\bar{x}), A^T \bar{u} = H^T \bar{r},$$

$$\text{and either } \bar{r} \in N(H\bar{x} + c | \text{dom } \psi) \text{ or}$$

$$\exists \bar{\mu} \geq 0, \bar{w} \in U \text{ such that } Hx + c \in B\bar{w} + N(\bar{w} | U) \text{ and } \bar{r} = \bar{\mu}\bar{w}.$$

**6.3.1. Vapnik penalty.** The Vapnik penalty

$$\rho(r) = \sup_{u \in [0,1]^{2n}} \left\{ \left\langle \begin{bmatrix} r - \epsilon \\ -r - \epsilon \end{bmatrix}, u \right\rangle \right\} = (r - \epsilon)_+ + (-r - \epsilon)_+$$

is an important example in the PLQ class which is most easily represented as the composition of an affine transformation with a PLQ function. The scalar version is shown in the right panel of Figure 6.1. In this case,

$$H = \begin{bmatrix} I \\ -I \end{bmatrix}, \quad c = - \begin{bmatrix} \epsilon \mathbf{1} \\ \epsilon \mathbf{1} \end{bmatrix}, \quad B = \mathbf{0} \in \mathbb{R}^{2n \times 2n}, \quad \text{and} \quad U = [0, 1]^{2n}.$$

In order to satisfy (6.32), we need to find a triple  $(\bar{x}, \bar{u}, \bar{w})$  with  $\bar{w} = [\bar{w}_1 \ \bar{w}_2]^T \in [0, 1]^{2n}$  so that  $\bar{u} \in \partial \rho(b - A\bar{x})$  and  $A^T \bar{u} = H^T \bar{w} = \bar{w}_1 - \bar{w}_2$ . We claim that either  $\bar{w}_1(i) = 0$  or  $\bar{w}_2(i) = 0$  for all  $i$ . To see this, observe that  $\bar{w} \in N(H\bar{x} + c | \text{lev}_\psi(\tau))$ , so

$$\left\langle \bar{w}, y - \begin{bmatrix} \bar{x} - \epsilon \\ -\bar{x} - \epsilon \end{bmatrix} \right\rangle \leq 0$$

whenever  $\psi(y) \leq \tau$ . Taking  $y$  first with  $-\epsilon$  as the only nonzero in the  $i$ th coordinate, and then with  $-\epsilon$  in the only nonzero in the  $(n + i)$ th coordinate, we get

$$\bar{w}_1(i)(-\bar{x}(i)) \leq 0 \quad \text{and} \quad \bar{w}_2(i)(\bar{x}(i)) \leq 0.$$

If  $x(i) < 0$ , from the first equation we get  $\bar{w}_1(i) = 0$ , while if  $x(i) > 0$ , we get  $\bar{w}_2(i) = 0$  from the second equation. If  $x(i) = 0$ , then taking  $y = 0$  gives

$$\bar{w}_1(i)\epsilon \leq 0 \quad \text{and} \quad \bar{w}_2(i)\epsilon \leq 0,$$

so  $\bar{w}_1(i) = \bar{w}_2(i) = 0$ . Since  $A^T \bar{u} = \bar{w}_1 - \bar{w}_2$  and  $\bar{w}_1(i)$  or  $\bar{w}_2(i)$  is 0 for each  $i$ , we get  $\mu = \gamma(\bar{w} | [0, 1]^{2n}) = \|A^T \bar{u}\|_\infty$ . Hence, the subdifferential  $\partial v$  is computed in precisely the same way for the Vapnik regularization as for the 1-norm.

**7. Numerical example: Robust nonnegative basis pursuit .** In this example, we recover a nonnegative undersampled sparse signal from a set of very noisy measurements using several formulations of  $\mathcal{P}$ . We compare the performance of three different penalty functions  $\rho$ : least-squares, Huber (see section 6.2.1), and a nonconvex penalty arising from the student’s t distribution (see, e.g., [5, 3]). The regularizing function  $\phi$  in all of the examples is the sum of the 1-norm and the indicator of the positive orthant (see section 6.1.1).

The formulations using Huber and Student’s t misfits are robust alternatives to the nonnegative basis pursuit problem [18]. The Huber misfit agrees with the quadratic penalty for small residuals but is relatively insensitive to larger residuals. The student’s t misfit is the negative likelihood of the student’s t distribution,

$$(7.1) \quad \rho_s(x) = \log(1 + x^2/\nu),$$

where  $\nu$  is the degrees of freedom parameter.

For each penalty  $\rho$ , our aim is to solve the problem

$$\underset{x \geq 0}{\text{minimize}} \quad \|x\|_1 \quad \text{subject to} \quad \rho(b - Ax) \leq \sigma$$

via a series of approximate solutions of  $\mathcal{P}$ . The 1-norm regularizer on  $x$  encourages a sparse solution. In particular, we solve the nonlinear equation (1.3), where  $v$  is the value function of  $\mathcal{P}$ . This is the approach used by the SPGL1 software package [12]; the underlying theory, however, does not cover the Huber function. Also,  $\phi$  is not everywhere finite valued, which violates [12, Assumption 3.1]. Finally, the student’s t misfit (7.1) is nonconvex; however, the inverse function relationship (cf. Theorem 2.1) still holds, so we can achieve our goal, provided we can solve the root-finding problem.

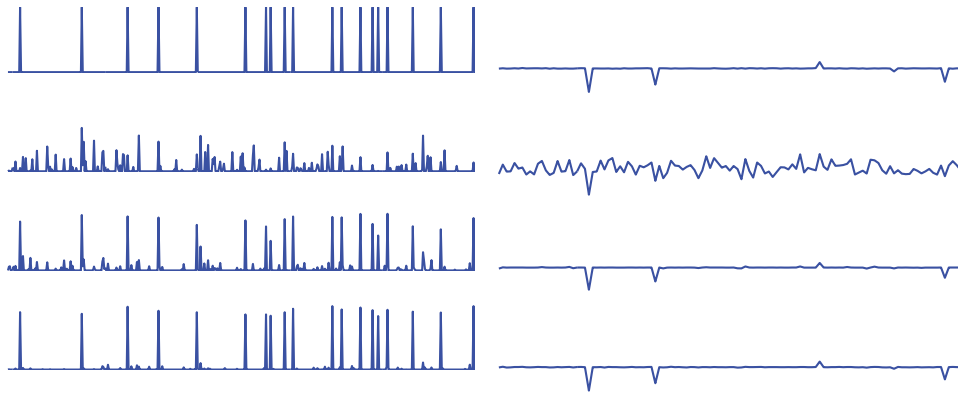


FIG. 7.1. Left, top to bottom: true signal, and reconstructions via least-squares, Huber, and student's  $t$ . Right, top to bottom: true errors, and least-squares, Huber, and student's  $t$  residuals.

Formula (6.14) computes the derivative of the value function associated with  $\mathcal{P}(b, \tau)$  for any convex differentiable  $\rho$ . The derivative requires  $\nabla\rho$ , evaluated at the optimal residual associated with  $\mathcal{P}(b, \tau)$ . For the Huber case, this is given by

$$(\nabla\rho(b - A\bar{x}))_i = \text{sign}(b_i - A_i\bar{x}) \cdot \min(|b_i - A_i\bar{x}|, \kappa).$$

The student's  $t$  misfit is also smooth, but nonconvex. Therefore, the formula (6.14) may still be applied—with the caveat that there is no guarantee of success. However, in all of the numerical experiments, we are able to find the root of (1.3).

We consider a common compressive sensing example: we want to recover a 20-sparse vector in  $\mathbb{R}_+^{512}$  from 120 measurements. We use a Gaussian measurement matrix  $A \in \mathbb{R}^{100 \times 1024}$ , where each entry is sampled from the distribution  $N(0, 1/10)$ . We generate measurements to test the BPDN formulation according to

$$b = Ax + w + \zeta,$$

where  $w \sim N(0, 0.005^2)$  is small Gaussian error, while  $\zeta$  contains five randomly placed large outliers sampled from  $N(0, 4)$ . For each penalty  $\rho$ , the  $\sigma$  parameter is the true measure of the error in that penalty, i.e.,  $\sigma_\rho = \rho(\zeta)$ . This allows a fair comparison between the penalties.

We expect the Huber function to out-perform the least squares penalty by budgeting the error level  $\sigma$  to allow a few large outliers, which will never happen with the quadratic. We expect the student's  $t$  penalty to work even better, because it is nonconvex and grows sublinearly as outliers increase. The results in Figure 7.1 demonstrate that this is indeed the case. In many instances the Huber function is able to do just as well as the student's  $t$ ; however, often the student's  $t$  does better (and never worse). Both robust penalties always do better than the least squares fit. The code is implemented in an extended version of SPGL1 and can be downloaded from <https://github.com/saravkin/spgl1>. The particular experiment presented here can be found in `tests/spgl1TestNN.m`.

## 8. Appendix: Proofs of results.

**Proof of Theorem 2.1.** Let  $\tau \in \mathcal{S}_{1,2}$  and set  $\sigma_\tau = v_1(\tau)$ . By assumption,  $\arg \min \mathcal{P}_{1,2}(\tau) \neq \emptyset$ . Let  $x_\tau \in \arg \min \mathcal{P}_{1,2}(\tau)$ , so that  $\psi_1(x_\tau) = \sigma_\tau$  and  $\psi_2(x_\tau) = \tau$ .

In particular,  $x_\tau$  is feasible for  $\mathcal{P}_{2,1}(\sigma_\tau)$ . Let  $\hat{x}$  be any other feasible point for  $\mathcal{P}_{2,1}(\sigma_\tau)$  so that  $\psi_1(\hat{x}) \leq \sigma_\tau = v_1(\tau) = \psi_1(x_\tau)$ . If  $\psi_1(\hat{x}) < \sigma_\tau = v_1(\tau)$ , then  $\psi_2(\hat{x}) > \tau$  since otherwise we contradict the definition of  $v_1(\tau)$ . If  $\psi_1(\hat{x}) = \sigma_\tau$ , then we claim that  $\psi_2(\hat{x}) \geq \tau$ . Indeed, if  $\psi_2(\hat{x}) < \tau$ , then  $\hat{x} \in \arg \min \mathcal{P}_{1,2}(\tau)$  but  $\psi_2(\hat{x}) < \tau$ , which contradicts the fact that  $\tau \in \mathcal{S}_{1,2}$ . Hence, every feasible point for  $\mathcal{P}_{2,1}(\sigma_\tau)$  has  $\psi_2(\hat{x}) \geq \tau$  with equality only if  $\psi_1(\hat{x}) = \sigma_\tau$ . But  $x_\tau$  is feasible for  $\mathcal{P}_{2,1}(\sigma_\tau)$  with  $\psi_2(x_\tau) = \tau$ . Therefore,  $x_\tau \in \arg \min \mathcal{P}_{2,1}(\sigma_\tau) \subset \{x \in X \mid \psi_1(x) = \sigma_\tau\}$ . Consequently,  $v_2(v_1(\tau)) = \tau$  and

$$(8.1) \quad \emptyset \neq \arg \min \mathcal{P}_{1,2}(\tau) \subset \arg \min \mathcal{P}_{2,1}(\sigma_\tau) \subset \{x \in X \mid \psi_1(x) = \sigma_\tau\}.$$

We now show that  $\arg \min \mathcal{P}_{2,1}(\sigma_\tau) \subset \arg \min \mathcal{P}_{1,2}(\tau)$ . Let  $\hat{x} \in \arg \min \mathcal{P}_{2,1}(\sigma_\tau)$ . In particular,  $\hat{x}$  is feasible for  $\mathcal{P}_{2,1}(\sigma_\tau)$ , so, by what we have already shown,  $\psi_2(\hat{x}) \geq \tau$  with equality only if  $\psi_1(\hat{x}) = \sigma_\tau$ . But, by our choice of  $\hat{x}$ ,  $\psi_2(\hat{x}) = v_2(v_1(\tau)) = \tau$ , so  $\psi_1(\hat{x}) = \sigma_\tau$ , i.e.,  $\hat{x} \in \arg \min \mathcal{P}_{1,2}(\tau)$ .

It remains to establish the final statement of the theorem. By (8.1), we already have that  $\{v_1(\tau) \mid \tau \in \mathcal{S}_{1,2}\} \subset \mathcal{S}_{2,1}$ , so we need only establish the reverse inclusion. For this, let  $\sigma \in \mathcal{S}_{2,1}$  and set  $\tau_\sigma = v_2(\sigma)$ . By interchanging the indices and applying the first part of the theorem, we have from (8.1) that

$$\emptyset \neq \arg \min \mathcal{P}_{2,1}(\sigma) \subset \arg \min \mathcal{P}_{1,2}(\tau_\sigma) \subset \{x \in X \mid \psi_2(x) = \tau_\sigma\}.$$

That is,  $\tau_\sigma \in \mathcal{S}_{1,2}$ , and, by (a),  $\sigma = v_1(v_2(\sigma)) = v_1(\tau_\sigma)$ .

**Proof of the inverse linear image (section 3.3).** For  $\lambda > 0$ , observe that

$$(8.2) \quad \begin{aligned} h^\pi(w, \lambda) &= \lambda \inf_{Ax = \lambda^{-1}w} p(x) \\ &= \lambda \inf_{As = w} p(\lambda^{-1}s) \quad (s := \lambda x) \\ &= \inf_{As = w} p^\pi(s, \lambda) \end{aligned}$$

$$(8.3) \quad = \inf \left\{ p^\pi(s, \zeta) \mid \hat{A} \begin{pmatrix} s \\ \zeta \end{pmatrix} = \begin{pmatrix} w \\ \lambda \end{pmatrix} \right\},$$

where

$$\hat{A} = \begin{bmatrix} A & 0 \\ 0 & 1 \end{bmatrix}.$$

Again by [34, Theorem 9.2] in conjunction with [34, Corollary 16.2.1], the function in (8.3) is closed if  $(\hat{A}^T)^{-1} \text{dom}(p^\pi)^* \neq \emptyset$ . Since, by [34, Corollary 13.5.1],  $\text{dom}(p^\pi)^* = \{(u, \eta) \mid p^*(u) \leq -\eta\}$ , we have

$$(\hat{A}^T)^{-1} \text{dom}(p^\pi)^* \neq \emptyset \quad \text{if and only if} \quad (A^T)^{-1} \text{dom} p^* \neq \emptyset.$$

Hence, by assumption, the function in (8.3) is closed, proper, and convex and equals  $h^\pi(w, \lambda)$  on the relative interior of its domain. Since  $h^\pi(w, \lambda)$  is closed, (8.2) implies that these functions must coincide.

**Proof of Lemma 4.1.** We first prove (4.1a). The conjugate of  $\delta((x, \tau) \mid \text{epi } h)$  is obtained as follows:

$$\begin{aligned} \delta^*((y, \mu) \mid \text{epi } h) &= \sup_{\tau, x} [\langle y, x \rangle + \mu\tau - \delta((x, \tau) \mid \text{epi } h)] \\ &= \sup_{\tau, x \in \text{dom } h} [\langle y, x \rangle + \mu\tau - \delta(h(x) - \tau \mid \mathbb{R}_-)] \\ &= \sup_{\omega, x \in \text{dom } h} [\langle y, x \rangle + \mu(h(x) - \omega) - \delta(\omega \mid \mathbb{R}_-)] \quad (\omega := h(x) - \tau) \\ &= \sup_{x \in \text{dom } h} [\langle y, x \rangle + \mu h(x) + \sup_{\omega} [-\mu\omega - \delta(\omega \mid \mathbb{R}_-)]] \\ &= \sup_{x \in \text{dom } h} [\langle y, x \rangle + \mu h(x) + \delta(\mu \mid \mathbb{R}_-)]. \end{aligned}$$

For  $\mu < 0$ , we obtain

$$\delta^*((y, \mu) \mid \text{epi } h) = -\mu \sup_x [\langle -\mu^{-1}y, x \rangle - h(x)] = -\mu h^*(-\mu^{-1}y).$$

Since  $h^*$  is necessarily a closed proper convex function, we obtain the result.

To see (4.1b), first note that the function

$$q(y) := \inf_{\mu > 0} [\tau\mu + \mu h^*(y/\mu)] = \inf_{\mu \geq 0} [\tau\mu + (h^*)^\pi(y, \mu)]$$

is the positively homogeneous function generated by the function  $y \mapsto \tau + h(y)$  [34, p. 35], and so it is convex in  $y$ . Next observe that the conjugate of  $q$  is given by

$$\begin{aligned} q^*(x) &= \sup_y \left[ \langle x, y \rangle - \inf_{\mu \geq 0} [\tau\mu + (h^*)^\pi(y, \mu)] \right] \\ &= \sup_{y, \mu \geq 0} [\langle x, y \rangle + \tau(-\mu) - (h^*)^\pi(y, \mu)] \\ &= \sup_{(y, \mu)} [\langle x, y \rangle + \tau\mu - (h^*)^\pi(y, -\mu)] \quad (\text{exchange } -\mu \text{ for } \mu) \\ &= \sup_{(y, \mu)} [\langle x, y \rangle + \tau\mu - \delta^*((y, \mu) \mid \text{epi } h)] \quad (\text{by (4.1a)}) \\ &= \delta((x, \tau) \mid \text{epi } h) = \delta(x \mid \text{lev}_h(\tau)). \end{aligned}$$

The result now follows from the biconjugate theorem [34, Theorem 12.2].

**Proof of Theorem 4.2.** Combining  $\mathcal{D}$  with (4.1b) and (4.2) gives

$$\begin{aligned} (8.4) \quad \hat{v}(b, \tau) &:= \sup_{u, \mu} [\langle b, u \rangle + \tau\mu - (\phi^*)^\pi(A^T u, -\mu) - \rho^*(u)] \\ &= \sup_u \left[ \langle b, u \rangle - \rho^*(u) - \inf_{\mu \leq 0} [\tau(-\mu) + (\phi^*)^\pi(A^T u, -\mu)] \right] \\ &= \sup_u \left[ \langle b, u \rangle - \rho^*(u) - \inf_{\mu \geq 0} [\tau(\mu) + (\phi^*)^\pi(A^T u, \mu)] \right] \\ &= \sup_u [\langle b, u \rangle - \rho^*(u) - \delta^*(A^T u \mid \text{lev}_\phi(\tau))], \end{aligned}$$

where the final equality follows from (4.1b). The equivalence (4.3a) follows from the definition of the conjugate, and the equivalence (4.3b) follows from [34, Theorems 16.3 and 16.4], which tell us that

$$\begin{aligned}
 g_\tau^*(b) &= \text{cl} \left( \rho \nabla \left[ \delta^* (A^T \cdot | \text{lev}_\phi(\tau)) \right]^* \right) (b) \\
 &= \text{cl} \left( \inf_{w_1+w_2=\cdot} \left[ \rho(w_1) + \inf_{Ax=w_2} \delta(x | \text{lev}_\phi(\tau)) \right] \right) (b) \\
 &= \text{cl} \left( \inf_{\phi(x) \leq \tau} \rho(\cdot - Ax) \right) (b) \\
 &= \text{cl} (v(\cdot, \tau)) (b).
 \end{aligned}$$

The uniqueness of  $u$  when  $\rho$  is differentiable follows from the essential strict convexity of  $\rho^*$  [34, Theorem 26.3].

**Proof of Lemma 4.3.**

*Part 1.* The inequality follows immediately from (4.1b). But it is also easily derived from the observation that if  $\mu > 0$  and  $x \in \text{lev}_\phi(\tau)$ , then

$$\begin{aligned}
 \tau\mu + \mu\phi^*(s/\mu) &\geq \tau\mu + \mu[\langle x, s/\mu \rangle - \phi(x)] \quad (\text{Fenchel–Young inequality}) \\
 &\geq \phi(x)\mu + \langle x, s \rangle - \mu\phi(x) \\
 &= \langle x, s \rangle.
 \end{aligned}$$

Taking the sup over  $x \in \text{lev}_\phi(\tau)$  gives the result.

*Part 2.* The proof uses the following three key facts:

- (i) By [34, Theorems 23.5 and 23.7], for any nonempty closed convex set  $U$  and  $\bar{u} \in U$ ,

$$(8.5) \quad \bar{v} \in N(\bar{u} | U) \iff \bar{u} \in \partial\delta^*(\bar{v} | U) = \arg \max_{u \in U} \langle \bar{v}, u \rangle.$$

- (ii) The Fenchel–Young inequality tells us that

$$(8.6) \quad \tau + \phi^*(\bar{s}) \geq \phi(\bar{x}) + \phi^*(\bar{s}) \geq \langle \bar{s}, \bar{x} \rangle.$$

- (iii) see [8, Lemma 26.17] or [45, Corollary 2.9.5]. Let  $g : \mathbb{R} \rightarrow \mathbb{R}$  be a convex function and  $\tau \in \mathbb{R}$  be such that  $\tau > \inf g$ . Then for every  $x \in \text{lev}_g(\tau)$

$$(8.7) \quad N(x | \text{lev}_g(\tau)) = \begin{cases} N(x | \text{dom } g) \cup \text{cone}(\partial g(x)) & \text{if } g(x) = \tau, \\ N(x | \text{dom } g) & \text{if } g(x) < \tau. \end{cases}$$

We divide the proof into two parts: (A) if  $S_1 \neq \emptyset$ , show  $S_1 \subset S_2$ , and (B) if  $S_2 \neq \emptyset$ , show  $S_2 \subset S_1$  and equality holds in (4.4). Combined, these implications establish Part 2 of the lemma.

(A) Let  $\bar{\mu} \in S_1$ . We show that  $\bar{\mu} \in S_2$ . First suppose  $\phi(\bar{x}) < \tau$ . By (8.7),  $N(\bar{x} | \text{lev}_\phi(\tau)) = N(\bar{x} | \text{dom } \phi)$ . Hence, by (8.5),  $\bar{s} \in N(\bar{x} | \text{dom } \phi)$ . Therefore, if  $\bar{\mu} = 0$ , we have  $\bar{\mu} \in S_2$ .

On the other hand, if  $\bar{\mu} > 0$ , by (8.5) and the fact that  $N(\bar{x} | \text{lev}_\phi(\tau)) = N(\bar{x} | \text{dom } \phi)$ , we have

$$\begin{aligned}
 \langle \bar{s}, \bar{x} \rangle &= \delta^*(\bar{s} | \text{dom } \phi) \\
 &= (\phi^*)^\infty(\bar{s}) && [34, \text{Theorem 13.3}] \\
 &= \tau + (\phi^*)^\pi(\bar{s}, 0) \\
 &\geq \tau\bar{\mu} + (\phi^*)^\pi(\bar{s}, \bar{\mu}) \\
 &> \bar{\mu}\phi(\bar{x}) + \bar{\mu}\phi^*(\bar{s}/\bar{\mu}) \\
 &\geq \bar{\mu}\langle \bar{s}/\bar{\mu}, \bar{x} \rangle && (\text{Fenchel–Young inequality}) \\
 &= \langle \bar{s}, \bar{x} \rangle.
 \end{aligned}$$

Since this cannot occur, it must be the case that  $\bar{\mu} = 0$  and  $\bar{\mu} \in S_2$ .

Now suppose that  $\phi(\bar{x}) = \tau$  and  $\bar{s} = 0$ . Then, for  $\mu > 0$ ,  $p_\tau(\bar{s}, \mu) = (\tau + \phi^*(0))\mu \geq 0$  by (8.6), and, for  $\mu = 0$ ,  $p_\tau(\bar{s}, \mu) = (\phi^*)^\infty(0) = 0$ . Therefore,  $0 = \inf_{0 \leq \mu} p_\tau(\bar{s}, \mu)$  with  $\mu = 0 \in S_1$ . But, in this case, it is also clear that  $0 \in S_2 \neq \emptyset$ , since  $\bar{s} = 0 \in N(\bar{x} | \text{dom } \phi)$ . Thus, if  $\bar{\mu} = 0$ , we have  $\bar{\mu} \in S_2$ . If  $\bar{\mu} > 0$ , then  $0 = \tau + \phi^*(0)$  since  $0 = p_\tau(0, \bar{\mu}) = (\tau + \phi^*(0))\bar{\mu}$ . But then, by (8.6),  $\phi(\bar{x}) + \phi^*(0) = \langle \bar{s}, \bar{x} \rangle = 0$  so that  $\bar{s} = 0 \in \partial\phi(\bar{x})$ . However,  $\phi(\bar{x}) = \tau > \inf \phi$ , so  $0 \notin \partial\phi(\bar{x})$  [34, Theorem 23.5(b)]. This contradiction implies that if  $\bar{s} = 0$ , then we must also have  $\bar{\mu} = 0$ , and, in particular, we have  $S_1 \subset S_2$ .

Finally, suppose that  $\phi(\bar{x}) = \tau$  and  $\bar{s} \neq 0$ . Then, by (8.7),

$$\text{either (a) } \bar{s} \in \text{cone}(\partial\phi(\bar{x})) \quad \text{or} \quad \text{(b) } \bar{s} \in N(\bar{x} | \text{dom } \phi).$$

Let us first suppose that  $\bar{s} \notin N(\bar{x} | \text{dom } \phi)$  so, in particular,  $\bar{s} \in \text{cone}(\partial\phi(\bar{x}))$ . As an immediate consequence, we have that  $S_2 \neq \emptyset$  and the only values of  $\mu$  for which  $\bar{s} \in \mu^+ \partial\phi(\bar{x})$  have  $\mu > 0$  since  $\bar{s} \notin N(\bar{x} | \text{dom } \phi)$ . Let  $0 < \hat{\mu} \in S_2$ . If  $\bar{\mu} = 0$ , then

$$\begin{aligned} \delta^*(\bar{s} | \text{dom } \phi) &= (\phi^*)^\infty(\bar{s}) \\ &= \inf_{0 \leq \mu} p_\tau(\bar{s}, \mu) \\ &\leq \tau\hat{\mu} + \hat{\mu}\phi^*(\bar{s}/\hat{\mu}) \\ &= \hat{\mu}\phi(\bar{x}) + \hat{\mu}[\langle \bar{s}/\hat{\mu}, \bar{x} \rangle - \phi(\bar{x})] \quad \left[ \begin{array}{l} \bar{s}/\hat{\mu} \in \partial\phi(\bar{x}) \text{ and} \\ [34, \text{Theorem 23.5(d)}] \end{array} \right] \\ &= \langle \bar{s}, \bar{x} \rangle \\ &\leq \delta^*(\bar{s} | \text{dom } \phi) \end{aligned}$$

so that  $\langle \bar{s}, \bar{x} \rangle = \delta^*(\bar{s} | \text{dom } \phi)$ , or, equivalently,  $\bar{s} \in N(\bar{x} | \text{dom } \phi)$ , contradicting the choice of  $\bar{s}$ . Hence, it must be the case that  $\bar{\mu} > 0$ . Again let  $0 < \hat{\mu} \in S_2$ . Then, by Part 1,

$$\begin{aligned} \delta^*(\bar{s} | \text{lev}_\phi(\tau)) &\leq p_\tau(\bar{s}, \bar{\mu}) \\ &= \inf_{0 \leq \mu} p_\tau(\bar{s}, \mu) \\ &\leq \tau\hat{\mu} + \hat{\mu}\phi^*(\bar{s}/\hat{\mu}) \\ &= \hat{\mu}\phi(\bar{x}) + \hat{\mu}[\langle \bar{s}/\hat{\mu}, \bar{x} \rangle - \phi(\bar{x})] \quad \left[ \begin{array}{l} \bar{s}/\hat{\mu} \in \partial\phi(\bar{x}) \text{ and} \\ [34, \text{Theorem 23.5(d)}] \end{array} \right] \\ &= \langle \bar{s}, \bar{x} \rangle \\ &\leq \delta^*(\bar{s} | \text{lev}_\phi(\tau)) \end{aligned}$$

so that  $\langle \bar{s}, \bar{x} \rangle = \bar{\mu}[\phi(\bar{x}) + \phi^*(\bar{s}/\bar{\mu})]$ , or, equivalently,  $\bar{s} \in \bar{\mu}\partial\phi(\bar{x})$ . Hence,  $\bar{\mu} \in S_2$ .

Finally, consider the case where  $0 \neq \bar{s} \in N(\bar{x} | \text{dom } \phi)$ . Then

$$\begin{aligned} \inf_{\mu \geq 0} p_\tau(\bar{s}, \mu) &\leq p_\tau(\bar{s}, 0) \\ &= (\phi^*)^\infty(\bar{s}) \\ &= \delta^*(\bar{s} | \text{dom } \phi) \quad [34, \text{Theorem 13.3}] \\ &= \langle \bar{s}, \bar{x} \rangle \quad (\text{by (8.5)}) \\ &= \delta^*(\bar{s} | \text{lev}_\phi(\tau)) \quad (\text{again by (8.5)}) \\ &\leq \inf_{\mu \geq 0} p_\tau(\bar{s}, \mu) \quad (\text{Part 1}), \end{aligned}$$



so  $0 \in S_1$  and  $0 \in S_2$ . If  $\bar{\mu} > 0$ , then this string of equivalences also implies that  $\langle \bar{s}, \bar{x} \rangle = p_\tau(\bar{s}, \bar{\mu}) = \bar{\mu}[\phi(\bar{x}) + \phi^*(\bar{s}/\bar{\mu})]$ , or, equivalently,  $\bar{s} \in \bar{\mu}\partial\phi(\bar{x})$  so that  $\bar{\mu} \in S_2$ . Putting this all together, we get that  $S_1 \subset S_2$ .

(B) Let  $\bar{\mu} \in S_2$ . If  $\bar{\mu} = 0$ , then

$$\begin{aligned} p_\tau(\bar{s}, 0) &= (\phi^*)^\infty(\bar{s}) \\ &= \delta^*(\bar{s} \mid \text{dom } \phi) \quad [34, \text{Theorem 13.3}] \\ &= \langle \bar{s}, \bar{x} \rangle \\ &\leq \delta^*(\bar{s} \mid \text{lev}_\phi(\tau)) \\ &\leq \inf_{\mu \geq 0} p_\tau(\bar{s}, \mu) \quad (\text{Part 1}). \end{aligned}$$

Therefore,  $\bar{\mu} = 0 \in S_1$  and equality holds in (4.4).

On the other hand, if  $\bar{\mu} > 0$ , then  $\bar{s}/\bar{\mu} \in \partial\phi(\bar{x})$ , and so

$$\begin{aligned} \tau\bar{\mu} + (\phi^*)^\pi(\bar{s}, \bar{\mu}) &= \bar{\mu}[\phi(\bar{x}) + \phi^*(\bar{s}/\bar{\mu})] \\ &= \bar{\mu} \langle \bar{x}, \bar{s}/\bar{\mu} \rangle \quad \left[ \begin{array}{l} \bar{s}/\bar{\mu} \in \partial\phi(\bar{x}) \text{ and} \\ [34, \text{Theorem 23.5(d)}] \end{array} \right] \\ &= \langle \bar{x}, \bar{s} \rangle \\ &\leq \delta^*(\bar{s} \mid \text{lev}_\phi(\tau)) \\ &\leq \inf_{\mu \geq 0} [\tau\mu + (\phi^*)^\pi(\bar{s}, \mu)] \quad (\text{Part 1}). \end{aligned}$$

Hence,  $\bar{\mu} \in S_1$  and equality holds in (4.4).

**Proof of Lemma 4.5.**

*Part 1.* The primal coercivity equivalence follows from [34, Theorems 8.4 and 8.7] since  $\text{hzn}(f(\cdot, b, \tau)) = \text{hzn}(\phi) \cap [-A^{-1}\text{hzn}(\rho)]$ .

*Part 2.* For the dual coercivity equivalence, let  $\hat{g}(u) = g_\tau(u) - \langle b, u \rangle$ , which is the objective of the reduced dual  $\mathcal{D}_\tau$ . By (4.3b),  $\hat{g}^*(0) = g_\tau^*(b) = \text{cl}(v(\cdot, \tau)) \leq v(b, \tau)$ . Therefore, the result follows from [34, Corollary 14.2.2] since by (8.8),  $\text{dom } v(\cdot, \tau) = \text{dom } \rho + A \text{dom } \phi$ .

**Proof of Theorem 5.1.** The expression for  $f^*$  is derived in (4.2). The weak and strong duality relationships as well as the expression for  $\partial v$  follow immediately from [36, Theorem 11.39].

Next, note that

$$(8.8) \quad \text{dom } f(\cdot, b, \tau) \neq \emptyset \iff \left[ \begin{array}{l} \exists x \in \text{lev}_\phi(\tau) \\ b - Ax \in \text{dom } \rho \end{array} \right] \iff b \in \text{dom } \rho + A \text{lev}_\phi(\tau).$$

Now assume that  $b \in \text{int}(\text{dom } \rho + A(\text{lev}_\phi(\tau)))$ . Recall from [34, Theorem 6.6 and Corollary 6.6.2] that

$$(8.9) \quad \text{int}(\text{dom } \rho + A(\text{lev}_\phi(\tau))) = \text{ri}(\text{dom } \rho) + A(\text{ri}(\text{lev}_\phi(\tau))).$$

Moreover, by [34, Theorem 7.6], for any convex function  $p$ ,

$$(8.10) \quad \text{ri}(\text{lev}_p(\tau)) = \{x \in \text{ri}(\text{dom } p) \mid p(x) < \tau\}.$$

Since  $b \in \text{int}(\text{dom } \rho + A(\text{lev}_\phi(\tau)))$ , (8.9)–(8.10) imply the existence of  $\bar{w} \in \text{ri}(\text{dom } \rho)$  and  $\bar{x} \in \text{ri}(\text{dom } \phi)$  with  $\phi(\bar{x}) < \tau$  such that  $b = \bar{w} + A\bar{x}$ . Since  $\phi$  is relatively

continuous on the relative interior of its domain [34, Theorem 10.1], there exists  $\delta > 0$  such that

$$\begin{aligned}(\bar{w} + \delta \mathcal{B}) \cap \text{dom } \rho &\subset \text{ri}(\text{dom } \rho), \\(\bar{x} + \delta \mathcal{B}) \cap \text{dom } \phi &\subset \text{ri}(\text{dom } \phi), \\ \phi(x) &< \frac{1}{2}(\phi(\bar{x}) + \tau) \quad \forall x \in (\bar{x} + \delta \mathcal{B}) \cap \text{dom } \phi.\end{aligned}$$

Set  $S_\rho = (\bar{w} + \delta \mathcal{B}) \cap \text{dom } \rho$  and  $S_\phi = (\bar{x} + \delta \mathcal{B}) \cap \text{dom } \phi$ . Since

$$\begin{aligned}\text{cone}(S_\rho + AS_\phi - b) &= \text{cone}(S_\rho - \bar{w}) + A\text{cone}(S_\phi - \bar{x}) \\ &= \text{span}(\text{dom } \rho - \bar{w}) + A\text{span}(\text{dom } \phi - \bar{x}) \\ &= \text{span}(\text{dom } \rho + A\text{dom } \phi - b) \\ &\supset \text{cone}(\text{dom } \rho + A\text{dom } \phi - b) \\ &= \mathbb{R}^m \quad (b \in \text{int}(\text{dom } \rho + A(\text{lev}_\phi(\tau))),\end{aligned}$$

we have  $0 \in \text{int}(S_\rho + AS_\phi - b)$ . Therefore, there exists an  $\epsilon > 0$  such that  $b + \epsilon \mathcal{B} \subset S_\rho + AS_\phi$ . Consequently, if  $\hat{b} \in b + \epsilon \mathcal{B}$  and  $|\hat{\tau} - \tau| < \frac{1}{2}(\tau - \phi(\bar{x}))$ , then  $\text{dom } f(\cdot, \hat{b}, \hat{\tau}) \neq \emptyset$  and so  $(\hat{b}, \hat{\tau}) \in \text{dom } v$ .

On the other hand, if  $(b, \tau) \in \text{int}(\text{dom } v)$ , then  $\text{dom } f(\cdot, \hat{b}, \hat{\tau}) \neq \emptyset$  for all  $(\hat{b}, \hat{\tau})$  near  $(b, \tau)$  so that  $\text{dom } f(\cdot, \hat{b}, \tau) \neq \emptyset$  for all  $\hat{b}$  near  $b$ . Hence,  $b \in \text{int}(\text{dom } \rho + A(\text{lev}_\phi(\tau)))$ .

### Proof of Theorem 5.2.

*Part 1.* First note that (5.1c) is equivalent to the optimality condition

$$(8.11) \quad 0 \in -A^T \partial \rho(b - A\bar{x}) + \partial \delta(\bar{x} \mid \text{lev}_\phi(\tau))$$

for the problem  $\mathcal{P}$ , and hence by [34, Theorem 23.8],  $\bar{x}$  solves  $\mathcal{P}$ . Moreover, by [34, Theorem 23.5], (5.1c) is equivalent to

$$b - A\bar{x} \in \partial \rho^*(\bar{u}), \quad \bar{x} \in \partial \delta^*(A^T \bar{u} \mid \text{lev}_\phi(\tau)),$$

or, equivalently,

$$(8.12) \quad b \in \partial \rho^*(\bar{u}) + A \partial \delta^*(A^T \bar{u} \mid \text{lev}_\phi(\tau)),$$

which by [34, Theorem 23.8] implies that  $\bar{u}$  solves the reduced dual  $\mathcal{D}_r$ .

*Part 2.* If  $\bar{x}$  solves  $\mathcal{P}$ , then

$$0 \in \partial[\rho(b - A(\cdot)) + \delta(\cdot \mid \text{lev}_\phi(\tau))](\bar{x}),$$

which by [34, Theorems 23.8, 23.9] and (5.1a) is equivalent to (8.11), which in turn is equivalent to (5.1c).

*Part 3.* If  $\bar{u}$  solves  $\mathcal{D}_r$ , then

$$b \in \partial[\rho^*(\cdot) + \delta^*(A^T(\cdot) \mid \text{lev}_\phi(\tau))](\bar{u}),$$

which by [34, Theorems 23.8, 23.9] and (5.1b) is equivalent to (8.12), which in turn is equivalent to (5.1c).

*Part 4.* The equivalence (5.1e) follows from (5.1d), Part 2 of Lemma 4.3, and the fact that  $A^T \bar{u} \in N(\bar{x} \mid \text{lev}_\phi(\tau))$  if and only if  $\bar{x} \in \partial \delta^*(A^T \bar{u} \mid \text{lev}_\phi(\tau))$ .

To see (5.1d), note that (4.5a), (5.1a), and Part 1 of Lemma 4.5 imply that the primal objective is coercive, so a solution  $\bar{x}$  exists. Hence, by Part 2, there exists  $\bar{u}$  so that  $(\bar{x}, \bar{u})$  satisfies (5.1c).

Analogously, (4.5b), (5.1b), and Part 2 of Lemma 4.5 imply that the solution  $\bar{u}$  to the dual exists, and so by Part 3, there exists  $\bar{x}$  such that the pair  $(\bar{x}, \bar{u})$  satisfies (5.1c). In either case, the subdifferential is nonempty and is given by (5.1d).

**Proof of Lemma 6.1.** Formula (6.3a) is just [34, Theorem 14.5]. The first equation in (6.3b) is obvious and the second follows from (6.3a) and the definition of the barrier cone. The formula (6.3c) is now obvious. Formulas (6.3d) and (6.3e) follow immediately from the definitions and [34, Corollary 8.3.3]. Formula (6.3f) follows from (6.3e), [34, Corollary 14.2.1], and [34, Corollary 16.4.2].

First note that (6.4) implies that  $\text{ri}(\text{cone}(U)) \cap \text{ri}(X) \neq \emptyset$ . Hence, the formula (6.5a) follows from [34, Theorem 16.4] and (6.3c). To see (6.5b), observe that the expression on the RHS is again an infimal convolution for which  $\inf = \min$  for the same reason as for (6.5a). The equivalence with  $(\phi^*)^\pi(z, \mu)$  follows from the calculus rules in section 3.3. For formula (6.5d), first note that

$$\begin{aligned} \inf_{\mu \geq 0} [\tau\mu + (\phi^*)^\pi(z, \mu)] &= \inf_{\mu \geq 0} \left[ \tau\mu + \inf_s [\delta^*(z - s \mid X) + \delta(s \mid \mu U^\circ)] \right] \\ &= \inf_s \left[ \delta^*(z - s \mid X) + \inf_{\mu \geq 0} [\tau\mu + \delta(s \mid \mu U^\circ)] \right] \\ &= \inf_s [\delta^*(z - s \mid X) + \tau\gamma(s \mid U^\circ)]. \end{aligned}$$

Again, the final infimum in this derivation is an infimal convolution for which  $\inf = \min$  for the same reasons as in (6.5a) since, by (6.3c) and [34, Theorem 14.5],

$$\text{dom}((\tau\gamma(\cdot \mid U^\circ))^*) = \text{dom}((\delta^*(\cdot \mid \tau U))^*) = \text{dom} \delta(\cdot \mid \tau U) = \tau U.$$

Therefore, an optimal  $\bar{s}$  in this infimal convolution exists giving  $\bar{\mu} = \gamma(\bar{s} \mid U^\circ)$  as the optimal solution to the first min in (6.5d).

Formula (6.5e) is an immediate consequence of (6.3d), (6.4), and [34, Corollary 23.8.1].

**Proof of Theorem 6.2.** By (6.3d) and the calculus rules for the relative interior [34, section 6], (5.1a) and (6.8) are equivalent. Similarly, by (6.3f) and [34, Theorem 6.3], (5.1b) and (6.9) are equivalent.

*Part 1.* Since (6.4) holds, the formula (6.5e) holds and so (6.10) and (5.1c) are equivalent. Hence, the result follows from Part 1 of Theorem 5.2.

*Part 2.* Since (5.1a) and (6.8) are equivalent, the result follows from Part 2 of Theorem 5.2.

*Part 3.* Since (5.1b) and (6.9) are equivalent, the result follows from Part 3 of Theorem 5.2.

*Part 4.* By (6.3e), (6.11) is equivalent to (4.5a) and (5.1a), and, by (6.3c), (6.12) is equivalent to (4.5b) and (5.1b). Therefore, by Theorem 5.2, (6.13) is equivalent to (5.1d) since  $\tau\gamma(s \mid U^\circ) = \inf_{\mu \geq 0} [\tau\mu + \delta(s \mid \mu U^\circ)]$ . The final equivalence is identical to that of Theorem 5.2.

**Proof of Lemma 6.3.** The formula for  $\text{dom } \phi$  follows from (6.17). Indeed, by (6.17),  $x \in \text{dom } \phi$  if and only if there exists  $s \in \mathbb{R}^k$  such that  $x - Ls \in \text{dom } \gamma(\cdot | U^\circ) = \text{cone}(U^\circ)$ , or, equivalently,  $x \in \text{cone}(U^\circ) + \text{Ran}(L) = \text{cone}(U^\circ) + \text{Ran}(B)$ . The formula for  $\text{hzn}(\phi)$  follows immediately from [34, Theorem 14.2] and (6.16). In particular,  $\phi$  is coercive if and only if  $\{0\} = \text{hzn}(\phi)$ , or, equivalently,  $\text{cone}(U) = \mathbb{R}^n$ , i.e.,  $0 \in \text{int}(U)$ .

Next we show that the  $\lambda$  given in (6.20) solves (6.18). First observe that the optimal  $\lambda$  must be greater than  $\gamma(w | U)$ , and from elementary calculus, the minimizer of the hyperbola  $\frac{1}{2\lambda}\|w\|_B^2 + \tau\lambda$  for  $\lambda \geq 0$  is given by  $\|w\|_B/\sqrt{2\tau}$ . Therefore, the minimizing  $\lambda$  is given by (6.20). Substituting this value into (6.18) gives (6.19).

It is now easily shown that the function in (6.19) is lower semicontinuous. Therefore, the equivalence in (6.18) follows from (4.1b).

**Proof of Theorem 6.4.** By [34, Theorem 7.6],

$$\text{ri}(\text{lev}_\phi(\tau)) = \{x | x \in \text{ri}(\text{dom } \phi), \phi(x) < \tau\}.$$

Hence, by Lemma 6.3, the equivalence between (5.1) and (6.21), (6.22), (6.24), (6.25), respectively, is easily seen. Therefore, Parts 1–4 follow immediately from Theorem 5.2.

**Proof of Corollary 6.5.** Condition (6.27a) occurs when  $\bar{\mu} = 0$  since  $0^+ \partial\phi(\bar{x}) = N(\bar{x} | \text{dom } \phi)$ . When  $\bar{\mu} > 0$ , by [34, Theorem 23.5],  $\partial\phi(x) = \arg \max_{w \in U} [\langle x, w \rangle - \frac{1}{2} \langle w, Bw \rangle]$ , so that  $w \in \partial\phi(x)$  if and only if  $x \in Bw + N(w | U)$ .

**Acknowledgments.** The authors are grateful to two anonymous referees for their remarkably careful reading of an intricate paper. Their detailed list of comments and suggestions led to many fixes and improvements.

#### REFERENCES

- [1] B. D. O. ANDERSON AND J. B. MOORE, *Optimal Filtering*, Prentice-Hall, Englewood Cliffs, N.J., 1979.
- [2] A. ARAVKIN, B. BELL, J. V. BURKE, AND G. PILLONETTO, *An  $\ell_1$ -Laplace robust Kalman smoother*, IEEE Trans. Automat. Control, 56 (2011), pp. 2898–2911.
- [3] A. ARAVKIN, J. BURKE, AND G. PILLONETTO, *Robust and trend following kalman smoothers using Student's t*, in Proceedings of the 16th IFAC Symposium on System Identification, 2012.
- [4] A. ARAVKIN, J. BURKE, AND G. PILLONETTO, *Sparse/robust estimation and kalman smoothing with nonsmooth log-concave densities: Modeling, computation, and theory*, J. Mach. Learn. Res., to appear.
- [5] A. ARAVKIN, M. P. FRIEDLANDER, F. HERRMANN, AND T. VAN LEEUWEN, *Robust inversion, dimensionality reduction, and randomized sampling*, Math. Program., 134 (2012), pp. 101–125.
- [6] N. ARONSZAJN, *Theory of reproducing kernels*, Trans. Amer. Math. Soc., 68 (1950), pp. 337–404.
- [7] A. AUSLENDER AND M. TEBoulLE, *Asymptotic Cones and Functions in Optimization and Variational Inequalities*, Springer, New York, 2003.
- [8] H. H. BAUSCHKE AND P. L. COMBETTES, *Convex Analysis and Monotone Operator Theory in Hilbert Spaces*, Springer, New York, 2011.
- [9] E. VANDEN BERG AND M. P. FRIEDLANDER, *SPGL1: A Solver for Large-Scale Sparse Reconstruction*. <http://www.cs.ubc.ca/labs/scl/index.php/Main/Spgl1> (June 2007).
- [10] E. VANDEN BERG AND M. P. FRIEDLANDER, *Probing the pareto frontier for basis pursuit solutions*, SIAM J. Sci. Comput., 31 (2008), pp. 890–912.
- [11] E. V. BERG AND M. P. FRIEDLANDER, *Theoretical and empirical results for recovery from multiple measurements*, IEEE Trans. Inform. Theory, 56 (2010), pp. 2516–2527.
- [12] E. V. BERG AND M. P. FRIEDLANDER, *Sparse optimization with least-squares constraints*, SIAM J. Optim., 21 (2011), pp. 1201–1229.

- [13] J. M. BORWEIN AND A. S. LEWIS, *Convex Analysis and Nonlinear Optimization*, Springer, New York, 2000.
- [14] S. BOYD AND L. VANDENBERGHE, *Convex Optimization*, Cambridge University Press, Cambridge, 2004.
- [15] R. BROCKETT, *Finite Dimensional Linear Systems*, Wiley, New York, 1970.
- [16] C. COMBARI, M. LAGHDIR, AND L. THIBAUT, *Sous-différentiels de fonctions convexes composées*, Ann. Sci. Math. Québec, 18 (1994), pp. 119–148.
- [17] F. CUCKER AND S. SMALE, *On the mathematical foundations of learning*, Bull. Amer. Math. Soc., 39 (2001), pp. 1–49.
- [18] D. L. DOHONO AND J. TANNER, *Sparse nonnegative solution of underdetermined linear equations by linear programming*, Proc. Natl. Acad. Sci. USA, 102 (2005), pp. 9446–9451.
- [19] D. DONOHO, *Compressed sensing*, IEEE Trans. Inform. Theory, 52 (2006), pp. 1289–1306.
- [20] B. EFRON, T. HASTIE, L. JOHNSTONE, AND R. TIBSHIRANI, *Least angle regression*, Annals Statist., 32 (2004), pp. 407–499.
- [21] I. EKELAND AND R. TEMAM, *Convex Analysis and Variational Problems*, Elsevier, New York, 1976.
- [22] T. EVGENIOU, M. PONTIL, AND T. POGGIO, *Regularization networks and support vector machines*, Adv. Comput. Math., 13 (2000), pp. 1–150.
- [23] S. FARAHMAND, G. GIANNAKIS, AND D. ANGELOSANTE, *Doubly robust smoothing of dynamical processes via outlier sparsity constraints*, IEEE Trans. Signal Process., 59 (2011), pp. 4529–4543.
- [24] J. GAO, *Robust  $l_1$  principal component analysis and its Bayesian variational inference*, Neural Comput., 20 (2008), pp. 555–572.
- [25] T. J. HASTIE AND R. J. TIBSHIRANI, *Generalized Additive Models*, Monogr. Statist. Appl. Probab. 43, Chapman and Hall, London, 1990.
- [26] T. J. HASTIE, R. J. TIBSHIRANI, AND J. FRIEDMAN, *The Elements of Statistical Learning: Data Mining, Inference and Prediction*, Springer, New York, 2001.
- [27] F. J. HERRMANN, M. P. FRIEDLANDER, AND O. YILMAZ, *Fighting the curse of dimensionality: compressive sensing in exploration seismology*, IEEE Signal Proc. Magazine, 29 (2012), pp. 88–100.
- [28] J. B. HIRIART-URRUTY AND C. LEMARÉCHAL, *Fundamentals of Convex Analysis*, Springer, New York, 2001.
- [29] P. J. HUBER, *Robust Statistics*, Wiley, New York, 1981.
- [30] D. MACKAY, *Bayesian non-linear modelling for the prediction competition.*, ASHRAE Trans., 100 (1994), pp. 3704–3716.
- [31] D. J. C. MACKAY, *Bayesian interpolation*, Neural Comput., 4 (1992), pp. 415–447.
- [32] H. M. MARKOWITZ, *Mean-Variance Analysis in Portfolio Choice and Capital Markets*, Wiley, New York, 1987.
- [33] M. PONTIL AND A. VERRI, *Properties of support vector machines*, Neural Comput., 10 (1998), pp. 955–974.
- [34] R. T. ROCKAFELLAR, *Convex Analysis*, Princeton Landmarks Math., Princeton University Press, Princeton, NJ, 1970.
- [35] R. T. ROCKAFELLAR, *Lagrange multipliers and optimality*, SIAM Rev., 35 (1993), pp. 183–238.
- [36] R. T. ROCKAFELLAR AND R. J.-B. WETS, *Variational Analysis*, Grundlehren Math. Wiss. 317, Springer, New York, 1998.
- [37] S. ROWEIS AND Z. GHAHRAMANI, *A unifying review of linear Gaussian models*, Neural Comput., 11 (1999), pp. 305–345.
- [38] S. SAITOH, *Theory of Reproducing Kernels and its Applications*, Longman, Longman, 1988.
- [39] B. SCHÖLKOPF, A. J. SMOLA, R. C. WILLIAMSON, AND P. L. BARTLETT, *New support vector algorithms*, Neural Comput., 12 (2000), pp. 1207–1245.
- [40] R. TAPIA, *The Isoperimetric Problem Revisited: Extracting a Short Proof of Sufficiency From Euler’s Approach to Necessity*, Technical report, Rice University, 2013.
- [41] R. TIBSHIRANI, *Regression shrinkage and selection via the Lasso*, J. Roy. Statist. Soc. Ser. B., 58 (1996), pp. 267–288.
- [42] M. TIPPING, *Sparse Bayesian learning and the relevance vector machine*, J. Mach. Learn. Res., 1 (2001), pp. 211–244.
- [43] V. VAPNIK, *Statistical Learning Theory*, Wiley, New York, 1998.
- [44] D. WIPF, B. RAO, AND S. NAGARAJAN, *Latent variable Bayesian models for promoting sparsity*, IEEE Trans. Inform. Theory, 57 (2011), pp. 6236–6255.
- [45] C. ZĂLINESCU, *Convex Analysis in General Vector Spaces*, World Scientific, River Edge, NJ, 2002.