

1 **PROXIMAL GRADIENT  $\mathcal{V}\mathcal{U}$ -METHOD WITH SUPERLINEAR CONVERGENCE**  
2 **FOR NONSMOOTH CONVEX OPTIMIZATION\***

3 SHUAI LIU<sup>†</sup>, CLAUDIA SAGASTIZÁBAL<sup>‡</sup>, AND MIKHAIL SOLODOV<sup>§</sup>

4 **Abstract.** The  $\mathcal{V}\mathcal{U}$ -theory for nonsmooth functions and the associated space decomposition have been used  
5 for studying the structure of nonsmoothness and for developing algorithms with superlinear convergence in those  
6 challenging (for fast convergence) settings. We extend the theory by defining a certain bivariate  $\mathcal{U}$ -Lagrangian  
7 function and the partial  $\mathcal{U}$ -Hessian. Utilizing smoothness properties of the new  $\mathcal{U}$ -Lagrangian we develop the  
8 Proximal Gradient  $\mathcal{V}\mathcal{U}$ -method for continuous nonsmooth convex optimization, and show its superlinear convergence  
9 under natural assumptions. The framework consists of a  $\mathcal{V}$ -step which is a prox-gradient step, and a  $\mathcal{U}$ -step which  
10 can be considered as a quasi-Newton step applied to the  $\mathcal{U}$ -Lagrangian. We show that partial  $\mathcal{U}$ - Hessians exist for  
11 most partly smooth functions. As an example, our method is applied to solving  $\ell_1$ -regularized problems. We exhibit  
12 the explicit process of constructing a basis of the  $\mathcal{U}$ -space and of calculating the  $\mathcal{U}$ -Hessian. We conclude with  
13 numerical results illustrating the method's performance.

14 **Key words.**  $\mathcal{V}\mathcal{U}$ -decomposition, proximal gradient  $\mathcal{V}\mathcal{U}$ -method, quasi-Newton method, proximal point,  
15 nonsmooth optimization, superlinear convergence

16 **MSC codes.** 49J52, 90C53, 49J53, 90C99, 58C20

17 **1. Introduction.** We consider the problem

18 (1.1) 
$$\min_{x \in \mathbb{R}^n} f(x), \text{ where } f: \mathbb{R}^n \rightarrow \mathbb{R} \text{ is a nonsmooth convex function.}$$

19 The  $\mathcal{V}\mathcal{U}$ -theory introduced in [20, 22] (closely related to partial smoothness [24]) has been  
20 used for the study of smooth structures in nonsmooth functions in [30, 21, 32, 34, 7, 2].  
21 As explained in [42, 28], nonsmoothness is particularly difficult for fast (i.e., superlinear)  
22 convergence. Despite this challenging context, the  $\mathcal{V}\mathcal{U}$ -theory provides a favorable setting  
23 for the development of superlinearly convergent algorithms [39, 35, 11, 13, 3]; see also [25, 40].  
24 The approach is to decompose the space  $\mathbb{R}^n$  into two orthogonal subspaces called  $\mathcal{V}$  and  $\mathcal{U}$ ,  
25 depending on a point  $\bar{x}$ . The  $\mathcal{V}$ -space is defined to be the subspace parallel to the affine  
26 hull of the subdifferential  $\partial f(\bar{x})$ , and  $\mathcal{U}$  consists of the directions such that the directional  
27 derivative  $f'(\bar{x}; \cdot)$  is linear. Roughly speaking, the  $\mathcal{V}$  and  $\mathcal{U}$  spaces are defined so that near  
28 the point  $\bar{x}$  the nonsmoothness of  $f$  is captured in the  $\mathcal{V}$ -space and the smoothness of  $f$  is  
29 captured in the  $\mathcal{U}$ -space. Through a parametrized Lagrangian defined on the  $\mathcal{U}$ -space, called  
30 the  $\mathcal{U}$ -Lagrangian, second-order Taylor expansions of  $f$  in  $\mathcal{U}$  can be obtained if a generalized  
31 Hessian (called  $\mathcal{U}$ -Hessian) exists for the  $\mathcal{U}$ -Lagrangian.

32 In the original  $\mathcal{V}\mathcal{U}$ -algorithm [22], the  $\mathcal{V}$ -step minimizes a prox-regularization of  $f$  in  
33 the  $\mathcal{V}$ -subspace, and the  $\mathcal{U}$ -step makes a Newton-type step in the  $\mathcal{U}$ -subspace of the  $\mathcal{U}$ -  
34 Lagrangian (where the  $\mathcal{U}$ -Lagrangian looks smooth). The superlinear convergence requires  
35 the existence of a positive definite  $\mathcal{U}$ -Hessian. This conceptual approach did not address how  
36 to identify the  $\mathcal{V}\mathcal{U}$ -geometry along an algorithmic procedure. In order to pass from theory

---

\*Submitted to the editors DATE.

**Funding:** Research of the first author was supported by the National Natural Science Foundation of China (Grant No. 12001208) and the São Paulo Research Foundation (FAPESP) Grant 2017/15936-2. The second author is supported by CNPq Grant 307509/2023-0 and by PRONEX–Optimization. The third author is supported in part by CNPq Grant 306775/2023-9, by FAPERJ Grant E-26/200.347/2023, and by PRONEX–Optimization.

<sup>†</sup>School of Software, South China Normal University, Shishan Town, Nanhai District, Foshan, Guangdong, 528225, China (shuai0liu@gmail.com).

<sup>‡</sup> Universidade Estadual de Campinas, Rua Sérgio Buarque de Holanda 651, Campinas, SP, 13083-859, Brazil sagastiz@unicamp.br).

<sup>§</sup> IMPA - Instituto de Matemática Pura e Aplicada, Estrada Dona Castorina 110, Rio de Janeiro, RJ, 22460-320, Brazil (solodov@impa.br).

37 to computational implementation, it is important to examine the structure of (nonsmooth)  
 38 functions. The specific  $\mathcal{V}\mathcal{U}$ -theory for finite-max functions and the numerical analysis of the  
 39 relevant  $\mathcal{V}\mathcal{U}$ -objects were considered in [29] and [11]. A more general class of functions  
 40 given in [30, 32], said to have primal-dual-gradient (PDG) structure, identifies a *fast track*,  
 41 which are points from where fast Newton-type steps are possible. Links with the  $\mathcal{V}\mathcal{U}$ -theory  
 42 of partly smooth functions are the subject of [15, 26, 10, 14]. In particular, [14, Theorem 3.2]  
 43 establishes a one-to-one correspondence with the fast track and the active manifold of a partly  
 44 smooth function.

45 An important step towards implementable  $\mathcal{V}\mathcal{U}$  algorithms is [31], where it is shown that  
 46 proximal points are on the fast track. This result suggests that the  $\mathcal{V}$ -step can be implemented  
 47 by a prox-step on  $f$  when it is easy to compute, or by a bundle method [18] that approximates  
 48 this step by computing proximal points of successive cutting-planes models of  $f$ . Fully imple-  
 49 mentable (fast)  $\mathcal{V}\mathcal{U}$ -algorithms for solving (1.1) are scarce, because they need to approximate  
 50 sufficiently well both the (exact proximal)  $\mathcal{V}$ -step and the (exact)  $\mathcal{U}$ -Newton direction. Gen-  
 51 erally, this involves solving at least two quadratic programming problems per iteration. Such  
 52 is the case of [35], the first  $\mathcal{V}\mathcal{U}$ -algorithm for problems like (1.1), where no specific structure  
 53 for  $f$  is required. The work [36] proposes two sequential Newtonian methods based on local  
 54 parameterizations obtained from relating  $\mathcal{V}\mathcal{U}$ -theory with Riemannian geometry. Like for the  
 55  $\mathcal{V}\mathcal{U}$ -theory, considering a family of functions with specific properties leads to more targeted  
 56 implementations. For maximum eigenvalue and convex finite-max functions, we mention [39]  
 57 and [13]. How the  $\mathcal{V}\mathcal{U}$ -decomposition can be iteratively constructed by bundle methods for  
 58 a certain class that includes max-functions is explored in [7].

59 When  $f$  in (1.1) has additive structure as in (3.1) further below, subtle geometrical  
 60 properties of the proximal operator allow [2] to asymptotically detect the correct  $\mathcal{V}$ -step by  
 61 means of a proximal gradient (PG) method [4]. Depending on the nonsmooth term, this  
 62 calculation is explicit, or entails solving a simple quadratic program. The  $\mathcal{U}$ -step corrects the  
 63 PG iterate by a certain Newton-like direction, computed by solving a (possibly another, second)  
 64 quadratic programming problem. When applied to the same class of functions, our proposal  
 65 eliminates the latter second quadratic program, thanks to a suitable shifting of the optimal  
 66 subgradient resulting from the PG iterate calculation. The full corresponding algorithm,  
 67 named Proximal Gradient  $\mathcal{V}\mathcal{U}$  method or PGVU for short, is given in Algorithm 3.1 below.

68 In order to analyze convergence of the PGVU method, an important extension of the  
 69  $\mathcal{V}\mathcal{U}$ -theory is needed. In all the mentioned studies, given a point  $\bar{x} \in \mathbb{R}^n$ , the  $\mathcal{U}$ -Lagrangian  
 70 is a single-variable function, defined considering a subgradient  $\bar{g} \in \partial f(\bar{x})$  as a parameter.  
 71 But in the algorithmic setting we have to deal with a sequence of subgradients ( $g^k \in \partial f(x^k)$   
 72 at iteration  $k$ ), that change the parameter defining the  $\mathcal{U}$ -Lagrangian along iterations. As  
 73 illustrated by Example 2.2 below, with more than one fast track converging to a minimizer  $\bar{x}$ ,  
 74 different subgradients yield  $\mathcal{U}$ -Lagrangians associated with different fast tracks. To prevent  
 75 possible oscillatory behaviour, in our  $\mathcal{U}$ -Lagrangian definition, the subgradient is no longer  
 76 a parameter, but another variable. Accordingly, we extend the theories to such bivariate  
 77  $\mathcal{U}$ -Lagrangian, defining a partial  $\mathcal{U}$ -Hessian as the general partial Hessian of the new  $\mathcal{U}$ -  
 78 Lagrangian. Properties that hold for the single-variable  $\mathcal{U}$ -Lagrangian are now shown to hold  
 79 for the bivariate  $\mathcal{U}$ -Lagrangian. Thanks to our extended  $\mathcal{V}\mathcal{U}$ -theory, the proposed Proximal  
 80 Gradient  $\mathcal{V}\mathcal{U}$ -method has superlinear convergence, requiring only the (natural, for potential  
 81 fast convergence) assumptions of the existence of a positive-definite  $\mathcal{U}$ -Hessian at a solution  
 82  $\bar{x}$  such that  $0 \in \text{ri } \partial f(\bar{x})$ . Moreover, we show that any convex partly smooth function that  
 83 satisfies  $0 \in \text{ri } \partial f(\bar{x})$  automatically has a partial  $\mathcal{U}$ -Hessian. Finally, we demonstrate the  
 84 constructive process through the application of PGVU to  $\ell_1$ -regularized minimization.

85 The rest of the paper is organized as follows. In the remaining part of Section 1 we  
 86 introduce the notation. In Section 2, we lay out the foundation of  $\mathcal{V}\mathcal{U}$ -theory for the

87 development of our Proximal Gradient  $\mathcal{V}\mathcal{U}$ -method. We give the definition and smoothness  
 88 properties of the bivariate  $\mathcal{U}$ -Lagrangian function and show that computing the proximal  
 89 point can serve as the  $\mathcal{V}$ -step. In Section 3, we give the details of our PGVU-method and  
 90 show its global convergence. The definition of a partial  $\mathcal{U}$ -Hessian is given in Section 4  
 91 and, under the assumption of a positive definite partial  $\mathcal{U}$ -Hessian and  $0 \in \text{ri } \partial f(\bar{x})$ , we  
 92 prove that PGVU is superlinearly convergent. In Section 4.2, we show that all convex partly  
 93 smooth functions satisfying  $0 \in \text{ri } \partial f(\bar{x})$  have a partial  $\mathcal{U}$ -Hessian at  $\bar{x}$ . Section 5 applies the  
 94 proposed Proximal Gradient  $\mathcal{V}\mathcal{U}$ -method to  $\ell_1$ -regularized minimization. We first verify the  
 95 existence of a  $\mathcal{U}$ -Hessian. Then we provide an inexact prox-step as the  $\mathcal{V}$ -step and construct  
 96 a basis for the  $\mathcal{U}$ -space. Numerical results reported at the end of this section show that PGVU  
 97 performs well both in terms of computational time and accuracy. Concluding remarks are  
 98 given in Section 6.

99 **Notation.** We mostly follow [41]. Let  $\bar{\mathbb{R}} = [-\infty, \infty]$ . By  $\partial f(x)$  we denote the limiting  
 100 subdifferential of  $f$  at  $x$ , and by  $\partial^\infty f$  the horizon subdifferential of  $f$ . This is needed to  
 101 refer to some cited results. Of course, for a convex finite-valued  $f$ ,  $\partial f(x)$  is the usual  
 102 subdifferential in Convex Analysis. The notation  $f'(x; d)$  is the directional derivative at  $x$  in  
 103 the direction  $d$ . For a smooth bivariate function  $f(x, y)$ ,  $\nabla_x f(x, y)$  and  $\nabla_{xx}^2 f(x, y)$  are the  
 104 partial gradient and partial Hessian of  $f$  with respect to the variable  $x$ . For given points  $\bar{x}$   
 105 and  $\bar{y}$ , the partial subdifferential  $\partial_x f(\bar{x}, \bar{y})$  is defined to be the subdifferential of  $f(\cdot, \bar{y})$  at  
 106  $\bar{x}$ . The indicator function of a convex set  $C$  is  $\delta_C(\cdot)$  and its interior and relative interior are  
 107 respectively  $\text{int } C$  and  $\text{ri } C$ . The distance of a point  $x$  to a set  $C$  is  $\text{dist}(x; C) := \inf_{z \in C} \|z - x\|$ .  
 108 The Euclidean closed ball in  $\mathbb{R}^n$  centered at  $\bar{x}$  with radius  $\epsilon \geq 0$  is denoted by  $B(\bar{x}, \epsilon)$  and  
 109 the ball in  $\mathbb{R}^m$  is  $B^m(\bar{x}, \epsilon)$ . For a function  $f$ , its minimal value is denoted by  $f^*$  and its  
 110 set of minimizers by  $S$ . The vector  $e^j \in \mathbb{R}^n$  has all of its components null, except for  
 111  $e_j^j = 1$ . Regarding convergence rates, the notation “little o” in [38] for scalars is used for  
 112 vectors, as follows. For vector sequences  $\mathbb{R}^n \supset \{x^k\} \rightarrow \bar{x}$  and  $\mathbb{R}^m \supset \{y^k\} \rightarrow \bar{y}$ , and  
 113 for  $\|\cdot\|$  the Euclidean norm in the corresponding space,  $x^k = o(y^k)$  is short hand for  
 114 “ $\forall \epsilon > 0, \exists K : \|x^k\| \leq \epsilon \|y^k\|$  for all  $k \geq K$ ”. The notation for “big O” term is used in a  
 115 similar manner. The class of twice continuously differentiable functions is  $C^2$ .

116 **2. Elements of the  $\mathcal{V}\mathcal{U}$ -theory.** We start with the definition of the two subspaces in  
 117 question.

118 **DEFINITION 2.1 ( $\mathcal{V}\mathcal{U}$ -decomposition).** *Given a convex function  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  and a point*  
 119  *$\bar{x}$ , the  $\mathcal{V}\mathcal{U}$ -decomposition of  $\mathbb{R}^n$  associated with  $f$  and  $\bar{x}$  is defined by the subspaces*

$$120 \quad \mathcal{V}(\bar{x}) = \text{span}(\partial f(\bar{x}) - g), \quad \mathcal{U}(\bar{x}) = \mathcal{V}(\bar{x})^\perp,$$

121 *where  $g$  is an arbitrary subgradient in  $\partial f(\bar{x})$ .*

122 The respective dimensions are  $m = \dim \mathcal{U}(\bar{x})$  and  $n - m = \dim \mathcal{V}(\bar{x})$ . As vector spaces,  
 123  $\mathcal{V}(\bar{x})$  and  $\mathcal{U}(\bar{x})$  are endowed with a scalar product and norm induced from  $\mathbb{R}^n$ . When clear  
 124 from the context, the short forms  $\mathcal{V}$  and  $\mathcal{U}$  are used below.

125 The algebraic form of the decomposition depends on two matrices:

$$126 \quad \begin{aligned} \mathbb{R}^{n \times m} \ni U &: = \text{a matrix whose columns form an orthonormal basis for } \mathcal{U} \\ \mathbb{R}^{n \times (n-m)} \ni V &: = \text{a matrix whose columns form a basis for } \mathcal{V}, \\ &\text{with Moore-Penrose pseudo-inverse } V^\dagger := (V^T V)^{-1} V^T. \end{aligned}$$

127 More specifically, the  $\mathcal{U}$ - and  $\mathcal{V}$ -components of any  $x \in \mathbb{R}^n$  are defined by

$$128 \quad x_u := U^T x, \quad x_v := V^\dagger x.$$

129 By the definitions of  $\mathcal{V}$  and  $\mathcal{U}$ , the set  $U^\top \partial f(\bar{x})$  is a singleton, and hence,

$$130 \quad (2.1) \quad \bar{g}_u := U^\top g \text{ for any } g \in \partial f(\bar{x}).$$

131 **2.1. The  $\mathcal{U}$ -Lagrangian.** Given  $\bar{g} \in \partial f(\bar{x})$ , the single-variable  $\mathcal{U}$ -Lagrangian of  $f$  is

$$132 \quad \mathbb{R}^m \ni u \mapsto L_U^{\bar{g}}(u) := \inf_{w \in \mathbb{R}^{n-m}} \{f(\bar{x} + Uu + Vw) - \langle \bar{g}_v, V^\top Vw \rangle\}.$$

133 The associated set of  $\mathcal{V}$ -space minimizers is

$$134 \quad W^{\bar{g}}(u) := \left\{ w \in \mathbb{R}^{n-m} : L_U^{\bar{g}}(u) = f(\bar{x} + Uu + Vw) - \langle \bar{g}_v, V^\top Vw \rangle \right\}.$$

135 By its definition, the  $\mathcal{U}$ -Lagrangian is finite-valued and convex on  $\mathbb{R}^n$ . When, in addition,  
136  $\bar{g}_v \in V^\top \text{ri } \partial f(\bar{x})$ , it is shown in [22, Theorem 3.2, Theorem 3.3(ii)] that

$$137 \quad (2.2) \quad W^{\bar{g}}(0) = \{0\}, L_U^{\bar{g}}(0) = f(\bar{x}), L_U^{\bar{g}} \text{ is differentiable at } 0 \text{ with } \nabla L_U^{\bar{g}}(0) = \nabla \hat{f}_{\bar{x}}(0) = \bar{g}_u.$$

138 Evaluating the  $\mathcal{U}$ -Lagrangian at some  $\mathcal{V}$ -minimizer yields the following special first-order  
139 expansion for  $f$ :

$$140 \quad \forall w^{\bar{g}}(u) \in W^{\bar{g}}(u), f(\bar{x} + Uu + Vw^{\bar{g}}(u)) = f(\bar{x}) + \langle \bar{g}_u, u \rangle + \langle \bar{g}_v, V^\top Vw^{\bar{g}}(u) \rangle + o(Uu).$$

141 When, in addition, second-order approximation for  $f$  exists along the  $\mathcal{U}$ -subspace, a Newton-  
142 like step is possible, opening the way to superlinearly convergent schemes; see Sections 2.3  
143 and 4 below.

144 In the sequel, we shall introduce an important advance with respect to the previous  $\mathcal{V}\mathcal{U}$ -  
145 literature. It has to do with the following considerations. Note that the original  $\mathcal{U}$ -Lagrangian  
146 from [22] was defined for some fixed  $\bar{g} \in \text{ri } \partial f(\bar{x})$ . While it is true that (2.1) guarantees  
147 that the  $\mathcal{U}$ -component  $\bar{g}_u$  is the same for all  $\bar{g} \in \partial f(\bar{x})$ , the argument is not valid for the  
148  $\mathcal{V}$ -component  $\bar{g}_v$ . If, as in the example below, the value of  $\bar{g}_v$  modifies the  $\mathcal{V}$ -minimizer,  
149 different  $\mathcal{U}$ -Lagrangians emerge from different  $\bar{g}_v$ .

150 **EXAMPLE 2.2** (A function with structured nonsmoothness). Given a scalar  $a > 0$ , for  
151  $(u, v) \in \mathbb{R}^2$  the function

$$152 \quad F(u, v) = \max \left\{ \frac{a}{2}u^2, |v| \right\} = \frac{a}{2}u^2 + \max \left\{ 0, |v| - \frac{a}{2}u^2 \right\}$$

153 is differentiable everywhere except for points satisfying the equation  $|v| = \frac{a}{2}u^2$ . Its unique  
154 minimizer is  $\bar{x} = (0, 0)$ , where the subdifferential is  $\partial F(\bar{x}) = \{0\} \times [-1, 1]$ . Figure 2.1 shows  
155 that the graph of  $F$  is  $\mathcal{U}$ -shaped along the  $u$ -axis and  $\mathcal{V}$ -shaped along the  $v$ -axis. The  $\mathcal{V}\mathcal{U}$ -  
156 decomposition at  $\bar{x}$  gives  $\mathcal{U} = \mathbb{R} \times 0$  and  $\mathcal{V} = 0 \times \mathbb{R}$ . Then for any  $\bar{g} \in \text{ri } \partial F(\bar{x}) = 0 \times (-1, 1)$   
157 we have  $\bar{g}_v \in (-1, 1)$ . (When clear that a point is in  $\mathbb{R} \times 0$  or  $0 \times \mathbb{R}$ , we omit the 0 component.)  
158 Working out the calculations of the three cases for the  $\mathcal{V}$ -minimizers, we obtain that

$$159 \quad W^{\bar{g}}(u) = \begin{cases} \left\{ \frac{a}{2}u^2 \right\}, & \text{if } \bar{g}_v \in (0, 1), \\ \left\{ v : |v| \leq \frac{a}{2}u^2 \right\}, & \text{if } \bar{g}_v = 0, \\ \left\{ -\frac{a}{2}u^2 \right\}, & \text{if } \bar{g}_v \in (-1, 0). \end{cases} \implies L_U^{\bar{g}}(u) = (1 - |\bar{g}_v|) \frac{a}{2}u^2.$$

160 Notice the dependence of the  $\mathcal{U}$ -Lagrangian on the chosen subgradient.

161 For functions with structured nonsmoothness, the  $\mathcal{V}\mathcal{U}$ -decomposition is useful to reveal  
162 hidden smoothness. For  $F$ , this relates to the trajectory below:

$$163 \quad \chi^{\bar{g}}(u) = \left\{ \bar{x} + (u, v^{\bar{g}}(u)) : \text{for } u \in \mathbb{R} \text{ and } v^{\bar{g}}(u) = \frac{a}{2} \text{sign}(\bar{g}_v)u^2 \in W^{\bar{g}}(u) \right\}.$$

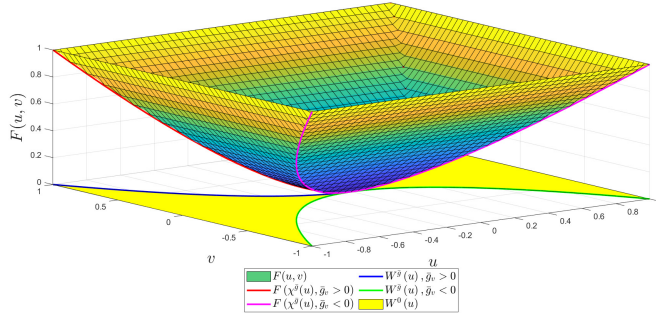


Fig. 2.1: Function  $F$  in Example 2.2, the associated  $\mathcal{V}$ -minimizers, and fast tracks

164 In the parlance of [31], this is a *fast track* along which  $F$  can be expanded up to second order:

165 
$$F(\chi^{\bar{g}}(u)) = \frac{a}{2}u^2, \quad \nabla_u F(\chi^{\bar{g}}(u)) = au, \quad \text{and} \quad \nabla_{uu}^2 F(\chi^{\bar{g}}(u)) = a.$$

166 The fast track  $\{(u, \frac{a}{2}u^2)\}$  in Figure 2.1 is obtained with  $\bar{g} = 0$ .

167 The situation illustrated by Example 2.2, with trajectories of smoothness depending on  
 168 the choice of the subgradient  $\bar{g}$ , motivates the consideration of a *bivariate*  $\mathcal{U}$ -Lagrangian,  
 169 introduced in this work.

170 **2.2. Two-variable  $\mathcal{U}$ -Lagrangian.** Pursuing further the analysis of the impact of  $g_v$   
 171 on the  $\mathcal{U}$ -objects, we now consider extending the  $\mathcal{V}\mathcal{U}$ -decomposition theory to a setting in  
 172 which  $g_v$  is an *argument* of the  $\mathcal{U}$ -Lagrangian. Rather than depending only on  $u$ , the function  
 173 has primal and dual variables, that is,  $(u, \bar{g}_v) \in \mathbb{R}^m \times \mathbb{R}^{n-m}$ .

174 **DEFINITION 2.3 ( $\mathcal{U}$ -Lagrangian with two variables).** *The bivariate  $\mathcal{U}$ -Lagrangian of  $f$*   
 175 *is defined from  $\mathbb{R}^m \times V^\dagger \partial f(\bar{x})$  to  $\bar{\mathbb{R}}$  as follows:*

176 
$$\mathbb{R}^m \times V^\dagger \partial f(\bar{x}) \ni (u, g_v) \mapsto L_U(u, g_v) := \inf_{w \in \mathbb{R}^{n-m}} \{f(\bar{x} + Uu + Vw) - \langle g_v, V^T Vw \rangle\},$$

177 and the associated set of  $\mathcal{V}$ -space minimizers is

178 
$$W(u, g_v) := \{w \in \mathbb{R}^{n-m} : L_U(u, g_v) = f(\bar{x} + Uu + Vw) - \langle g_v, V^T Vw \rangle\}.$$

179 The notation in this work,  $L_U(u, \bar{g}_v)$ , should not be confused with [31], where only  $u$  is a  
 180 variable, and the semicolon in  $L_U(u; \bar{g}_v)$  is used to expose that  $\bar{g}_v$  is a parameter.

181 The set  $\text{ri } \partial f(\bar{x}) := \{g \in \mathbb{R}^n : g + \text{int } B(0, \eta) \cap \mathcal{V} \subset \partial f(\bar{x}) \text{ for some } \eta > 0\}$  defines the  
 182 subdifferential relative interior. For each  $g \in \text{ri } \partial f(\bar{x})$ , we have  $U\bar{g}_u + Vg_v + \frac{\eta Vw}{\|Vw\|} \in \partial f(\bar{x})$   
 183 for all  $w \in \mathbb{R}^{n-m}$  and the convexity of  $f$  implies that, for any  $(u, w) \in \mathbb{R}^m \times \mathbb{R}^{n-m}$ , it holds  
 184 that

185 (2.3) 
$$f(\bar{x} + Uu + Vw) \geq f(\bar{x}) + \langle \bar{g}_u, u \rangle + \langle g_v, V^T Vw \rangle + \eta \|Vw\|.$$

186 In order to properly deal with variations on the  $\mathcal{U}$ -component in a manner that is uniform  
 187 relative to interior subgradients, for a small positive number  $\eta$ , we define the closed subset

188 (2.4) 
$$\eta\text{-ri } \partial f(\bar{x}) := \{g \in \mathbb{R}^n : g + B(0, \eta) \cap \mathcal{V} \subset \partial f(\bar{x})\} \text{ and let } G_v(\bar{x}) := V^\dagger \eta\text{-ri } \partial f(\bar{x}).$$

189 To show some continuity and smoothness properties of the  $\mathcal{U}$ -objects, we consider  $(u, g_v)$   
 190 as a perturbation parameter in a family of parametric minimization problems with value  
 191 function equal to the  $\mathcal{U}$ -Lagrangian, and solution mapping equal to the set of  $\mathcal{V}$ -minimizers.

192 LEMMA 2.4. *Given  $G_v(\bar{x})$  from (2.4), the mapping*

$$193 \quad \Phi(w, u, g_v) := f(\bar{x} + Uu + Vw) - \langle g_v, V^T Vw \rangle + \delta_{G_v(\bar{x})}(g_v),$$

194 *defined from  $\mathbb{R}^{n-m} \times \mathbb{R}^m \times \mathbb{R}^{n-m}$  to  $\bar{\mathbb{R}}$ , is proper, lsc, and level-bounded in  $w$  locally uniformly  
 195 in  $(u, g_v)$ ; see [41, Definition 1.16].*

196 *Proof.* Clearly,  $\Phi$  is proper because  $f$  is finite-valued and  $\Phi$  is lsc by the continuity  
 197 of  $f$  and closedness of  $G_v(\bar{x})$ , which follows from the definitions in (2.4). To show the  
 198 property of uniform level-boundedness, for all  $\alpha \in \mathbb{R}$  consider the mapping  $S_\alpha(u, g_v) := \{w :$   
 199  $\Phi(w, u, g_v) \leq \alpha\}$ . This mapping is nonempty (and can be unbounded) only when  $g_v \in G_v(\bar{x})$ ,  
 200 in which case

$$201 \quad S_\alpha(u, g_v) = \{w : f(\bar{x} + Uu + Vw) - \langle g_v, V^T Vw \rangle \leq \alpha\}.$$

202 In view of (2.3), for any such  $g_v$  and  $(u, w) \in \mathbb{R}^{m+p}$  with  $w$  an element of  $S_\alpha(u, g_v)$ ,  
 203 we have  $\alpha \geq f(\bar{x} + Uu + Vw) - \langle g_v, V^T Vw \rangle \geq f(\bar{x}) + \langle \bar{g}_u, u \rangle + \eta \|Vw\|$ , and therefore,  
 204  $S_\alpha(u, g_v) \subset T(u, g_v) := \{w : f(\bar{x}) + \langle \bar{g}_u, u \rangle + \eta \|Vw\| \leq \alpha\}$ . By [41, Example 5.17(b)],  
 205 it suffices to show that  $S_\alpha$  is uniformly bounded below. We first derive a uniform local  
 206 bound for the larger mapping  $T$ , by considering  $(u, g_v) \in B(\bar{u}, \epsilon) \times G_v(\bar{x}) \cap B(\bar{g}_v, \epsilon)$ , for  
 207 some  $(\bar{u}, \bar{g}_v) \in \mathbb{R}^m \times G_v(\bar{x})$ . As  $\bar{g}_u$  is fixed, the uniform bound follows, because  $\langle \bar{g}_u, u \rangle \geq$   
 208  $-\|\bar{g}_u\|(\epsilon + \|\bar{u}\|)$ , which in particular yields that  $\eta \|Vw\| \leq \alpha - f(\bar{x}) + \|\bar{g}_u\|(\epsilon + \|\bar{u}\|)$  for all  
 209  $w \in S_\alpha(u, g_v) \subset T(u, g_v)$ .  $\square$

210 As stated, minimizing the mappings in Lemma 2.4 in the first variable yields, for all  
 211  $(u, g_v) \in \mathbb{R}^m \times G_v(\bar{x})$ , the bivariate  $\mathcal{U}$ -objects in Definition 2.3:

$$212 \quad L_U(u, g_v) = \inf_w \Phi(w, u, g_v) \quad \text{and} \quad W(u, g_v) = \arg \min_w \Phi(w, u, g_v).$$

213 Thanks to Lemma 2.4, several important relations known for the single-variable  $\mathcal{U}$ -Lagrangian  
 214 hold in our new bivariate context.

215 THEOREM 2.5 (Smoothness of bivariate  $\mathcal{U}$ -objects). *Given  $G_v(\bar{x})$  from (2.4), the bivari-*  
 216 *ate  $\mathcal{U}$ -Lagrangian and the  $\mathcal{V}$ -space minimizer set from Definition 2.3 satisfy the following*  
 217 *properties.*

- 218 1.  $L_U$  is finite-valued on  $\mathbb{R}^m \times G_v(\bar{x})$ ;
- 219 2.  $W$  is outer semi-continuous and locally bounded on  $\mathbb{R}^m \times G_v(\bar{x})$ ;
- 220 3.  $W(0, g_v) = \{0\}$  and  $W$  is continuous at  $(0, g_v)$  for any  $g_v \in G_v(\bar{x})$ ;
- 221 4.  $L_U$  is locally Lipschitz continuous on the interior of  $\mathbb{R}^m \times G_v(\bar{x})$ ;
- 222 5.  $L_U$  is differentiable at  $(0, g_v)$  for any  $g_v$  in the interior of  $G_v(\bar{x})$ , with  
 223  $\nabla L_U(0, g_v) = (\bar{g}_u, 0)$ ;
- 224 6. For all  $(u, g_v) \in \mathbb{R}^m \times V^\dagger \text{ri } \partial f(\bar{x})$ , and  $w$  an arbitrary point in  $W(u, g_v)$ ,

$$225 \quad \partial_u L_U(u, g_v) = \{s_u : s \in \partial f(\bar{x} + Uu + Vw), s_v = g_v\}.$$

226 *Proof.* For notational convenience, we define  $G := G_v(\bar{x})$ . All the subsequent references  
 227 in this proof are from the book [41], noting that the assumptions in the invoked statements hold  
 228 thanks to Lemma 2.4. To see item (i), apply first Theorem 1.17(a), to show that  $L_U$  is proper  
 229 and lsc on  $\mathbb{R}^n$ . For each  $(u, g_v) \in \mathbb{R}^m \times G$ ,  $L_U(u, g_v) \leq \Phi(0, u, g_v) = f(\bar{x} + Uu) < +\infty$ .

230 Consequently,  $L_U$  is finite-valued on  $\mathbb{R}^m \times G$ . From Theorem 1.17(c), we have that  $L_U$  is  
 231 continuous on  $\mathbb{R}^m \times G$ , as for any  $\bar{w} \in W(\bar{u}, \bar{g}_v)$  the function  $\Phi(\bar{w}, \cdot)$  is continuous in  $(u, g_v)$   
 232 at  $(\bar{u}, \bar{g}_v)$  relative to  $\mathbb{R}^m \times G$ . Consequently, item (ii) follows from Theorem 7.41(b). Item  
 233 (iii) is derived from Example 5.22, by exhibiting a point  $(\bar{u}, \bar{g}_v) \in \mathbb{R}^m \times G$  such that  $W(\bar{u}, \bar{g}_v)$   
 234 is single valued with  $L_U(\bar{u}, \bar{g}_v)$  continuous relative to  $\mathbb{R}^m \times G$ . The latter condition is ensured  
 235 by item (i) while the former is achieved by taking  $\bar{u} = 0$  and  $\bar{g}_v \in G$  and applying (2.2), noting  
 236 that  $W(0, \bar{g}_v) = W^{\bar{g}}(0) = \{0\}$ . To show the next item, according to Corollary 10.14(a), we  
 237 need to verify that for any  $(u, g_v) \in \text{int } \mathbb{R}^m \times G$ , it holds that

$$238 \quad (2.5) \quad \bigcup_{w \in W(u, g_v)} \{(s^1, s^2) : (0, s^1, s^2) \in \partial^\infty \Phi(w, u, g_v)\} = \{(0, 0)\} .$$

239 To this end, define  $\Phi = h_1 + h_2$  for the following two functions from  $\mathbb{R}^{n-m} \times \mathbb{R}^m \times \mathbb{R}^{n-m}$  to  $\bar{\mathbb{R}}$ :

$$240 \quad h_1(w, u, g_v) := f(\bar{x} + Uu + Vw) + \delta_G(g_v) \quad \text{and} \quad h_2(w, u, g_v) := -\langle g_v, V^T V w \rangle + \delta_G(g_v) .$$

241 Since  $h_1$  is proper, lsc and convex and  $h_2$  is strictly differentiable at  $(w, u, g_v)$  for any  
 242  $g_v \in \text{int } G$ , applying Exercises 10.10 and 10.7, we obtain that

$$243 \quad \partial^\infty \Phi(w, u, g_v) = \partial^\infty h_1(w, u, g_v) = \{(V^T s, U^T s, 0) : s \in \partial^\infty f(\bar{x} + Uu + Vw)\} = \{(0, 0, 0)\} ,$$

244 where the last equality holds because the horizon subdifferential of the finite-valued convex  
 245 function  $f$  is null. As claimed, (2.5) holds and (iv) follows. Item (v) is derived from the result  
 246  $W(0, g_v) = \{0\}$ , item (iv), and Corollary 10.14(b), by showing that

$$247 \quad (2.6) \quad \{(y^1, y^2) : (0, y^1, y^2) \in \partial \Phi(0, 0, g_v)\} = \{(\bar{g}_u, 0)\} .$$

248 Once more, from Exercises 10.10 and 10.7, we get that

$$249 \quad \partial \Phi(0, 0, g_v) = \{(V^T y - V^T V g_v, U^T y, 0) : y \in \partial f(\bar{x})\} .$$

250 The expression in (2.6) then follows from (2.1), concluding the proof of item (v). Finally,  
 251 to see item (vi), fix any  $g_v \in V^\dagger \text{ri } \partial f(\bar{x})$ . For each  $u \in \mathbb{R}^m$ , the partial subdifferential  
 252  $\partial_u L_U(u, g_v)$  is defined to be  $\partial L_U^g(u)$ . For each  $g_v \in V^\dagger \text{ri } \partial f(\bar{x})$ ,  $L_U^g(u) = \inf_w \Phi_{g_v}(w, u)$ ,  
 253 where  $\Phi_{g_v}(w, u) := f(\bar{x} + Uu + Vw) - \langle g_v, V^T V w \rangle$ . In view of Lemma 2.4,  $\Phi_{g_v}(w, u)$   
 254 is level-bounded in  $w$  locally uniformly in  $u$  because in this case  $g_v$  is a parameter and  
 255 (2.3) holds for some  $\eta_{g_v}$ . Noting that  $L_U^g(u)$  is convex, we can apply Corollary 10.13  
 256 to obtain that  $\partial L_U^g(u) = \{y \in \mathbb{R}^m : (0, y) \in \partial \Phi_{g_v}(w, u)\}$  for any  $w \in W^{g_v}(u)$ . It is next  
 257 seen that  $\partial \Phi_{g_v}(w, u) = \{(V^T s - V^T V g_v, U^T s) : s \in \partial f(\bar{x} + Uu + Vw)\}$ . We also have that  
 258  $0 \in V^T s - V^T V g_v$  equivalent to  $(V^T V)^{-1} V^T s = g_v$ . Consequently, item (vi) holds.  $\square$

259 All the properties listed in (2.2), with only  $u$  considered a variable, can now be compared  
 260 with the statements in items (iii) and (iv) of Theorem 2.5, shown in the bivariate setting.

261 **2.3.  $\mathcal{V}$ -minimizers and proximal points.** When the special trajectory associated with  
 262  $\mathcal{V}$ -minimizers is identified, the function appears smooth along the resulting  $\mathcal{U}$ -subspace, and  
 263 a  $\mathcal{U}$ -Newton step is possible. For a  $\mathcal{V}\mathcal{U}$ -method to be superlinearly convergent, the fast  
 264  $\mathcal{U}$ -step should dominate over the  $\mathcal{V}$ -step. In this respect, the behavior of  $\mathcal{V}$ -minimizers in  
 265 the set  $W(u, g_v)$  from Definition 2.3 is crucial.

266 For the original  $\mathcal{U}$ -Lagrangian, [22, Corollary 3.5] shows that  $\mathcal{V}$ -minimizers are *tangent*  
 267 to the  $\mathcal{U}$ -subspace. The same important result holds for our new bivariate  $\mathcal{U}$ -Lagrangian, as  
 268 established next.

269 LEMMA 2.6 (Tangential trajectories). *Let  $\bar{g} \in \text{ri } \partial f(\bar{x})$ ,  $\bar{x} \in \mathbb{R}^n$ . With the notation and*  
 270 *assumptions of Lemma 2.4, we have that*

$$271 \quad (2.7) \quad L_U(u, g_v) = f(\bar{x}) + \langle \bar{g}_u, u \rangle + o\left(Uu + V(g_v - \bar{g}_v)\right).$$

272 *Proof.* By item (v) in Theorem 2.5, the function  $L_U$  is differentiable at  $(0, \bar{g}_v)$ , and hence,

$$273 \quad L_U(u, g_v) = L_U(0, \bar{g}_v) + \langle \nabla L_U(0, \bar{g}_v), (u, g_v) - (0, \bar{g}_v) \rangle + o\left(Uu + V(g_v - \bar{g}_v)\right).$$

274 Substituting  $L_U(0, \bar{g}_v)$  and its gradient by their explicit expressions gives (2.7).  $\square$

275 Together with the  $\mathcal{U}$ -Lagrangian given in Definition 2.3, from (2.7) we obtain the follow-  
 276 ing first-order expansion for  $f$ :

$$277 \quad f(\bar{x} + Uu + Vw(u, g_v)) = f(\bar{x}) + \langle \bar{g}_u, u \rangle + \langle g_v, V^T V w(u, g_v) \rangle + o\left(Uu + Vg_v\right).$$

278 When compared with the relation given in Section 2.1 for the single-variable setting, we see  
 279 that for fast convergence purposes, not only the  $\mathcal{U}$ -component ( $u$ ) should dominate eventually,  
 280 but also the  $\mathcal{V}$ -component of the subgradient ( $g_v$ ) should vanish. In Algorithm 3.1 below,  
 281 this is achieved by adding the  $\mathcal{U}$ -Newton step to the proximal gradient update.

282 Having highlighted the importance of  $\mathcal{V}$ -minimizers, we now show that they can be  
 283 identified by means of the proximal point mapping. For any function  $h: \mathbb{R}^n \rightarrow \mathbb{R}$  and a real  
 284 number  $\mu > 0$ , recall that the proximal mapping is given by

$$285 \quad \text{prox}_{h, \mu}(x) := \arg \min_{y \in \mathbb{R}^n} \left\{ h(y) + \frac{\mu}{2} \|y - x\|^2 \right\}.$$

286 Combined with Lemma 2.6, the following result provides a mechanism that makes the  $\mathcal{V}$ -step  
 287 be tangential to the  $\mathcal{U}$ -subspace.

288 LEMMA 2.7 (Characterization of  $\mathcal{V}$ -minimizers). *Let  $g \in \text{ri } \partial f(\bar{x})$ ,  $\bar{x} \in \mathbb{R}^n$ . For any*  
 289  *$p \in \mathbb{R}^n$  and the corresponding  $\mathcal{V}\mathcal{U}$ -components  $u(p) := (p - \bar{x})_u$  and  $v(p) := (p - \bar{x})_v$ , it holds*  
 290 *that  $v(p) \in W(u(p), g_v)$  if and only if  $g_v \in V^\dagger \partial f(p)$ , in which case  $g_u \in \partial_u L_U(u(p), g_v)$ .*  
 291 *If, in addition, there is  $g' \in \partial f(p)$  such that  $g'_v = g_v$ , then  $g'_u \in \partial_u L_U(u(p), g_v)$ .*

292 *Proof.* The convex function  $\mathbb{R}^{n-m} \ni v \mapsto h(v) := f(\bar{x} + Uu(p) + Vv)$  has the subdif-  
 293 ferential  $\partial h(v) = V^T \partial f(\bar{x} + Uu(p) + Vv)$ . The necessary and sufficient optimality condition  
 294 for  $v(p) \in W(u(p); g_v)$  is  $0 \in \partial h(v(p)) - V^T V g_v = V^T \partial f(\bar{x} + Uu(p) + Vv(p)) - V^T V g_v =$   
 295  $V^T \partial f(p) - V^T V g_v$ , and the equivalence follows from the definition of the pseudo-inverse  $V^\dagger$ .  
 296 To show that  $g_u \in \partial_u L_U(u(p), g_v)$ , note that if  $g' \in \partial f(p)$  and  $g'_v = g_v$ , then  $g_v \in V^\dagger \partial f(p)$ .  
 297 Hence, the expression of  $\partial_u L_U(\cdot, g_v)$  in Theorem 2.5 is verified by  $g'_u$ .  $\square$

298 Proximal points are related to  $\mathcal{V}$ -minimizers through Lemma 2.7, by taking the subgradient  
 299 in the optimality condition of the proximal point problem. Specifically, for given  $z \in \mathbb{R}^n$   
 300 and  $\mu > 0$ , the result is applied with  $p = \text{prox}_{f, \mu}(z)$  and  $g' = \mu(z - p)$ ; see Theorem 3.1(i)  
 301 below. In the algorithm given in next section, however, the  $\mathcal{V}$ -step does not compute exact  
 302 proximal points. Rather, having a model function for  $f$ , the proximal point of the model is  
 303 computed. By exploiting structural properties of the function to be minimized, given in (3.1)  
 304 below, we can rewrite the model proximal point as an *exact* proximal point of  $f$ , by shifting  
 305 the prox-center; see Theorem 3.1. Thanks to this shifting, the result in Lemma 2.7 applies.



306 **3. The algorithm and its global convergence.** We next focus our attention on the  
 307 function in (1.1) having the following additive structure:

$$\begin{array}{ll}
 f(x) \equiv q(x) + h(x) & \\
 \text{for } q: \mathbb{R}^n \rightarrow \mathbb{R} & \text{convex, } \mathcal{C}^2\text{-smooth} \\
 & \text{with gradient Lipschitz constant denoted by } \beta \\
 h: \mathbb{R}^n \rightarrow \mathbb{R} & \text{a continuous convex function, possibly nonsmooth,} \\
 & \text{with an easy-to-compute proximal point.}
 \end{array}
 \tag{3.1}$$

309 In the considered context, for the family of model functions

$$310 \tag{3.2} \quad m(x; y) := q(y) + \langle \nabla q(y), x - y \rangle + h(x),$$

311 it is easy to compute the proximal point of  $m(\cdot; y)$ , for any parameter  $y \in \mathbb{R}^n$ .

312 **3.1.  $\mathcal{V}$ -step: proximal gradient iterations.** To minimize a function  $f$  as in (3.1), the  
 313 well-known proximal gradient algorithm [37, 5] computes the proximal point of the model  
 314 (3.2). At iteration  $k$ , given the current iterate  $x^k$  and a prox-parameter  $\mu_k$ , the next point is  
 315  $x^{k+1} = \text{prox}_{m(\cdot; x^k), \mu_k}(x^k)$ . If  $\mu_k \geq \beta$  and  $f$  satisfies an error bound, the proximal gradient  
 316 iterates converge with linear rate [9, Theorems 3.1 and 5.5]; see also [1]. To achieve superlinear  
 317 speed, in our method those iterations are corrected by a suitable  $\mathcal{U}$ -step; see Algorithm 3.1  
 318 below. The proximal gradient iteration in Procedure 1 corresponds to our  $\mathcal{V}$ -step.

---

**Procedure 1:** Proximal Gradient(the  $\mathcal{V}$ -) step

---

**Input:**  $f$  as in (3.1),  $x^k \in \mathbb{R}^n$ ,  $\mu_k > 0$ , and  $p = \text{prox}_{m(\cdot; x^k), \mu_k}(x^k)$  for the model  
 (3.2).

**while**  $q(p) > q(x^k) + \langle \nabla q(x^k), p - x^k \rangle + \frac{\mu_k}{2} \|p - x^k\|^2$  **do**  
 | declare a null step: set  $\mu_k := 2\mu_k$   
 | compute  $p = \text{prox}_{m(\cdot; x^k), \mu_k}(x^k)$

**end**

**Output:**  $\mu_k$  and  $p^k = p$ .

---

319 Some comments regarding Procedure 1 are in order. The output  $p^k$  satisfies  $q(p^k) \leq$   
 320  $q(x^k) + \langle \nabla q(x^k), p^k - x^k \rangle + \frac{\mu_k}{2} \|p^k - x^k\|^2$ . With our assumptions in (3.1), this ensures that,  
 321 once  $\mu_k \geq \beta$ , the procedure will terminate (in the bundle methods terminology [18], the  
 322 sequence of null steps is always finite). Additionally, note that if in (3.1) there is no smooth  
 323 term, then  $q \equiv 0$  and  $h \equiv f$ . Since in this case the model is the same function (assuming the  
 324 prox-calculation of  $f$  is easy), the  $\mathcal{V}$ -step performs an exact proximal step for  $f$  and there are  
 325 no null steps.

326 The procedure output is the proximal point of the model (3.2) at  $x^k$ . In order to apply  
 327 Lemma 2.7, and in this way identify the output with a  $\mathcal{V}$ -minimizer,  $p^k$  needs to be the  
 328 proximal point of the function, and not of its model. This is shown in our next result, where  
 329 we exhibit  $p^k$  to be the *exact* proximal point of  $f$  at a certain shifted point.

330 **THEOREM 3.1** (Shifting proximal point of the model to exact proximal point of the func-  
 331 tion).

332 *Given the output  $p^k$  of Procedure 1, define*

$$333 \tag{3.3} \quad g^k := \mu_k(x^k - p^k) + \nabla q(p^k) - \nabla q(x^k) \quad \text{and} \quad z^k := p^k + \frac{1}{\mu_k} g^k.$$

334 Then it holds that

$$335 \quad (3.4) \quad p^k = \text{prox}_{f, \mu_k}(z^k), \quad z^k = x^k + \frac{1}{\mu_k}(\nabla q(p^k) - \nabla q(x^k)), \quad g^k \in \partial f(p^k).$$

336 Therefore, for  $\bar{x} \in S$  a minimizer of  $f$  in (1.1), the following holds.

- 337 (i) Setting  $u^k := (p^k - \bar{x})_u$ , the corresponding  $\mathcal{V}$ -component  $v^k := (p^k - \bar{x})_v \in$   
 338  $W(u^k, g_v^k)$  if and only if  $g_v^k \in V^\dagger \text{ri } \partial f(\bar{x})$ , in which case  $g_u^k \in \partial_u L_U(u^k, g_v^k)$ .  
 339 (ii) Furthermore, whenever  $\mu_k > \beta$ , it holds that

$$340 \quad \|p^k - \bar{x}\| \leq \frac{\mu + \beta}{\mu - \beta} \|x^k - \bar{x}\| \quad \text{and} \quad \|g^k\| \leq \frac{2\mu(\mu + \beta)}{\mu - \beta} \|x^k - \bar{x}\|.$$

341 *Proof.* We have that  $\mu_k(x^k - p^k) \in \partial m(p^k; x^k) = \nabla q(x^k) + \partial h(p^k)$ , which yields  
 342  $\mu_k(x^k - p^k) - \nabla q(x^k) \in \partial h(p^k)$ . Therefore,  $\mu_k(x^k - p^k) + \nabla q(p^k) - \nabla q(x^k) \in \nabla q(p^k) +$   
 343  $\partial h(p^k) = \partial f(p^k)$ . The remaining assertions in (3.4) follow from the optimality condition  
 344  $0 \in \partial f(p^k) + \mu_k(p^k - z^k)$ , i.e.,  $\mu_k(z^k - p^k) \in \partial f(p^k)$  which by (3.3) is just  $g^k \in \partial f(p^k)$ .  
 345 Item (i) follows from (3.4) and Lemma 2.7, written with  $p, g'$  therein replaced by  $p^k, g^k$ .

346 For the final item, first note that, by (3.3), it holds that

$$347 \quad (3.5) \quad \|g^k\| \leq (\mu_k + \beta) \|x^k - p^k\| \leq (\mu_k + \beta) (\|x^k - \bar{x}\| + \|\bar{x} - p^k\|),$$

348 Next, by (3.4) and the nonexpansiveness of the proximal operator,

$$349 \quad \|\bar{x} - p^k\| = \|\text{prox}_{f, \mu_k}(\bar{x}) - \text{prox}_{f, \mu_k}(z^k)\| \leq \|\bar{x} - z^k\|.$$

350 Using the expression for  $z^k$  in (3.4) and the bound for  $\nabla q$  from (3.1), we obtain that

$$351 \quad \|\bar{x} - p^k\| \leq \|\bar{x} - x^k\| + \frac{1}{\mu_k} \|\nabla q(p^k) - \nabla q(x^k)\| \leq \|\bar{x} - x^k\| + \frac{\beta}{\mu_k} \|p^k - x^k\|.$$

352 Adding  $0 = \pm \bar{x}$  in the right-most term, gives

$$353 \quad \|\bar{x} - p^k\| \leq \|\bar{x} - x^k\| + \frac{\beta}{\mu_k} \|p^k - \bar{x}\| + \frac{\beta}{\mu_k} \|\bar{x} - x^k\|.$$

354 After some rearrangements of terms we obtain that

$$355 \quad \|\bar{x} - p^k\| \leq \frac{1 + \frac{\beta}{\mu_k}}{1 - \frac{\beta}{\mu_k}} \|\bar{x} - x^k\| = \frac{\mu_k + \beta}{\mu_k - \beta} \|\bar{x} - x^k\|.$$

356 The last inequality in (3.5) yields the final result.  $\square$

357 The explicit shifting in Theorem 3.1 is possible thanks to the structure of  $f$  in (3.1).  
 358 Note that the tangential property depends on  $g_v^k$  eventually becoming (the  $V^\dagger$  component of)  
 359 an interior subgradient at a minimizer. To achieve this, the  $\mathcal{VU}$ -algorithm drives to zero  
 360 the subgradient  $g^k$  from (3.3). Accordingly, the Proximal Gradient  $\mathcal{VU}$ -method given in  
 361 Algorithm 3.1, stops when  $\|g^k\| \leq \text{TOL}$  for a given tolerance  $\text{TOL}$ .

362 **3.2.  $\mathcal{U}$ -step and the algorithm statement.** After the  $\mathcal{V}$ -step is done, the output of  
 363 Procedure 1 is corrected as follows:

$$364 \quad (3.6) \quad x^{k+1} = p^k - U_k Q_k U_k^\top g^k,$$

365 where  $U_k$  is a certain orthonormal matrix and  $Q_k$  is positive semidefinite. The purpose of this  
 366 correction is to eventually track a trajectory where the function behaves smoothly, through the  
 367 relation of  $\mathcal{V}$ -minimizers with proximal points. Thus, in (3.6)  $g^k$  is the shifted gradient from  
 368 (3.3),  $Q^k$  asymptotically approximates the so-called  $\mathcal{U}$ -Hessian (a second-order object related  
 369 to the  $\mathcal{U}$ -Lagrangian defined in Section 4), and the orthonormal matrix  $U_k$  is a basis for a  
 370 subspace  $\mathcal{U}_k$  that approximates  $\mathcal{U}(p^k)$ . For the latter, see [7], and also Section 5 concerning  
 371 the  $\ell_1$ -regularized setting which is our illustration in this paper.

372 The full proximal gradient  $\mathcal{V}\mathcal{U}$ -method is given in Algorithm 3.1, where the stopping  
 373 criterion is justified by the fact that  $g^k \rightarrow 0$ , shown in Theorem 3.3.

---

**Algorithm 3.1** Proximal Gradient  $\mathcal{V}\mathcal{U}$ -method (PGVU)

---

**Data:**  $f$  as in (3.1), starting point  $x^0 \in \mathbb{R}^n$ , prox-parameter  $\mu_0 > 0$ , and a stopping tolerance  $\text{TOL} \geq 0$ . Set  $k = 0$ .

**repeat**

- Obtain  $\mu_k$  and  $p^k$  from Procedure 1.
- Define (shift the subgradient)  $g^k := \mu_k(x^k - p^k) + \nabla q(p^k) - \nabla q(x^k)$ .
- Compute an orthonormal basis  $U_k \in \mathbb{R}^{n \times n_k}$ .
- Choose a symmetric positive semidefinite matrix  $Q_k \in \mathbb{R}^{n_k \times n_k}$ .
- Update  $x^{k+1} = p^k - U_k Q_k U_k^\top g^k$ , set  $k = k + 1$

**until**  $\|g^k\| \leq \text{TOL}$ ;

---

374 When  $f$  is differentiable,  $\mathcal{U}(p^k)$  is the whole space, its basis is the identity matrix, and  
 375  $Q_k$  can be defined as usual for quasi-Newton methods; see, e.g., [19]. Otherwise, the structural  
 376 properties for  $f$  are essential to define suitable matrices in (3.6). Namely, as  $\partial f = \nabla q + \partial h$ ,  
 377 in Definition 2.1 the  $\mathcal{V}$ -subspaces of  $f$ ,  $h$ , and the model  $m$  in (3.2) are all identical:

$$378 \quad \mathcal{V}(p) := \mathcal{V}f(p) = \mathcal{V}h(p) = \mathcal{V}m(p; y) \quad \text{for all } p, y \in \mathbb{R}^n.$$

379 Hence, also  $\mathcal{U}(p) := \mathcal{U}f(p) = \mathcal{U}h(p) = \mathcal{U}m(p; y)$ .

380 To give an insight/illustration, we consider again our function from Example 2.2, and  
 381 compare the performance of three methods, according to the possible choices.

382 **EXAMPLE 3.2** (Proximal, proximal gradient and proximal gradient VU algorithms). For  
 383 the function  $F$  from Example 2.2, the smooth function in (3.1) is  $q(u, v) = \frac{a}{2}u^2$  and the  
 384 Lipschitz constant is  $\beta = a = 2$ .

385 We consider minimizing  $F$  with the proximal point method (P), the proximal gradient  
 386 algorithm (PG), and PGVU as given in Algorithm 3.1, with the following specifications:

- 387 • The implementation of both P and PG follows Procedure 1, with respective model  
 388 functions  $F(\cdot)$  and  $m(\cdot, x^k)$  from (3.2), and stopping test  $\max(\|x^k\|, \|p^k\|) \leq \text{TOL}$ .
- 389 • For PGVU, there is the additional  $\mathcal{U}$ -step, requiring the matrices in (3.6).
  - 390 – Outside of the fast track, that is when  $p^k \notin \{(p_1, p_2) : |p_2| = \frac{a}{2}(p_1)^2\}$ , the  
 391 function is differentiable,  $\mathcal{U}(p^k)$  is the whole space, and  $U_k$  is just the identity  
 392 matrix. The matrix  $Q_k$  is set to  $\begin{bmatrix} a & 0 \\ 0 & 0 \end{bmatrix}^\dagger$  if  $|p_2^k| < \frac{a}{2}(p_1^k)^2$ , and  $Q_k = 0$  if  
 393  $|p_2^k| > \frac{a}{2}(p_1^k)^2$ .
  - 394 – When  $p^k$  is on the fast track, its  $\mathcal{V}(p^k)$  subspace is generated by  $\partial F(p^k) =$   
 395  $\{(\xi a p_1^k, 1 - \xi)^\top : \xi \in [0, 1]\}$ . In this case,  $U_k$  is the orthonormal vector gen-  
 396 erated from  $(\text{sign}(p_1^k), \text{sign}(p_2^k) a p_1^k)$  and  $Q_k = a^{-1}$ .

397 The methods were run with the same initial  $\mu_0$  and  $x^0$ , until the termination tolerance  
 398  $\text{TOL} = 10^{-10}$  or a maximum number of iterations set to 50 achieved. Table 3.1 reports the

399 number of iterations and accuracy of the three algorithms, for two different initial values of  
400  $\mu_0$  and  $x^0$ .

Table 3.1: Methods' performance

$\mu^0 = 0.18$	$x^0 = (-1.2, 2.2)^\top$			$x^0 = (-1.1, 0.9)^\top$		
	P	PG	PGVU	P	PG	PGVU
Iterations	10	16	11	8	14	11
Digits	9	11	11	7	10	10
$\mu^0 = 10$	$x^0 = (-1.2, 2.2)^\top$			$x^0 = (-1.1, 0.9)^\top$		
	P	PG	PGVU	P	PG	PGVU
Iterations	50	50	45	50	50	16
Digits	2	3	23	3	4	23

401 When  $\mu_0$  is small, the performances of P, PG, and PGVU are similar. This is explained  
402 by the “null-step” inner loop in Procedure 1, which makes all the three methods increase  $\mu_k$   
403 until a value larger than  $\beta = 2.0$  is attained. By contrast, the runs with the large value of  
404  $\mu_0$  are troublesome for P and PG. Procedure 1 always accepts the large value  $\mu_k = \mu_0 = 10$ ,  
405 there is no backtracking mechanism to reduce  $\mu_k$ . This is prejudicial for P and PG, and can  
406 be explained by observing the plots in Fig. 3.1, with the three trajectories generated for one  
407 of the starting points.

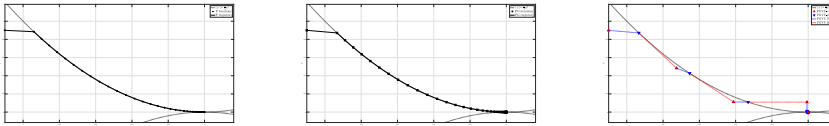


Fig. 3.1: P, PG, and PGVU iterations (left, middle, right), when minimizing  $F$  from  $x^0 = (-1.1, 0.9)^\top$ , starting with  $\mu_0 = 10$ . The dotted parabola is the fast track, rapidly identified by all the algorithms. Because  $\mu_k = \mu_0$  is too large, P and PG take very small steps along the fast track, which hinders their performance. By contrast, thanks to its  $\mathcal{U}$ -step, PGVU approaches the fast track tangentially and reaches rapidly a much better degree of accuracy.

408 Figure 3.1 highlights the following important point: *for a nonsmooth optimization method*  
409 *to achieve high accuracy, identifying the fast track is necessary, but it is not sufficient.* Both P  
410 and PG iterates land soon on the smooth trajectory defined by the fast track. This is consistent  
411 with the theory in [31]; see also [3, Theorem 3.1]. Once on the fast track, P and PG both  
412 remain there. But because  $\mu_k = \mu_0$  is too large, the progress between consecutive iterates  
413 becomes too small and, as illustrated by the left and middle plots in Fig. 3.1, those two methods  
414 see their performance severely impaired. On the other hand, on the right plot we see that the  
415 too large  $\mu_k = \mu_0$  also makes PGVU take a too small  $\mathcal{V}$ -step, but afterwards its  $\mathcal{U}$ -step is  
416 long, in the “right direction” to solution, but also driving the iterate far from the fast track.  
417 The remedy is that the subsequent  $\mathcal{V}$ -step takes the iterate to the fast track again, and a new  
418 fast  $\mathcal{U}$ -step is possible again. The overall iterative process can be interpreted to work in a  
419 “predictor-corrector” fashion. Observe that PGVU achieves a much higher accuracy, in less  
420 iterations. Noting, of course, that one PGVU iteration involves more computational work.

421 The conclusion is that when accuracy is important, PGVU is the right approach. If accuracy  
 422 is not a concern, simpler techniques would be preferred.

423 **3.3. Global convergence and some other asymptotic results.** In Theorem 3.3 below  
 424 about the global convergence result, we do not show that the sequence  $\{x^k\}$  is bounded (has  
 425 accumulation points). This is in part related to the comment above that in some naturally  
 426 included special cases which do not involve nonsmoothness, our method can actually reduce  
 427 to the usual smooth quasi-Newton technique. It is known that proving boundedness of the  
 428 usual quasi-Newton updates (say, BFGS) without very technical complex modifications to  
 429 the update, is not possible; see, e.g., [19]. Naturally, it is the same in our setting (also,  
 430 proving properties about accumulation points of algorithmic sequences is quite standard in  
 431 the literature, separately from their existence).

432 **THEOREM 3.3** (Global convergence). *In Algorithm 3.1, let the matrices  $\{U_k\}$  and  $\{Q_k\}$*   
 433 *be bounded, and the prox-parameters satisfy  $\mu_{\max} \geq \mu_k \geq \beta$  for all  $k$  sufficiently large.*

434 *Then every accumulation point  $\bar{x}$  of the sequence  $\{x^k\}$  generated by Algorithm 3.1 is a*  
 435 *minimizer of  $f$ .*

436 *Proof.* Let  $\bar{x}$  be an accumulation point of  $\{x^k\}$ , i.e., there exists a convergent subsequence  
 437 of  $\{x^k\}$ , indexed by  $k' \in \mathcal{K}$ , such that  $\{x^{k'}\}_{k' \in \mathcal{K}} \rightarrow \bar{x}$ . We claim that  $\{p^k - x^k\}_{k' \in \mathcal{K}} \rightarrow 0$   
 438 and  $\{g_{k' \in \mathcal{K}}^k\} \rightarrow \bar{g} = 0$ .

439 Let  $U$ ,  $Q$  and  $\bar{\mu}$  denote limit points of the corresponding subsequences of matrices  
 440 and prox-parameters. Define  $H = UQU^T$ . With the assumptions in (3.1), the models in  
 441 (3.2) converge continuously in the second argument:  $m(\cdot; x^k) \rightarrow m(\cdot; \bar{x})$ . Then, combining  
 442 [41, Thms. 7.11 and 12.35], the models epi-converge to  $m(\cdot; \bar{x})$ , and the proximal mappings  
 443  $\text{prox}_{m(\cdot; x^k), \mu_k}(\cdot)$  converge to  $\text{prox}_{m(\cdot; \bar{x}), \bar{\mu}}(\cdot)$  uniformly on bounded sets. By epi-convergence  
 444 of the models and [41, Theorem 7.14], the proximal point sequence converges continuously:  
 445  $p^k \rightarrow \text{prox}_{m(\cdot; \bar{x}), \bar{\mu}}(\bar{x}) =: \bar{p}$ . Passing onto the limit in (3.6) with  $g^k$  from (3.3),

$$446 \quad \bar{x} = \bar{p} + H(\bar{\mu}(\bar{p} - \bar{x}) + \nabla q(\bar{x}) - \nabla q(\bar{p})).$$

447 Since the Hessian  $\nabla^2 q(\cdot)$  is positive semidefinite by assumption, by the mean-value theorem  
 448  $\nabla q(\bar{x}) - \nabla q(\bar{p}) = \nabla^2 q(\bar{y})(\bar{x} - \bar{p})$  for some intermediate point  $\bar{y}$ . Then, after some direct  
 449 algebraic manipulations, we obtain that

$$450 \quad (I + \bar{\mu}H - H\nabla^2 q(\bar{y}))(\bar{p} - \bar{x}) = (I + H[\bar{\mu}I - \nabla^2 q(\bar{y})])(\bar{p} - \bar{x}) = 0.$$

451 Because  $\mu_k \geq \beta$  for large  $k$ , the matrix  $\bar{\mu}I - \nabla^2 q(\bar{y})$  is positive semidefinite and by positive  
 452 semidefiniteness of  $Q$ , the matrix  $H$  is positive semidefinite. As a result,  $I + H[\bar{\mu}I - \nabla^2 q(\bar{y})]$   
 453 is positive definite and, hence,  $\bar{p} = \bar{x}$ . The definition of  $g^k$  in (3.3) and the continuity of  $\nabla q$ ,  
 454 readily give  $g^k \rightarrow 0$ . By (3.4), this means that  $\bar{p} = \bar{x}$  is a minimizer of  $f$ , as stated.  $\square$

455 Thanks to the fact that  $g^k \rightarrow \bar{g} = 0$ , we are now in position to show that convergent  
 456 subsequences eventually generate  $\mathcal{V}$ -minimizers and identify smooth trajectories associated  
 457 with the  $\mathcal{U}$ -Lagrangian.

458 **COROLLARY 3.4** (Asymptotic  $\mathcal{U}$ -Lagrangian identification and rates). *Under the as-*  
 459 *sumptions in Theorem 3.3, suppose  $\bar{g} = 0 \in \text{ri } \partial f(\bar{x})$ . Then, for  $k \in \mathcal{K}$  sufficiently large and*  
 460  *$(u^k, v^k)$  defined in Theorem 3.1, the following holds.*

- 461 (i)  $g_u^k \in \partial_u L_{\mathcal{U}}(u^k, g_v^k)$ ;
- 462 (ii) *If, in addition,  $\mu_k > \beta$ , then  $Vv^k = o(x^k - \bar{x})$ .*

463 *Proof.* Throughout the proof,  $k \in \mathcal{K}$ . Notice that, by Theorem 3.3, the subsequence  
 464  $p^k \rightarrow \bar{p} = \bar{x}$ , and, thus, both  $u^k$  and  $g^k \rightarrow 0$ . By (2.4),  $g_v^k \in G_v(\bar{x}) \subset V^\dagger \text{ri } \partial f(\bar{x})$ , so

465 eventually  $g_v^k$  lies in the relative interior of  $V^\dagger \partial f(\bar{x})$ , and the statement in (i) corresponds to  
 466 Theorem 3.1(i), where it is also shown that  $v^k \in W(u^k, g_v^k)$ .

467 To show the final result recall that, since  $\|Uu^k\| \leq \|p^k - \bar{x}\|$  and  $\|Vg_v^k\| \leq \|g^k\|$ , by  
 468 Theorem 3.1(iii) we obtain that

$$469 \quad (3.7) \quad \|Uu^k + Vg_v^k\| \leq \frac{3\mu(\mu + \beta)}{\mu - \beta} \|x^k - \bar{x}\|.$$

470 The  $\mathcal{U}$ -Lagrangian Definition 2.3 combined with the expression (2.7) with  $\bar{g} = 0$ , yields

$$471 \quad f(\bar{x} + Uu^k + Vv^k) - \langle g_v^k, V^\top Vv^k \rangle = L_U(u^k, g_v^k) = f(\bar{x}) + o\left(Uu^k + Vg_v^k\right).$$

472 The inequality (2.3) written with  $w = v^k$  and  $\bar{g}_u = 0$  gives the lower bound  $f(\bar{x}) + \eta\|Vv^k\|$   
 473 for the left-hand side in the relations above. Hence,  $\eta\|Vv^k\| \leq o\left(Uu^k + Vg_v^k\right)$ , and (3.7)  
 474 concludes the proof.  $\square$

475 **4.  $\mathcal{U}$ -Hessians and superlinear convergence.** The  $\mathcal{V}$ -minimizers exhibit first-order  
 476 expansions for  $f$ . To proceed further, a generalized notion of a Hessian [17] is needed.

477 **4.1. A partial second-order object for the bivariate  $\mathcal{U}$ -Lagrangian.** The  $\mathcal{U}$ -Hessian  
 478 introduced in [22, § 3.3] for the single-variable  $\mathcal{U}$ -Lagrangian is the basis to define a partial  
 479  $\mathcal{U}$ -Hessian, obtained when differentiating the bivariate  $\mathcal{U}$ -Lagrangian in the first variable.

480 **DEFINITION 4.1** (partial  $\mathcal{U}$ -Hessian). *Given  $\bar{x}$  and  $\bar{g} \in \text{ri } \partial f(\bar{x})$ , we say that  $f$  has at  $\bar{x}$   
 481 a partial  $\mathcal{U}$ -Hessian  $H(\bar{x}; \bar{g}_v)$  associated with  $\bar{g}_v$  if*

$$482 \quad (4.1) \quad \partial_u L_U(u, \bar{g}_v + z) \subset \bar{g}_u + H(\bar{x}; \bar{g}_v)u + B^m(0, o(u, z)),$$

483 where  $\bar{g}_u$  is defined in (2.1) and  $z \in \mathbb{R}^{n-m}$  is such that  $\bar{g}_v + z \in V^\dagger \partial f(\bar{x})$ .

484 **LEMMA 4.2** (relation with single-variable  $\mathcal{U}$ -Hessian). *The partial  $\mathcal{U}$ -Hessian in (4.1) is  
 485 also the  $\mathcal{U}$ -Hessian associated with single-variable  $\mathcal{U}$ -Lagrangian, with  $\bar{g}$  being a parameter:*

$$486 \quad \partial L_U^{\bar{g}}(u) \subset \nabla L_U^{\bar{g}}(0) + H^{\bar{g}}(\bar{x})u + B^m(0, o(\|u\|)), \text{ where } H^{\bar{g}}(\bar{x}) = H(\bar{x}; \bar{g}_v).$$

487 Furthermore, if  $0 \in \text{ri } \partial f(\bar{x})$  and  $f$  has a partial  $\mathcal{U}$ -Hessian at  $\bar{x}$  associated with  $\bar{g}_v = 0$ ,  
 488 then the following holds for  $\bar{H} = H(\bar{x}, 0)$ .

- 489 (i) *Under the assumptions in Theorem 3.3, the shifted gradients in Algorithm 3.1 satisfy  
 490 the inclusion  $g_u^k \in \bar{H}u^k + o(u^k)$ ;*  
 491 (ii) *For small  $d \in \mathbb{R}^n$ ,  $f(\bar{x} + d) = f(\bar{x}) + \frac{1}{2} \langle \bar{H}d_u, d_u \rangle + o(\|d_u^2\|)$ . As a result, if  $\bar{H}$  is  
 492 positive definite,*

$$493 \quad \exists c > 0 : d \in \mathbb{R}^n \text{ small} \implies f(\bar{x} + d) \geq f(\bar{x}) + \frac{c}{2} \|d\|^2.$$

494 In particular,  $\bar{x}$  is the unique solution in (1.1).

495 *Proof.* For the identification with the single-variable  $\mathcal{U}$ -Hessian it suffices to recall (2.2)  
 496 and write (4.1) with  $z \equiv 0$ . Item (i) also follows from (4.1), because with our assumptions  
 497  $g_u^k \in \partial_u L_U(u^k, g_v^k)$  by Corollary 3.4(i). To show (ii), note that if  $0 \in \text{ri } \partial f(\bar{x})$  then  $\bar{g}_u = 0$  and  
 498 we can take  $\bar{g}_v = 0$ . By the identification between the partial and single-variable  $\mathcal{U}$ -Hessians,  
 499  $\bar{H} = H^{\bar{g}=0}(\bar{x})$ . Thus,  $f$  has a  $\mathcal{U}$ -Hessian at  $\bar{x}$ , and [22, Theorem 3.9] gives the second-order  
 500 expansion. If  $\bar{H}$  is positive definite, then [21, Corollary 1] gives the lower bound for all small  
 501  $d$ . Uniqueness of  $\bar{x}$ , called a strong minimizer in [35], follows from  $f$ 's convexity.  $\square$

502 The lower bound for  $f$  that holds when the partial  $\mathcal{U}$ -Hessian is positive definite is called  
 503 *local subdifferential convexity* in [9, Theorem 6.2]. The property is shown to be equivalent  
 504 to both tilt stability and to strong metric regularity of the subdifferential at  $\bar{x}$ , a stable strong  
 505 minimizer for (1.1); see also [8].

506 **4.2. The partial  $\mathcal{U}$ -Hessian of partly smooth functions.** When the function in (1.1) is  
 507 partly smooth [24, Definition 2.7], it is shown in [3] that Riemannian Newton-like methods can  
 508 be combined with proximal gradient steps to boost convergence speed. Under the assumption  
 509 of partial smoothness, we now study conditions for the existence of a partial  $\mathcal{U}$ -Hessian.

510 **DEFINITION 4.3.** *A convex function  $f$  is said to be partly smooth at  $x$  relative to a set  $\mathcal{M}$*   
 511 *if  $\mathcal{M}$  is a manifold around  $x$  and the following three properties hold:*

- 512 (i) *(restricted smoothness) in a neighbourhood  $X$  of  $x$ , the restriction of  $f$  to  $\mathcal{M}$ ,  $f|_{\mathcal{M} \cap X}$ ,*  
 513 *is of class  $C^2$  ;*
- 514 (ii) *(normals parallel to subdifferential)  $N_{\mathcal{M}}(x) = \mathcal{V}(x)$ ;*
- 515 (iii) *(subgradient continuity) the subdifferential  $\partial f$  is continuous at  $x$  relative to  $\mathcal{M}$ .*

516 The function  $F$  in Example 2.2 has the partial  $\mathcal{U}$ -Hessian  $\bar{H} = a > 0$ , corresponding to  
 517 the bivariate  $\mathcal{U}$ -Lagrangian  $L(u, \bar{g}_v) = (1 - |\bar{g}_v|) \frac{a^2}{2}$ . However,  $F$  is not partly smooth, because  
 518 near  $\bar{x} = 0$  there are two distinct activity manifolds  $\mathcal{M}$ . These are the two fast tracks displayed  
 519 in Figure 2.1, generated by the different  $\mathcal{V}$ -minimizers that emanate from taking a positive  
 520 or a negative  $\bar{g}_v$  in Definition 2.3, i.e.,  $W(u, \bar{g}_v) = \{\frac{a}{2} \text{sign}(\bar{g}_v) u^2\}$ . By contrast, the simple  
 521 modification of Example 2.2 given by  $\tilde{F}(u, v) = \frac{a}{2} u^2 + \max\{0, v - \frac{a}{2} u^2\}$ , is partly smooth  
 522 at  $\bar{x}$ . The fundamental difference is that the  $\mathcal{V}$ -minimizers of  $\tilde{F}$  are  $W(u, \bar{g}_v) = \{\frac{a}{2} u^2\}$ , *the*  
 523 *same* for all  $\bar{g}_v$ . Now the  $\mathcal{U}$ -Lagrangian  $L_U(u, \bar{g}_v) = (1 - \bar{g}_v) \frac{a}{2} u^2$  provides the single activity  
 524 manifold  $\mathcal{M} := \{(v, u) : v = \frac{a}{2} u^2\}$ , where the partial  $\mathcal{U}$ -Hessian is again  $\bar{H} = a > 0$ .

525 Our next result states a similar relation in the general setting, by connecting partial  
 526 smoothness and  $\mathcal{V}\mathcal{U}$ -analysis, thanks to [24, Theorem 6.1]. We associate the manifold of  
 527 partial smoothness with a special  $\mathcal{V}$ -minimizer that is a  $C^2$ -function and *is the same* for all  
 528 interior subgradients (the same function was considered in [34, Theorem 6] to characterize  
 529 fast tracks for prox-regular functions).

530 **THEOREM 4.4** (Special  $\mathcal{V}$ -minimizers from partial smoothness). *Let  $f$  be a convex*  
 531 *function that is partly smooth at the point  $\bar{x}$  relative to a non-singleton set  $\mathcal{M} \subset \mathbb{R}^n$ . Then,*  
 532 *for all small  $u$ , there exists a  $C^2$  function  $v_{\partial f}$  such that*

$$533 \quad \forall g \in \text{ri } \partial f(\bar{x}), v_{\partial f}(u) \in W(u, g_v), \text{ and } v_{\partial f}(u) = O(\|u\|^2).$$

534 *As a result, there exist a neighborhood  $X \subset \mathbb{R}^n$  of  $\bar{x}$ , a neighborhood  $Y \subset \mathbb{R}^m$  of 0 such that*

$$535 \quad (4.2) \quad L_U(u, g_v) = f|_{\mathcal{M} \cap X}(u) - \langle g_v, V^T V v_{\partial f}(u) \rangle .$$

$$536 \quad (4.3) \quad \mathcal{M} \cap X = \{\bar{x} + Uu + Vv_{\partial f}(u) : u \in Y\} ,$$

537 *where  $f|_{\mathcal{M} \cap X}$  is considered a composite function of  $u$ .*

538 *Proof.* From the property (ii) in Definition 4.3, when the set  $\mathcal{M}$  is not a singleton, the  
 539 subspaces tangent and normal to the manifold at  $\bar{x}$  coincide respectively with  $\mathcal{U}$  and  $\mathcal{V}$ . Then,  
 540 by [24, Theorem 6.1], there exist a neighborhood  $X \subset \mathbb{R}^n$  of  $\bar{x}$ , a neighborhood  $Y \subset \mathbb{R}^m$  of 0,  
 541 and a function  $v_{\partial f} : \mathbb{R}^m \rightarrow \mathbb{R}^{n-m}$  such that for all  $u \in Y$ ,

$$542 \quad v_{\partial f}(u) \text{ is of class } C^2, v_{\partial f}(u) = O(\|u\|^2), \text{ and } \mathcal{M} \cap X = \{\bar{x} + Uu + Vv_{\partial f}(u) : u \in Y\} .$$

543 From the last relation, the restriction in Definition 4.3(i) has the expression  $f|_{\mathcal{M} \cap X} = f(\bar{x} +$   
 544  $Uu + Vv_{\partial f}(u))$ . Again by [24, Theorem 6.1], for all  $g \in \text{ri } \partial f(\bar{x})$ , the function  $h^g(w) :=$   
 545  $f(\bar{x} + Uu + Vw) - \langle g, \bar{x} + Uu + Vw \rangle$  has  $v_{\partial f}$  as a sharp minimizer. The identity (4.2) follows  
 546 because  $\langle g_v, Vw \rangle = \langle g, \bar{x} + Uu + Vw \rangle$ , which shows that  $h^g(w)$  is the minimand defining the  
 547 bivariate  $\mathcal{U}$ -Lagrangian.  $\square$

548 When specializing Theorem 4.4 to the setting (3.1), it is possible to derive an explicit  
 549 expression for the  $\mathcal{U}$ -Hessian when the nonsmooth function in (3.1) is polyhedral, as in  
 550 regularized regression problems.

551 **COROLLARY 4.5** (explicit partial  $\mathcal{U}$ -Hessian). *When, under the assumptions in Theo-*  
 552 *rem 4.4,  $0 \in \text{ri } \partial f(\bar{x})$ , a partial  $\mathcal{U}$ -Hessian of  $f$  at  $\bar{x}$  associated with  $\bar{g} = 0$  is given by the*  
 553 *Hessian restriction:  $\bar{H} = \nabla_{uu}^2 f|_{\mathcal{M} \cap X}(0)$ .*

554 *If, in addition, in (3.1) the nonsmooth function  $h$  is finite-valued and polyhedral, then*  
 555  *$\bar{H} = U^\top \nabla^2 q(\bar{x}) U$ .*

556 *Proof.* In (4.2), the linear term defines the bivariate function on  $Y \times V^\dagger \text{ri } \partial f(\bar{x})$  given by

$$557 \quad P(u, g_v) := \langle g_v, V^\top V v(u) \rangle = g^\top V v_{\partial f}(u),$$

558 by definition of the  $\mathcal{V}$ -components. The Jacobian of this function on the first component is

$$559 \quad (4.4) \quad \mathcal{J}_u P(u, g_v) = \mathcal{J} v_{\partial f}(u)^\top V^\top g = \mathcal{J} v_{\partial f}(u)^\top V^\top V g_v.$$

560 Since  $W(u, g_v) \ni v_{\partial f}(u) = O(\|u\|^2)$  by Theorem 4.4, combining Theorem 2.5(iii) and the  
 561 fact that

$$562 \quad \mathcal{J} v(u) = \mathcal{J} v(0) + O(u) = O(u),$$

563 gives the following expansion for the gradient of the  $\mathcal{U}$ -Lagrangian from (4.2):

$$564 \quad \nabla_u L_U(u, g_v + z) = \nabla_u L_U(0, g_v) + \nabla_{uu}^2 L_U(0, g_v) u + \frac{\partial^2}{\partial u \partial g_v} L_U(0, g_v) z + o(u, z),$$

565 for any  $z \in \mathbb{R}^{n-m}$  small enough such that  $g_v + z \in V^\dagger \text{ri } \partial f(\bar{x})$ . In this expansion, by (4.4) and  
 566 (4.2), the cross-derivative has the form

$$567 \quad \frac{\partial^2}{\partial u \partial g_v} L_U(0, g_v) = -\frac{\partial^2}{\partial u \partial g_v} P(0, g_v) = -\mathcal{J} v_{\partial f}(0)^\top V^\top V = 0.$$

568 Recalling that  $\nabla_u L_U(0, g_v) = \bar{g}_u$ , gives the desired expression for the partial  $\mathcal{U}$ -Hessian. In  
 569 view of (4.4),

$$570 \quad (4.5) \quad \mathcal{J}_u P(u, g_v) = o(u, g_v).$$

571 In particular,  $\mathcal{J}_u P(0, g_v) = 0$ . Hence, by (4.2) and the fact that  $f|_{\mathcal{M} \cap X}$  is a composite function  
 572 of  $u$ , we obtain that  $\nabla_u f|_{\mathcal{M} \cap X}(0) = \nabla_u L_U(0, g_v) + \nabla_u P(0, g_v) = \bar{g}_u$ . Then the smoothness of  
 573  $f|_{\mathcal{M} \cap X}$  yields  $\nabla_u f|_{\mathcal{M} \cap X}(u) = \bar{g}_u + \nabla_{uu}^2 f|_{\mathcal{M} \cap X}(0) u + o(u)$ . Consequently, by (4.2) and (4.5),

$$574 \quad \nabla_u L_U(u, g_v) = \nabla_u f|_{\mathcal{M} \cap X}(u) - \nabla_u P(u, g_v) = \bar{g}_u + \nabla_{uu}^2 f|_{\mathcal{M} \cap X}(0) u + o(u) - o(u, g_v).$$

575 As  $o(u) - o(u, g_v) = o(u, g_v)$ , we see from Prop. 4.2 that  $\nabla_{uu}^2 f|_{\mathcal{M} \cap X}(0)$  is a partial  $\mathcal{U}$ -Hessian  
 576 of  $f$  at  $\bar{x}$  associated with 0.

577 Now consider the special setting of  $f$  in (3.1), with  $h$  finite-valued polyhedral, so that

$$578 \quad h(x) = \max_{i \in I} \{ \langle a^i, x \rangle + b^i \} \quad \text{for some finite index set } I \neq \emptyset.$$

579 Then, for the ‘‘active’’ index set  $I(x) = \{ i \in I : \langle a^i, x \rangle + b^i = f(x) \}$ ,

$$580 \quad (4.6) \quad \partial h(x) = \left\{ \sum_{i \in I(x)} \alpha_i a^i : \sum_{i \in I(x)} \alpha_i = 1, \alpha_i \geq 0 \ (i \in I(x)). \right\}$$



581 The function  $h$  is partly smooth at any  $\bar{x}$  relative to  $\mathcal{M}_{\bar{x}} := \{x \in \mathbb{R}^n : I(x) = I(\bar{x})\}$ , [24,  
582 Example 3.4]. From [24, Corollary 4.7] we see that  $f$  is partly smooth at  $\bar{x}$  relative to  $\mathcal{M}_{\bar{x}}$   
583 and that  $h$  is partly smooth at  $\bar{x}$  relative to  $\mathcal{M}$ . By [12, Corollary 4.2], the active manifold in  
584 the definition of partial smoothness is unique. Hence, near  $\bar{x}$  we have  $\mathcal{M} \equiv \mathcal{M}_{\bar{x}}$ .

585 Next, we show that whenever a vector  $\bar{v} \in \mathbb{R}^{n-m}$  satisfies  $\bar{x} + Uu + V\bar{v} \in \mathcal{M} \cap X$ , it  
586 must hold that  $\bar{v} = 0$ ; so  $\mathcal{M} \cap X = \{\bar{x} + Uu : u \in Y\}$ , for  $Y$  a neighbourhood of  $0 \in \mathbb{R}^m$ . To  
587 show the claim, consider  $i \in I(\bar{x})$  and note that, because  $\bar{x} + Uu + V\bar{v} \in \mathcal{M}_{\bar{x}}$ , it must be  
588 that  $I(\bar{x} + Uu + V\bar{v}) = I(\bar{x})$ , that is  $h(\bar{x} + Uu + V\bar{v}) = \langle a^i, \bar{x} + Uu + V\bar{v} \rangle = \langle a^i, \bar{x} \rangle + b^i +$   
589  $\langle a^i, Uu + V\bar{v} \rangle = h(\bar{x}) + \langle a^i, Uu + V\bar{v} \rangle$ . Therefore,

$$590 \quad (4.7) \quad h(\bar{x} + Uu + V\bar{v}) - h(\bar{x}) = \langle a^i, Uu + V\bar{v} \rangle .$$

591 Because  $0 \in \partial f(\bar{x})$ , we have that  $-\nabla q(\bar{x}) \in \partial h(\bar{x})$  and  $a^i + \nabla q(\bar{x}) \in \mathcal{V}$ . As a result,  
592  $\langle a^i + \nabla q(\bar{x}), Uu \rangle = 0$  and  $\langle a^i, Uu \rangle = -\langle \nabla q(\bar{x}), Uu \rangle$ . And (4.7) yields  $\langle V\bar{v}, a^i \rangle = h(\bar{x} + Uu +$   
593  $V\bar{v}) - h(\bar{x}) + \langle \nabla q(\bar{x}), Uu \rangle$ . The expression  $-\nabla q(\bar{x}) = \sum_{i \in I(\bar{x})} \bar{\alpha}_i a^i$  with  $\sum_{i \in I(\bar{x})} \alpha_i = 1$ ,  $\bar{\alpha}_i \geq$   
594  $0$  ( $i \in I(\bar{x})$ ) implies that  $\langle V\bar{v}, -\nabla q(\bar{x}) \rangle = h(\bar{x} + Uu + V\bar{v}) - h(\bar{x}) + \langle \nabla q(\bar{x}), Uu \rangle$ , and hence,  
595  $\langle V\bar{v}, a^i + \nabla q(\bar{x}) \rangle = 0$ . Because  $\mathcal{V} = \text{span}(\partial h(\bar{x}) + \nabla q(\bar{x})) = \text{span}\{a^i + \nabla q(\bar{x}) : i \in I(\bar{x})\}$ ,  
596 we actually have that  $\langle V\bar{v}, z \rangle = 0$  for all  $z \in \mathcal{V}$  and our claim that  $\bar{v} = 0$  follows.

597 Consider  $u \in Y$ . Since  $\mathcal{M} \cap X = \{\bar{x} + Uu : u \in Y\} = \mathcal{M}_{\bar{x}} \cap X$  and  $I(\bar{x} + Uu) = I(\bar{x})$ ,  
598 from the characterization of  $\partial h(x)$  in (4.6), it holds that  $\partial h(\bar{x} + Uu) = \partial h(\bar{x})$ . Consequently,  
599  $U^\top \partial h(\bar{x} + Uu) = U^\top \partial h(\bar{x}) = \bar{g}_u$ . On the other hand, in view of (4.3) in Thm. 4.4 we  
600 can take  $X$  and  $Y$  sufficiently small such that  $v_{\partial f}(u) \equiv 0$ . Consequently, the restriction  
601 of  $f$  on  $\mathcal{M} \cap X$  is  $f(\bar{x} + Uu) = q(\bar{x} + Uu) + h(\bar{x} + Uu)$  and, therefore,  $\nabla_u f(\bar{x} + Uu) =$   
602  $U^\top \nabla q(\bar{x} + Uu) + U^\top \partial h(\bar{x} + Uu)$ . Because  $U^\top \partial h(\bar{x} + Uu) = \bar{g}_u$ , this completes the proof, as  
603 then  $\nabla_{uu}^2 f(\bar{x} + Uu) = U^\top \nabla^2 q(\bar{x} + Uu)U$ .  $\square$

604 Partly smooth functions with the structure in (3.1) are considered in [3] to show that the  
605 proximal gradient method can identify the smooth manifold at a minimizer. This manifold is  
606 actually the fast track, which has been shown to be equivalent objects, for convex functions in  
607 [12] and for prox-regular functions in [26]. In the method proposed in [3], after identifying  
608 the manifold via the proximal gradient mapping, certain Riemannian gradient and Hessian are  
609 employed to compute a  $\mathcal{U}$ -Newton direction.

610 **4.3. Superlinear convergence of the PGVU method.** For superlinear convergence,  
611 naturally, properties of the partial  $\mathcal{U}$ -Hessian at  $\bar{x}$  associated with  $\bar{g}_v = 0$  are important. This  
612 matrix is assumed to be positive definite. Also, the Dennis-Moré-type condition below, typical  
613 in quasi-Newton methods (see, e.g., [19]), is required:

$$614 \quad (4.8) \quad (U_k Q_k U_k^\top - U \bar{W} U^\top) g^k = o(U g_u^k), \text{ where } \bar{W}^{-1} = \bar{H} := H(\bar{x}; 0) .$$

615 Recall that the matrix  $U$  spans the  $\mathcal{U}(\bar{x})$ -subspace.

616 **THEOREM 4.6 (Superlinear rate).** *Suppose  $f$  has a positive definite partial  $\mathcal{U}$ -Hessian*  
617  *$\bar{H}$  at  $\bar{x}$  associated with  $\bar{g}_v = 0$  and that (4.8) holds. Under the assumptions in Theorem 3.3,*  
618 *let  $\bar{x}$  be an accumulation point of  $\{x^k\}$  such that  $\bar{g} = 0 \in \text{ri } \partial f(\bar{x})$ . Then  $\bar{x}$  is the unique*  
619 *minimizer of  $f$ , and both  $\{x^k\}$  and  $\{p^k\}$  converge to  $\bar{x}$ . Furthermore, if  $\mu_k > \beta$ , then*  
620  *$\|x^{k+1} - \bar{x}\| = o(x^k - \bar{x})$ , i.e., the iterates generated by Algorithm 3.1 converge superlinearly.*

621 *Proof.* Because  $0 \in \text{ri } \partial f(\bar{x})$  and  $f$  has a positive definite partial  $\mathcal{U}$ -Hessian at  $\bar{x}$ , we  
622 have from Lemma 4.2 that  $\bar{x}$  is the unique minimizer of  $f$ . Recall from Theorem 3.3 that  
623 every accumulation point of the sequence  $\{x^k\}$  generated by Algorithm 3.1 is a minimizer  
624 of  $f$ . Consequently,  $\bar{x}$  is the unique accumulation point of  $\{x^k\}$ , with both  $\{x^k\}$  and  $\{p^k\}$   
625 converging to  $\bar{x}$ , as stated.

626 Next, using (3.6), adding  $\pm U\bar{W}U^\top g^k$ , and recalling the definition of  $(u^k, v^k)$  in Theo-  
627 rem 3.1:

$$\begin{aligned}
x^{k+1} - \bar{x} &= p^k - U_k Q_k U_k^\top g^k - \bar{x} \\
&= -U_k Q_k U_k^\top g^k + p^k - \bar{x} \pm U\bar{W}U^\top g^k \\
628 \quad &= \left( U\bar{W}U^\top - U_k Q_k U_k^\top \right) g^k - U\bar{W}U^\top g^k + (p^k - \bar{x})_u + (p^k - \bar{x})_v \\
&= \left( U\bar{W}U^\top - U_k Q_k U_k^\top \right) g^k + \left( Uu^k - U\bar{W}U^\top g^k \right) + v^k.
\end{aligned}$$

629 The third term above is of the order  $o(x^k - \bar{x})$ , by item (ii) in Corollary 3.4. The same holds  
630 for the first term, as  $\left( U\bar{W}U^\top - U_k Q_k U_k^\top \right) g^k = o(g^k)$  because of (4.8), and  $g^k = O(x^k - \bar{x})$ , by  
631 Theorem 3.1(ii). To conclude the proof, it remains to show that  $T_2 := \left( Uu^k - U\bar{W}U^\top g^k \right) =$   
632  $o(x^k - \bar{x})$ . Since  $T_2 = U \left( u^k - \bar{W}U^\top g^k \right) = U \left( u^k - \bar{W}g_u^k \right)$ , after multiplying on the left by  
633  $\bar{H}U^\top$ , we see that  $\bar{H}U^\top T_2 = \bar{H}u^k - g_u^k$ . Lemma 4.2(i) then ensures that  $\bar{H}U^\top T_2 = o(u^k)$ . The  
634 result follows, because  $u^k = O(x^k - \bar{x})$ , by Theorem 3.1(ii).  $\square$

635 The following two examples of polyhedral functions that are partly smooth relative to  
636 an affine or linear set from [43, Sec. 3.1], are common illustrations in regularized regression  
637 problems. The corresponding  $\mathcal{U}$ -subspaces have then an explicit expression. If

$$638 \quad (4.9) \quad h(x) = \|x\|_1 \text{ then } \mathcal{U}(p) = \text{lin} \{ e^j : j \in J(p) \} \text{ for } J(p) = \{ i \leq n : |p_i| > 0 \}.$$

639 Likewise, if  $h(x) = \|x\|_\infty$ , then  $\mathcal{U}(p) = \{ x : x_J = k \text{ sign}(p_J), k \in \mathbb{R} \}$  for the activity index  
640 set  $J = J(p) = \{ i \leq n : |p_i| = \|p\|_\infty \}$ .

641 For our numerical validation, we now consider problems having  $h$  as in (4.9).

642 **5. Application to  $\ell_1$ -regularized minimization.** We now apply the PGVU method  
643 (Algorithm 3.1), as its illustration, to solve  $\ell_1$ -regularized problems. So,

$$644 \quad (5.1) \quad \text{in (3.1) the nonsmooth function is } h(x) = \lambda \|x\|_1, \text{ for a positive parameter } \lambda,$$

645 and the proximal points for  $h$  are easy to compute. Accordingly, the  $\mathcal{V}$ -step in Procedure 1  
646 computes

$$647 \quad (5.2) \quad p^k = \max \left\{ 0, w^k - \frac{\lambda}{\mu_k} \right\} - \max \left\{ 0, -w^k - \frac{\lambda}{\mu_k} \right\}, \text{ for } w^k = x^k - \frac{1}{\mu_k} \nabla q(x^k).$$

648 The remaining two calculations that need to be specified in Algorithm 3.1 refer to the  
649 orthonormal basis  $U_k$ , and the positive semidefinite matrix  $Q_k$ .

650 **5.1. Defining  $Q_k$  and  $U_k$ .** For PGVU global convergence in Theorem 3.3, the matrices  
651 only need to be bounded. The superlinear rate in Theorem 4.6 requires a quite standard  
652 Dennis-Moré-type condition, natural in quasi-Newton frameworks. There are various choices  
653 that are compatible with the theory. For better numerical performance, it is preferable that  
654 matrices do not change too abruptly along consecutive iterations.

655 Regarding the second-order information along the  $\mathcal{U}$ -subspace, [43, Example 10] shows  
656 that the  $\ell_1$ -norm is partly smooth at any  $p \in \mathbb{R}^n$  relative to  $\mathcal{U}(p)$  defined in (4.9). As the  
657 same holds for  $f$  when  $h$  is as in (5.1), by Theorem 4.4, if  $0 \in V^\dagger \text{ri } \partial f(\bar{x})$  then

$$658 \quad \bar{H} := U^\top \nabla^2 q(\bar{x}) U$$

659 is a partial  $\mathcal{U}$ -Hessian of  $f$  at  $\bar{x}$  associated with  $\bar{g}_v = 0$ . Thus, a natural choice for the  
660 quasi-Newton matrices in the  $\mathcal{U}$ -step is to take  $Q_k^{-1} = U_k^\top \nabla^2 q(x^k) U_k$ .

661 The choice of the matrices  $U_k$  is more delicate. Since  $\partial f(x) = \nabla q(x) + \lambda \partial h(x)$ , the  
 662  $\mathcal{V}\mathcal{U}$ -decomposition is determined by the  $\mathcal{V}$ -subspace associated with  $h$ . Then, at first glance,  
 663 defining the basis matrix  $U_k$  for the subspaces  $\mathcal{U}(p^k)$  from (4.9) might appear straightforward.  
 664 Nevertheless, given  $\varepsilon \geq 0$ , we consider instead the smaller subspaces

665 (5.3) 
$$\mathcal{U}_\varepsilon(x) = \text{span} \{e^j : j \in J_\varepsilon(x)\} \text{ for } J_\varepsilon(x) := \left\{i \leq n : |x_i| > \frac{\varepsilon}{2}\right\} .$$

666 This  $\mathcal{V}_\varepsilon \mathcal{U}_\varepsilon$ -decomposition was introduced in [27] to deal with the lack of continuity of the  
 667 subdifferential as a multifunction. Unlike the  $\mathcal{V}\mathcal{U}$ -decomposition, the following important  
 668 continuity property, [27, equation (5.13)], holds for the  $\varepsilon$ -counterpart:

669 
$$\lim_{\substack{(x,\varepsilon) \rightarrow (\bar{x},0) \\ \varepsilon \geq 0}} \mathcal{V}_\varepsilon(x) = \mathcal{V}(\bar{x}) \quad \text{and} \quad \lim_{\substack{(x,\varepsilon) \rightarrow (\bar{x},0) \\ \varepsilon \geq 0}} \mathcal{U}_\varepsilon(x) = \mathcal{U}(\bar{x}) .$$

670 Thanks to this property, taking  $U_k$  as an orthonormal basis for  $[e^j : j \in J_{\varepsilon^k}(x^k)]$  with  $\varepsilon^k \rightarrow 0$   
 671 and  $x^k \rightarrow \bar{x}$ , ensures that  $U_k \rightarrow U$ , as needed to satisfy the Dennis-Moré condition (4.8) in  
 672 Theorem 4.6.

673 The choice of the parameter  $\varepsilon$  should ensure that it is driven to zero by the algorithmic  
 674 process. We discuss the impact of such choices on our previous example function.

675 **EXAMPLE 5.1 (Choosing  $\varepsilon$ ).** For  $F$  from Example 2.2, we run PGVU considering for the  
 676  $\mathcal{U}$ -step three different (natural) choices for  $\varepsilon$  in (5.3):

677 
$$\varepsilon_0^k = 0, \quad \varepsilon_1^k = f(x^k) - f(p^k) - \mu_k \|p^k - x^k\|^2, \quad \text{and} \quad \varepsilon_2^k = \mu_k \|p^k - x^k\|^2 .$$

678 The first option corresponds to the PGVU runs in Example 3.2. The second one transports  
 679  $\mu_k(x^k - p^k)$ , a subgradient of the model at  $p^k$ , to  $\mu_k(x^k - p^k) \in \partial_{\varepsilon_k} f(x^k)$ , [18, Prop. XI.4.2.2].  
 680 Since computing  $\varepsilon_1^k$  can be time consuming (it requires to evaluate  $f$  at two points), the third  
 681 option appears a good alternative (Theorem 3.3 shows that  $x^k - p^k \rightarrow 0$ ).

682 For the two initial values of  $\mu_0$  and  $x^0$  in Table 3.1, we run each PGVU variant, respectively  
 683 labeled PGVU-0,1,2, in a reference to the value of  $\varepsilon^k$  employed in (5.3). The comparison of  
 684 the number of iterations and digits of accuracy indicate  $\varepsilon_2^k$  as the best option, as illustrated by  
 the trajectories in Figure 5.1, generated with  $\varepsilon_0^k$  on the left plot and with  $\varepsilon_2^k$  on the right.

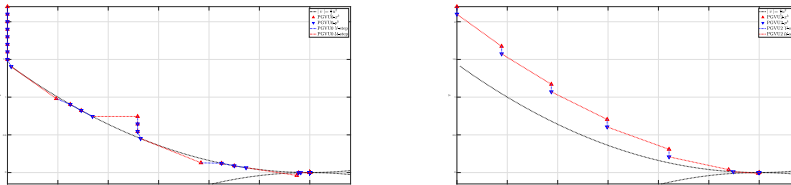


Fig. 5.1: Trajectories of PGVU-0 and PGVU-2 iterations, when minimizing  $F$  from  $x^0 = (-1.2, 2.2)^T$ , starting with  $\mu_0 = 10$ . Both variants stopped having reached more than 20 digits of accuracy, but PGVU-2 needed much less iterations. When the  $\mathcal{U}_{\varepsilon^k}(p^k)$  subspace is determined with  $\varepsilon^k = 0$ , as on the left plot, only 4  $\mathcal{U}$ -steps are done, and PGVU-0 needed 45 iterations to trigger the stopping test. For the trajectory on the right, by contrast, that was generated with  $\varepsilon_2^k$ , it sufficed to perform 9 iterations that involved 8  $\mathcal{U}$ -steps.

685

---

**Algorithm 5.1** Proximal Gradient  $\mathcal{V}\mathcal{U}$ -method for  $\ell_1$ -regularized minimization (PGVU-2)

---

**Data:**  $f$  from (5.1), starting point  $x^0$ , prox-parameter  $\mu_0$ , and a stopping tolerance  $\text{TOL}$ . Set  $k = 0$ .

**repeat**

Apply Procedure 1 with  $p^k$  defined in (5.2).  
 Shift the gradient  $g^k = \mu_k(x^k - p^k) + \nabla q(p^k) - \nabla q(x^k)$ .  
 Compute  $U_k = \left[ e^j : |x_j^k| > \frac{\varepsilon_k}{2}, 1 \leq j \leq n \right]$  for  $\varepsilon_k := \max(\text{TOL}, \mu_k \|p^k - x^k\|)$ .  
 Obtain the direction  $d^k = -W^k U_k^\top g^k$  for  $W^k \approx (U_k^\top \nabla^2 q(x^k) U_k)^\dagger$ .  
 Update  $x^{k+1} = p^k + U_k d^k$ , set  $k = k + 1$

**until**  $\|g^k\| \leq \text{TOL}$ ;

---

686 **5.2. Algorithm statement and numerical experiments.** We are now ready to introduce  
 687 the algorithm.

688 Being a special instance of PGVU, Algorithm 5.1 has global and superlinear convergence  
 689 if the conditions in Theorems 3.3 and 4.6 are satisfied. With our definitions, such is the case  
 690 if  $\mu_k > \beta$  and, for  $\bar{x}$  an accumulation point of  $\{x^k\}$  generated by Algorithm 5.1, it holds that  
 691  $0 \in \text{ri } \partial f(\bar{x})$  and  $\bar{H} = U^\top \nabla^2 q(\bar{x}) U$  is positive definite.

692 **5.2.1. Test functions and parameters.** The performance of Algorithm 5.1 is assessed  
 693 on regularized least-square problems:

694 in (1.1),  $f(x) = \frac{1}{2} \|Ax - b\|^2 + \lambda \|x\|_1$  for  $A \in \mathbb{R}^{m \times n}$ ,  $b \in \mathbb{R}^m$ , and  $\lambda = 0.1 \|A^\top b\|_\infty$ .

695 There are two sets of problems: 5000 mid-size randomly generated instances, and 95 large-  
 696 size statistical classification and regression instances in the webpage<sup>1</sup> of LIBSVM, a library  
 697 for support vector machines. The mid-size problems, referred below as QUAD, have random  
 698 dimensions  $m \in [10, 1000]$  and  $n \in [0.1m, 2m]$ . For half of the QUAD problems, that is 2500  
 699 runs, we set  $A = -\frac{1}{\sqrt{2n}} A'$  for a random matrix  $A'$ . The outcome vector  $b = Ab' + 0.0001\xi$  for  
 700 random  $\xi \in [0, 1]$  uniformly distributed. The vector  $b'$  has non-null components set to  $\pm 1$ ,  
 701 depending on a sparsity parameter randomly chosen. The second half of the QUAD problems,  
 702 sets  $A$  as a random matrix with normalized columns and  $b = Ab' + \sqrt{0.002}\xi$  for  $b'$  and  $\xi$  as  
 703 above. The support vector machine (SVM) problems are all scaled to  $[-1, 1]$  or  $[0, 1]$ .

704 In Procedure 1, we set  $x^0 = 0$ ,  $\sigma = 10^{-4}$  and  $\mu_0 = \frac{\nabla f(x^0)^\top \nabla f(x^0)}{2 \max\{1, |\nabla f(x^0)|\}}$ . The maximal  
 705 number of iterations was set to 100, and the stopping tolerance is  $\text{TOL} = 10^{-6}$ . After the  
 706  $\mathcal{U}$ -step, the prox-parameter is updated according to [23], i.e.,  $\mu_{k+1} = \frac{y^{k\top} y^k \mu_k}{y^{k\top} y^k + \mu_k y^{k\top} s^k}$  where  
 707  $y^k := \nabla f(x^{k+1}) - \nabla f(x^k)$  and  $s^k := x^{k+1} - x^k$ . In the  $\mathcal{U}$ -step, to compute  $d^k$ , since  
 708  $\nabla^2 q(x) = A^\top A$  for all  $x$ , we let  $\text{Id}$  denote the identity matrix of order  $n$  and define

709 
$$W^k = \begin{cases} (A^\top A + \text{TOL Id})^{-1} & \text{for QUAD,} \\ \text{diag}(A^\top A + \text{TOL Id})^{-1} & \text{for SVM.} \end{cases}$$

710 **5.2.2. Solvers and figures with evaluation measures.** The benchmark compares MAT-  
 711 LAB implementations of the following solvers:

- 712 1. PGVU-2, as in Algorithm 5.1.
- 713 2. SpaRSA 2.0, the sparse reconstruction by separable approximation[44]<sup>2</sup>.

<sup>1</sup><https://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/>

<sup>2</sup><http://www.lx.it.pt/~mtf/SpaRSA/>

- 714 3. FISTA, the fast iterative shrinkage-thresholding algorithm [4] <sup>3</sup>.
- 715 4. ADMM, the alternating direction method of multipliers[6], with parameters  $(\rho, \alpha) =$
- 716  $(0.0001, 1.5)$  <sup>4</sup>.
- 717 5. qNVU, the  $\mathcal{V}\mathcal{U}$ -algorithm from [35].

718 The experiments were performed on an Intel Core i7 computer with 12 cores and 32 GB  
719 RAM, running under Ubuntu 22.02.2 LTS.

720 The performance is measured by comparing the accuracy and the computing time of each  
721 solver, separately for the QUAD and SVM problems. We proceed as follows. First, for the  
722 accuracy criterion, if `best` is the smallest functional value found on a given instance for all  
723 the solvers, the accuracy of solver  $s$  is

$$724 \text{acc}(s) = \min \left( -\log \left( \frac{|f^*(s) - \text{best}|}{|\text{best}|} \right), 16 \right).$$

725 Solvers having achieved at least `cutoff`  $\in \{2, 4, 6\}$  digits of accuracy are considered success-  
726 ful. The corresponding values are reported in Table 5.1, where the high achieved accuracy, the  
727 differential of both  $\mathcal{V}\mathcal{U}$ -based methods, becomes evident. The accuracy achieved by FISTA  
728 is also impressive, being slightly inferior to  $\mathcal{V}\mathcal{U}$  for the large SVM instances.

Table 5.1: Successful runs for each solver

		PGVU-2	SpaRSA	FISTA	ADMM	qNVU
QUAD	<code>ACC</code> $\geq 2$	5000	4998	5000	516	5000
mid-size	<code>ACC</code> $\geq 4$	4999	4933	5000	5	5000
(5000 runs)	<code>ACC</code> $\geq 6$	4998	4493	5000	1	5000
SVM	<code>ACC</code> $\geq 2$	86	91	85	8	86
large-size	<code>ACC</code> $\geq 4$	81	51	79	2	86
(95 runs)	<code>ACC</code> $\geq 6$	74	26	70	0	86

729 SpaRSA has good accuracy for the mid-size instances, but performs less well for the  
730 SVM problems. On these runs, and with the considered parameters, ADMM did not perform  
731 well.

732 Regarding computing times, for low accuracy (`cutoff`=2), SpaRSA is always the fastest  
733 solver. The profile in Figure 5.2 compares computing times among the successful runs, for  
734 the value of `cutoff`=4. In the right plot in Figure 5.2, qNVU exhibits a slower performance.  
735 This is because at each iteration the qNVU method [35] solves two quadratic programming  
736 problems, a computationally expensive calculation for the large SVM instances. For the  
737 mid-size instances, SpaRSA remains the fastest solver, but not for the SVM problems. A  
738 solver-to-solver comparison clarifies this situation with four plots reported in Figure 5.3,  
739 comparing PGVU-2 to SpaRSA and FISTA, when SVM problems were solved with at least 2  
740 or 4 digits.

741 On the top left plot in Figure 5.3, we notice that to reach 2 digits of accuracy, SpaRSA  
742 is faster than PGVU-2, which is in turn faster than FISTA (right top plot). To get 4 digits,  
743 the bottom plots show that PGVU-2 always wins, reaching the accuracy level faster than both  
744 SpaRSA and FISTA.

745 **6. Concluding remarks.** We have extended the  $\mathcal{V}\mathcal{U}$ -theory for convex functions by  
746 defining the two variable  $\mathcal{U}$ -Lagrangian and the partial  $\mathcal{U}$ -Hessian. We showed that  $\mathcal{V}$ -

<sup>3</sup><https://github.com/tiepvupsu/FISTA>

<sup>4</sup><https://web.stanford.edu/~boyd/papers/admm/lasso/lasso.html>

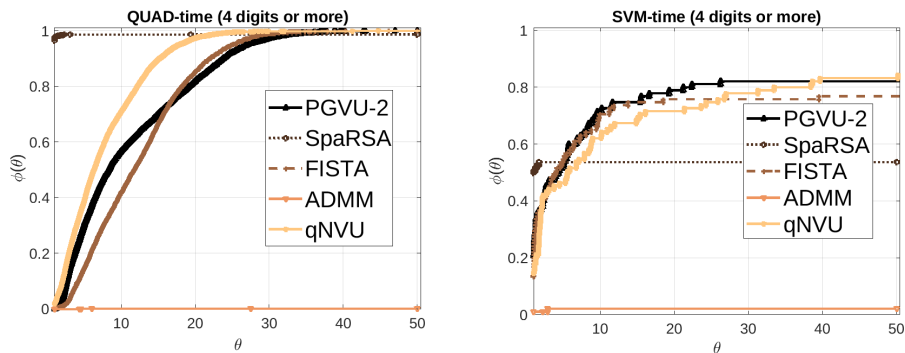


Fig. 5.2: Profiles for computing time for the successful runs for all solvers.

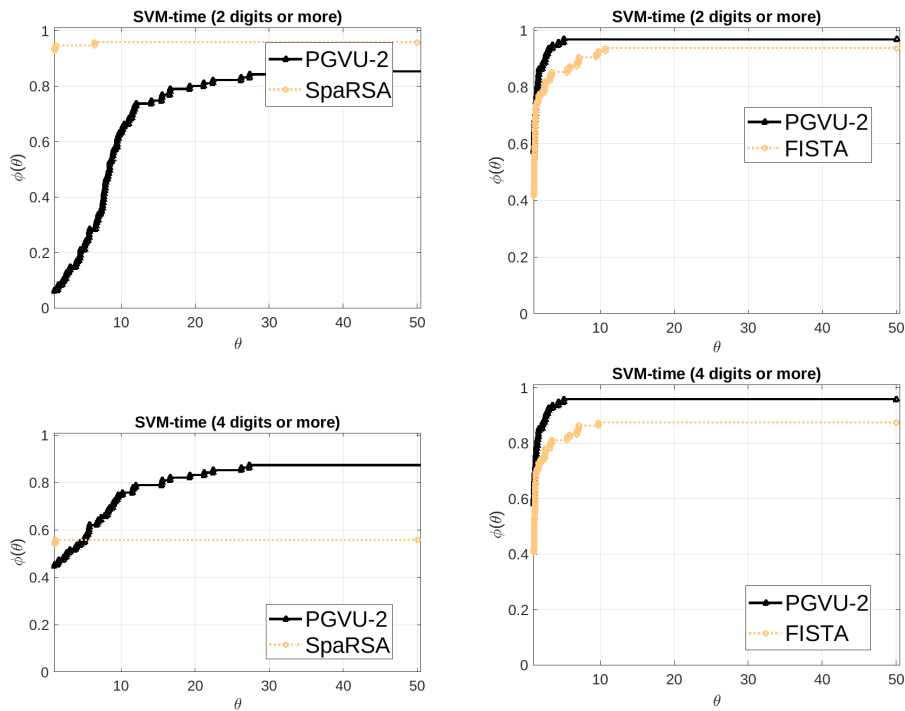


Fig. 5.3: Solver-to-solver time comparisons over SVM instances for different accuracy.

747 minimizers are *tangent* to the  $\mathcal{U}$ -subspace, an important property leading to superlinear  
 748 convergence of the Proximal Gradient  $\mathcal{V}\mathcal{U}$ -method, under natural assumptions.

749 For PDG-structured functions (including  $\ell_1$ -regularization), the Hessian of the single-  
 750 variable  $\mathcal{U}$ -Lagrangian exists along a certain fast-track [33, Theorem 4.1]. We extend this  
 751 result to our bivariate  $\mathcal{U}$ -Lagrangian, so that a Newtonian step can be performed as the  $\mathcal{U}$ -step.  
 752 In particular, we proved that partly smooth functions satisfying  $0 \in \text{ri } \partial f(\bar{x})$  always have a  
 753 partial  $\mathcal{U}$ -Hessian at  $\bar{x}$ .

754 We introduced the Proximal Gradient  $\mathcal{V}\mathcal{U}$  method, applicable to various structured

755 convex optimization problems, with superlinear convergence despite the nonsmoothness.  
 756 Numerical experiments verify that the method is particularly useful when high accuracy is  
 757 desired.

758 Originally defined for convex functions, the  $\mathcal{V}\mathcal{U}$ -theory has been generalized to the  
 759 nonconvex setting [33, 16, 26]. In [26] a localized version of  $\mathcal{U}$ -Lagrangian and the notion  
 760 of fast track are defined for a type of nonconvex functions called prox-regular functions [41],  
 761 and the correspondence between an active manifold of a partly smooth function and a fast  
 762 track is given. In [10], under the condition called tilt stability, the smoothness properties of  
 763 the function  $f$  restricted to the fast track are shown. Combining those theoretical results with  
 764 a suitable line-search, developing nonconvex versions of PGVU might be a subject for future  
 765 research.

766

## REFERENCES

- 767 [1] F. ATENAS, C. SAGASTIZÁBAL, P. J. S. SILVA, AND M. SOLODOV, *A unified analysis of descent sequences in weakly*  
 768 *convex optimization, including convergence rates for bundle methods*, SIAM Journal on Optimization,  
 769 33 (2023), pp. 89–115, <https://doi.org/10.1137/21M1465445>.
- 770 [2] G. BAREILLES, F. IUTZELER, AND J. MALICK, *Harnessing structure in composite nonsmooth minimization*,  
 771 SIAM Journal on Optimization, 33 (2023), pp. 2222–2247, <https://doi.org/10.1137/22M1505827>, <https://doi.org/10.1137/22M1505827>, <https://arxiv.org/abs/https://doi.org/10.1137/22M1505827>.
- 772 [3] G. BAREILLES, F. IUTZELER, AND J. MALICK, *Newton acceleration on manifolds identified by proximal*  
 773 *gradient methods*, Mathematical Programming, 200 (2023), pp. 37–70, <https://doi.org/10.1007/s10107-022-01873-w>, <https://doi.org/10.1007/s10107-022-01873-w>.
- 774 [4] A. BECK AND M. TEOULLE, *Fast gradient-based algorithms for constrained total variation image denoising*  
 775 *and deblurring problems*, IEEE Transactions on Image Processing, 18 (2009), pp. 2419–2434.
- 776 [5] A. BECK AND M. TEOULLE, *A fast iterative shrinkage-thresholding algorithm for linear inverse problems*,  
 777 SIAM Journal on Imaging Sciences, 2 (2009), pp. 183–202, <https://doi.org/10.1137/080716542>, <https://doi.org/10.1137/080716542>.
- 778 [6] S. BOYD, N. PARIKH, E. CHU, B. PELEATO, AND J. ECKSTEIN, *Distributed optimization and statistical learning*  
 779 *via the alternating direction method of multipliers*, Found. Trends Mach. Learn., 3 (2011), p. 1–122,  
 780 <https://doi.org/10.1561/2200000016>, <https://doi.org/10.1561/2200000016>.
- 781 [7] A. DANILIDIS, C. SAGASTIZÁBAL, AND M. SOLODOV, *Identifying structure of nonsmooth convex functions by*  
 782 *the bundle technique*, SIAM Journal on Optimization, 20 (2009), pp. 820–840, <https://doi.org/10.1137/080729864>.
- 783 [8] D. DRUSVYATSKIY AND A. S. LEWIS, *Tilt stability, uniform quadratic growth, and strong metric regularity*  
 784 *of the subdifferential*, SIAM Journal on Optimization, 23 (2013), pp. 256–267, <https://doi.org/10.1137/120876551>.
- 785 [9] D. DRUSVYATSKIY AND A. S. LEWIS, *Error bounds, quadratic growth, and linear convergence of proximal*  
 786 *methods*, Mathematics of Operations Research, 43 (2018), pp. 919–948, <https://doi.org/10.1287/moor.2017.0889>, <https://doi.org/10.1287/moor.2017.0889>.
- 787 [10] A. C. EBERHARD, Y. LUO, AND S. LIU, *On partial smoothness, tilt stability and the  $\mathcal{V}\mathcal{U}$ -decomposition*,  
 788 Mathematical Programming, 175 (2019), pp. 155–196, <https://doi.org/10.1007/s10107-018-1238-8>.
- 789 [11] W. HARE, *Numerical analysis of  $V\mathcal{U}$ -decomposition,  $U$ -gradient, and  $U$ -hessian approximations*, SIAM Journal  
 790 on Optimization, 24 (2014), pp. 1890–1913.
- 791 [12] W. HARE AND A. LEWIS, *Identifying active constraints via partial smoothness and prox-regularity*, J. Convex  
 792 Anal., 11 (2004), pp. 251–266.
- 793 [13] W. HARE, C. PLANIDEN, AND C. SAGASTIZÁBAL, *A derivative-free  $V\mathcal{U}$ -algorithm for convex finite-max prob-*  
 794 *lems*, Optimization Methods and Software, 0 (2019), pp. 1–39, <https://doi.org/10.1080/10556788.2019.1668944>, <https://doi.org/10.1080/10556788.2019.1668944>.
- 795 [14] W. HARE, C. PLANIDEN, AND C. SAGASTIZÁBAL, *The chain rule for  $V\mathcal{U}$ -decompositions of nonsmooth functions*,  
 796 Journal of Convex Analysis, 27 (2020), pp. 335–360.
- 797 [15] W. L. HARE, *Nonsmooth optimization with smooth substructure*, PhD thesis, Simon Fraser University, 2004.
- 798 [16] W. L. HARE AND R. POLIQUIN, *The quadratic sub-Lagrangian of a prox-regular function*, Nonlinear Analysis:  
 799 Theory, Methods & Applications, 47 (2001-08), pp. 1117–1128, [https://doi.org/10.1016/s0362-546x\(01\)00251-6](https://doi.org/10.1016/s0362-546x(01)00251-6).
- 800 [17] J. B. HIRIART-URRUTY, *The approximate first-order and second-order directional derivatives for a convex*  
 801 *function*, in Mathematical Theories of Optimization, J. P. Cecconi and T. Zolezzi, eds., Berlin, Heidelberg,  
 802 1983, Springer Berlin Heidelberg, pp. 144–177.
- 803 [18] J.-B. HIRIART-URRUTY AND C. LEMARÉCHAL, *Convex Analysis and Minimization Algorithms II*, vol. 306

- of Grundlehren der mathematischen Wissenschaften, Springer Berlin Heidelberg, New York, 1993, <https://doi.org/10.1007/978-3-662-06409-2>.
- [19] A. F. IZMAILOV AND M. V. SOLODOV, *Newton-type methods for optimization and variational problems*, Springer Series in Operations Research and Financial Engineering, Springer, Cham, 2014, <https://doi.org/10.1007/978-3-319-04247-3>.
- [20] C. LEMARÉCHAL AND C. SAGASTIZÁBAL, *Practical aspects of the Moreau-Yosida regularization: theoretical preliminaries*, SIAM Journal on Optimization, 7 (1997), pp. 367–385, <http://link.aip.org/link/?SJE/7/367/1>.
- [21] C. LEMARÉCHAL AND F. OUSTRY, *Growth conditions and U-Lagrangians*, Set-Valued Analysis, 9 (2001), pp. 123–129, <https://doi.org/10.1023/A:1011267019516>.
- [22] C. LEMARÉCHAL, F. OUSTRY, AND C. SAGASTIZÁBAL, *The U-Lagrangian of a convex function*, Transactions of the American Mathematical Society, 352 (2000), pp. 711–729.
- [23] C. LEMARÉCHAL AND C. SAGASTIZÁBAL, *Variable metric bundle methods: From conceptual to implementable forms*, Mathematical Programming, 76 (1997), pp. 393–410, <https://doi.org/10.1007/BF02614390>.
- [24] A. S. LEWIS, *Active sets, nonsmoothness, and sensitivity*, SIAM Journal on Optimization, 13 (2002), pp. 702–725, <https://doi.org/10.1137/s1052623401387623>.
- [25] J. LIANG, J. FADILI, AND G. PEYRÉ, *Activity identification and local linear convergence of forward-backward-type methods*, SIAM Journal on Optimization, 27 (2017), pp. 408–437, <https://doi.org/10.1137/16M106340X>, <https://arxiv.org/abs/https://doi.org/10.1137/16M106340X>.
- [26] S. LIU, A. EBERHARD, AND Y. LUO, *The U-Lagrangian, fast track, and partial smoothness of a prox-regular function*, Set-Valued and Variational Analysis, 28 (2020), pp. 369–394, <https://doi.org/10.1007/s11228-019-00518-z>.
- [27] S. LIU, C. SAGASTIZÁBAL, AND M. SOLODOV, *Subdifferential enlargements and continuity properties of the VU-decomposition in convex optimization*, in Nonsmooth Optimization and Its Applications, S. Hosseini, B. S. Mordukhovich, and A. Uschmajew, eds., Springer International Publishing, Cham, 2019, pp. 55–87, [https://doi.org/10.1007/978-3-030-11370-4\\_4](https://doi.org/10.1007/978-3-030-11370-4_4).
- [28] R. MIFFLIN AND C. SAGASTIZÁBAL, *Optimization Stories*, vol. Extra Volume ISMP 2012, ed. by M. Grötschel, DOCUMENTA MATHEMATICA, 2012, ch. A Science Fiction Story in Nonsmooth Optimization Originating at IIASA, p. 460.
- [29] R. MIFFLIN AND C. SAGASTIZÁBAL, *VU-decomposition derivatives for convex max-functions*, Ill-Posed Variational Problems and Regularization Techniques, 477 (1999), pp. 167–186.
- [30] R. MIFFLIN AND C. SAGASTIZÁBAL, *On VU-theory for functions with primal-dual gradient structure*, SIAM Journal on Optimization, 11 (2000), pp. 547–571.
- [31] R. MIFFLIN AND C. SAGASTIZÁBAL, *Proximal points are on the fast track*, Journal of Convex Analysis, 9 (2002), pp. 563–580.
- [32] R. MIFFLIN AND C. SAGASTIZÁBAL, *Primal-dual gradient structured functions: second-order results; links to epi-derivatives and partly smooth functions*, SIAM Journal on Optimization, 13 (2003), pp. 1174–1194.
- [33] R. MIFFLIN AND C. SAGASTIZÁBAL, *UV-smoothness and proximal point results for some nonconvex functions*, Optimization Methods and Software, 19 (2004), pp. 463–478, <https://doi.org/10.1080/10556780410001704902>.
- [34] R. MIFFLIN AND C. SAGASTIZÁBAL, *Relating U-Lagrangians to second order epiderivatives and proximal tracks*, Journal of Convex Analysis, 12 (2005), pp. 81–93.
- [35] R. MIFFLIN AND C. SAGASTIZÁBAL, *A VU-algorithm for convex minimization*, Mathematical Programming, 104 (2005), pp. 583–608.
- [36] S. A. MILLER AND J. MALICK, *Newton methods for nonsmooth convex minimization: connections among U-Lagrangian, Riemannian Newton and SQP methods*, Mathematical Programming, 104 (2005), pp. 609–633.
- [37] Y. NESTEROV, *Smooth Convex Optimization*, Springer US, 2004, pp. 51–110, [https://doi.org/10.1007/978-1-4419-8853-9\\_2](https://doi.org/10.1007/978-1-4419-8853-9_2).
- [38] J. ORTEGA AND W. RHEINOLDT, *Iterative Solution of Nonlinear Equations in Several Variables*, Academic Press, New York, 1970.
- [39] F. OUSTRY, *A second-order bundle method to minimize the maximum eigenvalue function*, Mathematical Programming, Series B, 89 (2000), pp. 1–33, <https://doi.org/10.1007/s101070000166>, <https://www.scopus.com/inward/record.uri?eid=2-s2.0-0002928322&partnerID=40&md5=4ff600667e238f59008a3320f8e82dae>.
- [40] C. PLANIDEN AND T. RAJAPAKSHA, *Linear convergence of the derivative-free proximal bundle method on convex nonsmooth functions, with application to the derivative-free VU-algorithm*, Set-Valued and Variational Analysis, 32 (2024), <https://doi.org/10.1007/s11228-024-00718-2>, <https://doi.org/10.1007/s11228-024-00718-2>.
- [41] R. T. ROCKAFELLAR AND R. J.-B. WETS, *Variational Analysis*, vol. 317 of Grundlehren der mathematischen Wissenschaften, Springer, Berlin, 1998.



- 874 [42] C. SAGASTIZÁBAL, *A  $\mathcal{V}\mathcal{U}$ -point of view of nonsmooth optimization*, in Proceedings of the International  
875 Congress of Mathematicians 2018- Invited Lectures, vol. 3, 2018, pp. 3785–3806.
- 876 [43] S. VAITER, C. DELEDALLE, J. FADILI, G. PEYRÉ, AND C. DOSSAL, *The degrees of freedom of partly smooth*  
877 *regularizers*, Annals of the Institute of Statistical Mathematics, 69 (2017), pp. 791–832, [https://doi.org/](https://doi.org/10.1007/s10463-016-0563-z)  
878 [10.1007/s10463-016-0563-z](https://doi.org/10.1007/s10463-016-0563-z).
- 879 [44] S. J. WRIGHT, R. D. NOWAK, AND M. A. T. FIGUEIREDO, *Sparse reconstruction by separable approximation*,  
880 IEEE Transactions on Signal Processing, 57 (2009), pp. 2479–2493, [https://doi.org/10.1109/TSP.2009.](https://doi.org/10.1109/TSP.2009.2016892)  
881 [2016892](https://doi.org/10.1109/TSP.2009.2016892).