



**HAL**  
open science

# Harnessing structure in composite nonsmooth minimization

Gilles Bareilles, Franck Iutzeler, Jérôme Malick

► **To cite this version:**

Gilles Bareilles, Franck Iutzeler, Jérôme Malick. Harnessing structure in composite nonsmooth minimization. 2022. hal-03706958v1

**HAL Id: hal-03706958**

**<https://hal.science/hal-03706958v1>**

Preprint submitted on 28 Jun 2022 (v1), last revised 20 Mar 2023 (v3)

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

1  
2

# HARNESSING STRUCTURE IN COMPOSITE NONSMOOTH MINIMIZATION\*

3 GILLES BAREILLES<sup>†</sup>, FRANCK IUTZELER<sup>†</sup>, AND JÉRÔME MALICK<sup>‡</sup>

4 **Abstract.** We consider the problem of minimizing the composition of a nonsmooth function  
5 with a smooth mapping in the case where the proximity operator of the nonsmooth function can  
6 be explicitly computed. We first show that this proximity operator can provide the exact smooth  
7 substructure of minimizers, not only of the nonsmooth function, but also of the full composite  
8 function. We then exploit this proximal identification by proposing an algorithm which combines  
9 proximal steps with sequential quadratic programming steps. We show that our method identifies  
10 the optimal smooth substructure and converges locally quadratically. We illustrate its behavior on  
11 two problems: the minimization of a maximum of quadratic functions and the minimization of the  
12 maximal eigenvalue of a parametrized matrix.

13 **Key words.** Nonsmooth optimization, proximal operator, partial smoothness, manifold identi-  
14 fication, maximum eigenvalue minimization, sequential quadratic programming.

15 **AMS subject classifications.** 65K10, 90C26, 49Q12, 90C55.

16 **1. Introduction.**

17 **1.1. Context: structured nonsmooth optimization.** In this paper, we con-  
18 sider nonsmooth optimization problems of the form

19 (1.1) 
$$\min_{x \in \mathbb{R}^n} F(x) \triangleq g(c(x)),$$
  
20

21 where the inner mapping  $c : \mathbb{R}^n \rightarrow \mathbb{R}^m$  is smooth and the outer function  $g : \mathbb{R}^m \rightarrow$   
22  $\mathbb{R} \cup \{+\infty\}$  is nonsmooth and may be nonconvex, but admits an explicit proximity  
23 operator. Such composite nonsmooth optimization problems appear in a variety of  
24 applications in signal processing, machine learning, and control, such as robust nonlin-  
25 ear regression, phase synchronization, nonsmooth penalty functions; see *e.g.* [20, 27]  
26 and the references therein.

27 Throughout the paper, we illustrate our developments on two classes of functions:  
28 the pointwise maximum of  $m$  smooth real-valued functions  $c_i$

29 (1.2) 
$$F(x) = \max_{i=1, \dots, m} (c_i(x))$$
  
30

31 and the maximum eigenvalue of a parametrized symmetric real matrix  $c$

32 (1.3) 
$$F(x) = \lambda_{\max}(c(x)).$$
  
33

34 In these two examples and many others, subgradients of  $F$  can be computed and thus  
35 the composite function can be minimized using standard nonsmooth optimization  
36 algorithms (*e.g.* subgradient methods, gradient sampling [7], nonsmooth BFGS [19],  
37 or bundle methods [13]). Nevertheless, these methods do not exploit the fact that  
38  $F$  is a composition of a smooth mapping  $c$ , which can hinder their performance.  
39 In contrast, the so-called prox-linear methods leverage this composite expression by  
40 introducing an extension of the proximity operator where the nonlinear mapping  $c$

---

\*Submitted to the editors June 27, 2022.

<sup>†</sup>Univ. Grenoble Alpes, LJK, France ([gilles.bareilles@univ-grenoble-alpes.fr](mailto:gilles.bareilles@univ-grenoble-alpes.fr), [franck.iutzeler@univ-grenoble-alpes.fr](mailto:franck.iutzeler@univ-grenoble-alpes.fr)).

<sup>‡</sup>Univ. Grenoble Alpes, CNRS, Grenoble INP, LJK ([jerome.malick@univ-grenoble-alpes.fr](mailto:jerome.malick@univ-grenoble-alpes.fr)).

41 is iteratively replaced by a first-order Taylor approximation [20]. These methods  
 42 benefit from theoretical convergence guarantees, and nicely generalize to Taylor-like  
 43 approximations [9, 3]. However these methods are not always directly implementable  
 44 because the prox-linear step may be hard to compute, as in (1.3).

45 In this paper, we propose an optimization algorithm for solving (1.1) exploiting  
 46 that the nonsmooth objective function  $F = g \circ c$  writes as a composition with a simple  
 47 nonsmooth function  $g$ , which displays some smooth substructure, as discussed below.

48 **1.2. Smooth substructure, identification, and existing algorithms.** For  
 49 many composite functions, including (1.2) and (1.3), the nondifferentiability points  
 50 locally organize into *smooth manifolds over which  $F$  evolves smoothly*. We illustrate  
 51 in Figure 1 such a smooth substructure for a maximum of two functions.

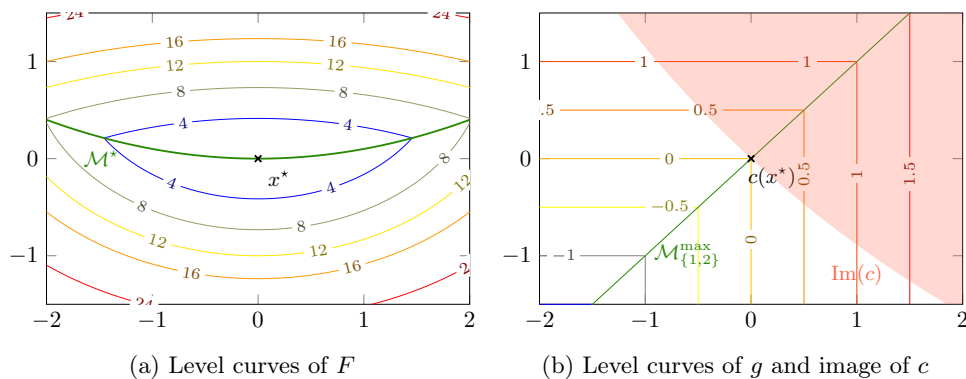


Fig. 1: Smooth substructure on a simple example ( $n = m = 2$ ). The figures show the level curves of  $g(y) = \max(y_1, y_2)$  (on the right, in the intermediate space) and of  $F = g \circ c$ , with two quadratic functions  $c_1(x)$ ,  $c_2(x)$  (on the left, in the input space). The manifolds of non-differentiability are in green; the image of  $c$  is the red area.

52 The smooth substructure of  $F$  can help in solving (1.1). Indeed, if the optimal  
 53 solution  $x^*$  belongs to a manifold  $\mathcal{M}^*$  that is known beforehand, then minimizing  
 54 the nonsmooth function  $F$  over  $\mathbb{R}^n$  boils down to minimizing the smooth restriction  
 55  $F|_{\mathcal{M}^*}$  over this smooth *optimal manifold*  $\mathcal{M}^*$ . This would enable to solve (1.1) by  
 56 smooth constrained optimization algorithms, such as Sequential Quadratic Program-  
 57 ming (SQP) methods (see *e.g.* [23, 5]). The main difficulty in practice is that *we do*  
 58 *not know  $\mathcal{M}^*$  in advance*.

59 Thus, the algorithms exploiting this smooth substructure require two ingredients:

- 60 i) a mechanism to identify the optimal manifold;
- 61 ii) an efficient method to minimize  $F$  restricted to this manifold.

62 For general convex functions, the algorithm of [21] mixes a proximal bundle it-  
 63 eration (as a heuristic for identification) and a so-called  $\mathcal{U}$ -Newton iteration (which  
 64 interprets as an SQP step; see [22, Sec. 5]). The obtained superlinear rate hinges on  
 65 the identification of the optimal manifold.

66 For max-of-smooth functions (1.2), the paper [29] pioneered the idea of seeking  
 67 the optimal manifold and using it to make second-order steps. Their identification  
 68 heuristic uses the indices of the maximal function along a descent direction. Recently,  
 69 [17, 10] investigate a related setting and propose bundle-like algorithms incorporat-  
 70 ing high-order information that converge (super)linearly on max-of-smooth functions  
 71 when the optimal manifold is known.

72 For the maximum eigenvalue of a parametrized matrix (1.3), a specific version of  
 73 the  $\mathcal{U}$ -Newton method discussed above is studied by [24]. Again, the identification  
 74 mechanism is a heuristic determining the multiplicity of the maximal eigenvalue and  
 75 the optimization step is an SQP iteration.

76 None of these methods guarantee identification of the optimal manifold: they  
 77 either assume that the optimal manifold is known in advance, or rely on heuristics for  
 78 identification. Here, we aim at further harnessing the smooth substructure of  $F = g \circ c$   
 79 to have *guaranteed local identification* of the optimal manifold and then *guaranteed*  
 80 *quadratic convergence* when using SQP iterations.

81 **1.3. Contributions and outline.** We propose a second-order algorithm for  
 82 solving the nonsmooth composite problem (1.1) that identifies the optimal manifold  
 83 of non-differentiability. The two main ingredients of our algorithm are the following:

- 84 i) we use the explicit proximal operator of  $g$  with chosen stepsizes to provide a  
 85 guaranteed identification procedure;
- 86 ii) for a candidate manifold  $\mathcal{M}$ , we make an SQP iteration minimizing a smooth  
 87 extension of  $F|_{\mathcal{M}}$  subject to the constraint of belonging to  $\mathcal{M}$ .

88 The fact proximal-based operators have identification properties around mini-  
 89 mizers is well-known: the proximal operator [11], the proximal gradient operator [1],  
 90 approximate variable-metric proximal gradient operators [14], and prox-linear oper-  
 91 ators [20] locally identify the optimal manifold under some natural geometrical as-  
 92 sumptions. Here, we only have access to the proximity operator of  $g$ , and in order to  
 93 exploit the structure it provides, we face the double challenge of, first, identifying the  
 94 smooth structure around a point which is not a minimizer for  $g$ , and, second, deducing  
 95 the corresponding structure of  $F = g \circ c$ . Thus, our main technical contribution is to  
 96 establish that  $\mathbf{prox}_{\gamma g}$  maps a point  $y$  close to  $c(x^*)$  to  $c(\mathcal{M}^*)$ . The step  $\gamma$  should be  
 97 carefully chosen, in particular larger than the distance of  $y$  to  $c(x^*)$ . Mathematically,  
 98 we study the range of steps for which the curve  $\gamma \mapsto \mathbf{prox}_{\gamma g}(y)$  belongs to  $c(\mathcal{M}^*)$ .  
 99 This analysis shows connections with recent works in nonsmooth analysis, such as the  
 100 modulus of identifiability appearing in [16].

101 We combine this new identification result with standard SQP-steps to propose an  
 102 algorithm for minimizing the composite function  $F$ . We pay a special attention to  
 103 prevent the quadratic convergence of SQP from jeopardizing identification: we prove  
 104 that, for a well-chosen stepsize policy, the method identifies the optimal structure  
 105 and locally converges quadratically. We illustrate numerically these properties on  
 106 problems of the form (1.2) and (1.3).

107 The outline of the remainder of the paper is as follows. First, in Section 2,  
 108 we introduce the technical tools to describe the manifold identification brought by  
 109 proximity operators (including prox-regularity and partial smoothness). Furthermore,  
 110 we lay out two technical properties needed for proximal identification in the composite  
 111 setting. In Section 3, we show our main result consisting in a description of a stepsize  
 112 range for which the proximity operator of  $g$  identifies the optimal manifold locally  
 113 around a minimizer. In Section 4 we detail the proposed method combining SQP-  
 114 steps and proximal identification steps. Finally, we present in Section 5 numerical  
 115 illustrations of our method and of the identification result.

116 **2. Setting and assumptions.** Let us start by representing schematically the  
 117 type of functions we consider:

$$118 \quad \mathbb{R}^n \xrightarrow[\text{smooth mapping}]{c} \text{Im}(c) \subset \mathbb{R}^m \xrightarrow[\text{nonsmooth function}]{g} \mathbb{R} \cup \{+\infty\}.$$

Throughout the paper, we denote by  $x$  points in the *input space*  $\mathbb{R}^n$  and by  $y$  points in the *intermediate space*  $\mathbb{R}^m$ .

In all the results presented in this paper, we make the following assumption that describes the minimal global properties on  $g$  and  $c$  to conduct our reasoning.

ASSUMPTION 2.1. *The mapping  $c : \mathbb{R}^n \rightarrow \mathbb{R}^m$  is  $\mathcal{C}^2$ , the function  $g : \mathbb{R}^m \rightarrow \mathbb{R} \cup \{+\infty\}$  is proper and lower semi-continuous.*

We work with the set of (*general*) *subgradients* (see [26, Def. 8.3]), defined at a point  $\bar{y}$  where  $g(\bar{y})$  is finite as:

$$\partial g(\bar{y}) \triangleq \left\{ \lim_r v_r : v_r \in \widehat{\partial} g(y_r), y_r \rightarrow \bar{y}, g(y_r) \rightarrow g(\bar{y}) \right\},$$

where  $\widehat{\partial} g(\bar{y})$  denotes the set of *regular (or Fréchet) subgradients*, defined as

$$\widehat{\partial} g(\bar{y}) \triangleq \{v : g(y) \geq g(\bar{y}) + \langle v, y - \bar{y} \rangle + o(\|y - \bar{y}\|) \text{ for all } y \in \mathbb{R}^m\}.$$

These two subdifferentials match if (and only if)  $g$  is (Clarke) regular at  $\bar{y}$ . Closed convex functions are regular everywhere and these subdifferentials match the usual convex subdifferential (see [26, Chap. 8.11-12] for details).

In the remainder of this section, we provide quick recalls and definitions about the two important objects of our analysis: the proximity operator in [Subsection 2.1](#) and the structure manifolds in [Subsection 2.2](#). We illustrate them on our running examples (1.2) and (1.3).

**2.1. Proximity operator.** The proximity operator of a function  $g$  with step  $\gamma > 0$  at  $y \in \mathbb{R}^m$  is defined as the set-valued mapping

$$\mathbf{prox}_{\gamma g}(y) \triangleq \operatorname{argmin}_{u \in \mathbb{R}^m} \left\{ g(u) + \frac{1}{2\gamma} \|u - y\|^2 \right\}.$$

This operator is well-defined when  $g$  is prox-regular and prox-bounded; see *e.g.* [26, 13.37]. We quickly introduce these two notions and recall a result on the uniqueness and characterization of the prox operator, which is important in our developments.

A function  $g$  is *prox-regular* at a point  $\bar{y}$  for a subgradient  $\bar{v} \in \partial g(\bar{y})$  if  $g$  is finite, locally lower semi-continuous at  $\bar{y}$ , and there exists  $r > 0$  and  $\epsilon > 0$  such that

$$g(y') \geq g(y) + \langle v, y' - y \rangle - \frac{r}{2} \|y' - y\|^2$$

whenever  $v \in \partial g(y)$ ,  $\|y - \bar{y}\| < \epsilon$ ,  $\|y' - \bar{y}\| < \epsilon$ ,  $\|v - \bar{v}\| < \epsilon$ , and  $g(y) < g(\bar{y}) + \epsilon$ . When this holds for all  $\bar{v} \in \partial g(\bar{y})$ , we say that  $g$  is prox-regular at  $\bar{y}$  [26, Def. 13.27].

A function  $g$  is *prox-bounded* if there exists  $R \geq 0$  such that the function  $g + \frac{R}{2} \|\cdot\|^2$  is bounded below. The corresponding *threshold* (of prox-boundedness) is the smallest  $r_{pb} \geq 0$  such that  $g + \frac{R}{2} \|\cdot\|^2$  is bounded below for all  $R > r_{pb}$ . In this case,  $g + \frac{R}{2} \|\cdot - \bar{y}\|^2$  is bounded below for any  $\bar{y}$  and  $R > r_{pb}$  [26, Def. 1.23, Th. 1.25].

We can now recall a relevant result on the characterization of proximal points.

PROPOSITION 2.2 ([12, Th. 1]). *Suppose that the function  $g$  is prox-regular at  $\bar{y}$  for  $\bar{v} \in \partial g(\bar{y})$  with parameter  $r_{pr}$ , and prox-bounded with threshold  $r_{pb}$ . Then, for any  $\gamma < \min(r_{pr}^{-1}, r_{pb}^{-1})$  and all  $y$  near  $\bar{y} + \gamma \bar{v}$ , the proximal operator is:*

- *single-valued and locally Lipschitz continuous;*
- *uniquely determined by the relation*

$$p = \mathbf{prox}_{\gamma g}(y) \Leftrightarrow y - p \in \gamma \partial g(p).$$

165 In addition to its existence and characterization provided by the result above, the  
 166 proximity operator has a closed-form expression in our running examples.

167 **EXAMPLE 2.3** (Maximum). *The subdifferential of  $g(y) = \max(y_1, \dots, y_m)$  is*

$$168 \quad \partial \max(y) = \text{Conv} \{e_i : y_i = \max(y)\},$$

170 where  $e_i$  is the  $i$ -th element of the Cartesian basis of  $\mathbb{R}^m$ . This function is convex, thus  
 171 globally prox-regular and prox-bounded everywhere (with parameters 0). Its proximity  
 172 operator is given (coordinate-wise) by

$$173 \quad [\mathbf{prox}_{\gamma \max}(y)]_i = \begin{cases} s & \text{if } y_i > s \\ y_i & \text{else} \end{cases}$$

175 where  $s$  is the unique real number such that  $\sum_{\{i: y_i > s\}} (y_i - s) = \gamma$ .

176 **EXAMPLE 2.4** (Maximum eigenvalue). *Denote the eigenvalue decomposition of a*  
 177 *point  $y \in \mathbb{S}_m$  as  $y = E \text{Diag}(\lambda) E^\top$ , where  $\lambda \in \mathbb{R}^m$  is a vector with decreasing entries*  
 178 *and  $E \in \mathbb{R}^{m \times m}$  an orthogonal matrix. The subdifferential of the maximum eigenvalue*  
 179 *at  $y$  writes [18, Ex. 3.6]*

$$180 \quad \partial \lambda_{\max}(y) = \{E_{1:r} Z E_{1:r}^\top, Z \in \mathbb{S}_r, Z \succeq 0, \text{trace } Z = 1\}$$

182 where  $r$  is the multiplicity of the maximum eigenvalue of  $y$ . This function is convex,  
 183 thus prox-regular and prox-bounded (with parameters 0). Its proximity operator can  
 184 be expressed using the one of the max function as

$$185 \quad \mathbf{prox}_{\gamma \lambda_{\max}}(y) = E \text{Diag}(\mathbf{prox}_{\gamma \max}(\lambda)) E^\top.$$

187 **2.2. Structure manifolds.** We now specify the notion of structure manifold in  
 188 relation with a nonsmooth function  $g$ .

189 A subset  $\mathcal{M}$  of  $\mathbb{R}^n$  is said to be a  $p$ -dimensional  $\mathcal{C}^2$ -submanifold of  $\mathbb{R}^n$  around  
 190  $\bar{x} \in \mathcal{M}$  if there exists a  $\mathcal{C}^2$  manifold-defining map  $h : \mathbb{R}^n \rightarrow \mathbb{R}^{n-p}$  with a surjective  
 191 derivative at  $\bar{x} \in \mathcal{M}$  that satisfies for all  $x$  close enough to  $\bar{x}$ :  $x \in \mathcal{M} \Leftrightarrow h(x) = 0$ .  
 192 We define the tangent and normal spaces at a point  $x \in \mathcal{M}$  as follows:

$$193 \quad T_x \mathcal{M} = \ker D h(x) \quad N_x \mathcal{M} = \text{Im } D h(x)^*$$

195 The important notion of *structure manifolds of  $g$*  can be defined as a manifold  
 196  $\mathcal{M}^g$  where  $g$  is nondifferentiable. More precisely, at a point  $\bar{y} \in \mathcal{M}^g$ , we require  $g$   
 197 to be prox-regular and *partly smooth*. This property of ( $\mathcal{C}^2$ -)partial smoothness is  
 198 verified at a point  $\bar{y}$  for a function  $g$  relatively to a set  $\mathcal{M}^g$  containing  $\bar{y}$  if  $\mathcal{M}^g$  is a  
 199  $\mathcal{C}^2$  manifold around  $\bar{y}$  and if

- 200 • (smoothness) the restriction of  $g$  to  $\mathcal{M}^g$  is a  $\mathcal{C}^2$  function near  $\bar{y}$ ;
- 201 • (regularity)  $g$  is (Clarke) regular at all points  $y \in \mathcal{M}^g$  near  $\bar{y}$ , with  $\partial g(y) \neq \emptyset$ ;
- 202 • (sharpness) the affine span of  $\partial g(\bar{y})$  is a translate of  $N_{\bar{y}} \mathcal{M}^g$ ;
- 203 • (sub-continuity) the mapping  $\partial g$  restricted to  $\mathcal{M}^g$  is continuous at  $\bar{y}$ .

204 The concept of partial smoothness, introduced in [18], captures (locally) well-  
 205 behaved nonsmoothness by requiring  $g$  to be smooth along a manifold and non-smooth  
 206 across it. In addition, the prox-regularity of  $g$  ensures unicity of the structure manifold  
 207 near  $\bar{y}$  [11, Corollary 4.2, Example 7.1]. To highlight the relation between the manifold  
 208 and the function  $g$ , we use the notation  $\mathcal{M}^g$  for the structure manifold related to  $g$ .

209 EXAMPLE 2.5. *The structure manifolds of max are*

$$210 \quad \mathcal{M}_I^{\max} = \{y \in \mathbb{R}^m : y_i = \max(y) \text{ for } i \in I\},$$

212 *where  $I \subset \{1, \dots, m\}$ . A smooth manifold-defining map for  $\mathcal{M}_I^{\max}$  is  $h : \mathbb{R}^m \rightarrow \mathbb{R}^{|I|-1}$*   
 213 *such that  $h(y)_l = y_{i_l} - y_{i_{|I|}}$ , where  $|I|$  denotes the size of  $I$  and  $i_l$  the  $l$ -th element of*  
 214  *$I$  (with some ordering). As required, this map is surjective. At any point  $y \in \mathbb{R}^m$ , the*  
 215 *maximum is partly smooth relative to  $\mathcal{M}_I^{\max}$ , where  $I = \{i : y_i = \max(y)\}$ .*

216 EXAMPLE 2.6. *The structure manifolds of  $\lambda_{\max}$  in  $\mathbb{S}_m$  consist of all matrices*  
 217 *having a largest eigenvalue with fixed multiplicity  $r$ :*

$$218 \quad \mathcal{M}_r^{\lambda_{\max}} = \{y \in \mathbb{S}_m : \lambda_1(y) = \dots = \lambda_r(y)\}.$$

220 *A manifold-defining map of  $\mathcal{M}_r^{\lambda_{\max}}$  is described in [28] and  $\lambda_{\max}$  is partly smooth*  
 221 *relative to  $\mathcal{M}_r^{\lambda_{\max}}$  at any point  $y \in \mathcal{M}_r^{\lambda_{\max}}$ .*

222 In view of the expression of the proximity operators in our examples, their output  
 223 naturally lie on the structure manifolds described above. More precisely,  $\mathbf{prox}_{\gamma \max}(y)$   
 224 belongs to the structure manifold  $\mathcal{M}_I^{\max}$ , where  $I$  collects the indices of the  $k$  largest  
 225 entries of  $y$  and  $k$  grows as  $\gamma$  increases. Similarly,  $\mathbf{prox}_{\gamma \lambda_{\max}}(y)$  belongs to the  
 226 structure manifold  $\mathcal{M}_r^{\lambda_{\max}}$ , where  $r$  increases as  $\gamma$  does. This observation is at the  
 227 core of the ability of proximal operators to identify neighboring structure manifolds.

228 **2.3. Structure Identification.** It is well-known that the proximity operator  
 229 identifies structure locally around critical points (see *e.g.* [11, Th. 4]): all points near  
 230 a minimizer are mapped to the manifold containing the minimizer. Furthermore, this  
 231 structure is revealed during the computation of the operator.<sup>1</sup>

232 In the situation we consider, the proximity operator of  $F$  cannot be explicitly  
 233 computed. However,  $\mathbf{prox}_{\gamma g}$  is available and can provide some structure in the inter-  
 234 mediate space  $\mathbb{R}^m$  that we would like to exploit. To do so, we introduce two properties  
 235 (holding on our two running examples), that will allow us to retrieve the structural  
 236 information in the intermediate space near points that are not minimizers of  $g$ .

237 The first property holds at point  $\bar{y} \in \mathcal{M}^g$  if the nonsmooth function  $g$  strictly  
 238 increases on all directions on which it is nonsmooth.

239 **PROPERTY 2.7 (Normal ascent).** *A function  $g$  satisfies the normal ascent prop-*  
 240 *erty at point  $\bar{y}$  if  $0$  lies in the relative interior of the projection of  $\partial g(\bar{y})$  on the normal*  
 241 *space at  $\bar{y}$ , that is:  $0 \in \text{ri proj}_{N_{\bar{y}}\mathcal{M}^g} \partial g(\bar{y})$ .*

242 **Remark 2.8 (Positive directional derivative).** In a “nice” setting where  $g$  is Lip-  
 243 schitz continuous and regular at  $\bar{y}$ , **Property 2.7** implies that the directional derivative  
 244 of  $g$  along any normal direction  $d \in N_{\bar{y}}\mathcal{M}^g$  is positive. Indeed, in that case one-sided  
 245 directional derivatives are well-defined [26, p. 358, Th. 9.16], and the derivative along  
 246 direction  $w$  equals  $\max_{v \in \partial g(x)} \langle v, w \rangle$ . Along a normal direction  $d \in N_{\bar{y}}\mathcal{M}^g$ , by par-  
 247 tial smoothness the directional derivative writes  $\max_{v_n \in \text{proj}_{N_{\bar{y}}\mathcal{M}^g}(\partial g(x))} \langle v_n, d \rangle$ . **Prop-**  
 248 **erty 2.7** ensures the existence of  $\alpha > 0$  such that  $\alpha d \in \text{proj}_{N_{\bar{y}}\mathcal{M}^g}(\partial g(x))$ , making the  
 249 derivative positive.

250 The second property is more technical and controls the velocity of a curve on the  
 251 manifold  $\mathcal{M}^g$ .

<sup>1</sup>Computing exactly the *structure* of the output point of the operator, as can be done for the prox, is opposed to merely observing the structure of the output after its computation. This last option is not desirable in our opinion as it entails delicate numerical questions such as testing equality between reals for the maximum, or computing the multiplicity of the maximal eigenvalue of a matrix.



252 **PROPERTY 2.9** (Curve property). *A function  $g$  partly smooth at  $\bar{y}$  relative to  $\mathcal{M}^g$*   
 253 *satisfies the curve property at  $\bar{y}$  when there exists a neighborhood  $\mathcal{N}_{\bar{y}}$  of  $\bar{y}$  and  $T > 0$*   
 254 *such that any smooth application  $e : \mathcal{N}_{\bar{y}} \times [0, T] \rightarrow \mathcal{M}^g$  such that  $e(y, 0) = \text{proj}_{\mathcal{M}^g}(y)$ ,*  
 255  *$\frac{d}{dt}e(y, t)|_{t=0} = -\text{grad } g(\text{proj}_{\mathcal{M}^g}(y))$  satisfies*

$$256 \quad \|\text{proj}_{N_{e(y,t)}\mathcal{M}^g}(e(y, t) - y)\| \leq \text{dist}_{\mathcal{M}^g}(y) + \tilde{L} t^2 \quad \text{for all } y \in \mathcal{N}_{\bar{y}}, t \in [0, T],$$

258 *where  $\text{dist}_{\mathcal{M}^g}(y) \triangleq \|y - \text{proj}_{\mathcal{M}^g}(y)\|$  is the distance between  $\mathcal{M}^g$  and  $y$ , and  $\text{grad } g(p) \in$*   
 259  *$T_p\mathcal{M}^g$  denotes the Riemannian gradient of  $g$  obtained as  $\text{proj}_{T_p\mathcal{M}^g}(\partial g(p))$ .*

260 The idea behind this property is to ensure that the differential of the projection  
 261 of the (time dependent) normal space is (uniformly) negligible at time 0. Note that  
 262 for affine spaces, we trivially have  $\|\text{proj}_{N_{e(y,t)}\mathcal{M}^g}(y - e(y, t))\| = \text{dist}_{\mathcal{M}^g}(y)$  for all  $t$   
 263 near 0: the normal spaces are equal at all points of the manifold.

264 These two properties are satisfied at any structured point for the two nonsmooth  
 265 functions  $\max$  and  $\lambda_{\max}$  of our running examples as detailed in the following lemma.  
 266 The proofs for the two functions are rather direct but require precise technical de-  
 267 scriptions; we defer them to [Appendix A](#).

268 **LEMMA 2.10.** *Consider either:*

- 269 •  $g = \max$ ,  $\bar{y} \in \mathbb{R}^m$ , and the structure manifold  $\mathcal{M}_T^{\max}$  (of [Example 2.5](#));
- 270 •  $g = \lambda_{\max}$ ,  $\bar{y} \in \mathbb{S}_m$ , and the structure manifold  $\mathcal{M}_T^{\lambda_{\max}}$  (of [Example 2.6](#)).

271 *Then, [Properties 2.7](#) and [2.9](#) hold at  $\bar{y}$ .*

272 Finally, the structure provided by  $\text{prox}_{\gamma g}$  lies in the intermediate space  $\mathbb{R}^m$ , while  
 273 the optimization variable lives in  $\mathbb{R}^n$ . In order to transfer the structure information to  
 274 the input space, we will also require the smooth map  $c : \mathbb{R}^n \rightarrow \mathbb{R}^m$  to be transversal  
 275 to  $\mathcal{M}^g \subset \mathbb{R}^m$  at some point  $\bar{x} \in \mathbb{R}^n$ , which holds when  $\mathcal{M}^g$  is a manifold around  $c(\bar{x})$   
 276 and the following (equivalent) conditions hold:

$$277 \quad \ker(\text{D } c(\bar{x})^*) \cap N_{c(\bar{x})}\mathcal{M}^g = \{0\} \quad \text{or} \quad T_{c(\bar{x})}\mathcal{M}^g + \text{Im } \text{D } c(\bar{x}) = \mathbb{R}^m.$$

279 In that case, the set  $c^{-1}(\mathcal{M}^g)$  is a submanifold of  $\mathbb{R}^n$  [[15](#), Th. 6.30], whose normal  
 280 space has the same dimension as the one of  $\mathcal{M}^g$ . Furthermore, we have [[15](#), Ex. 6-10]

$$281 \quad N_x c^{-1}(\mathcal{M}^g) = \text{D } c(x)^* N_{c(x)}\mathcal{M}^g \quad \text{and} \quad T_x c^{-1}(\mathcal{M}^g) = \text{D } c(x)^{-1} T_{c(x)}\mathcal{M}^g.$$

283 **3. Collecting structure with the proximity operator.** We show in this  
 284 section how to exactly detect the optimal structure manifold of the composite function  
 285  $F = g \circ c$  around a point  $\bar{x}$  using the proximity operator of  $g$ .

286 In our nonconvex and nonsmooth setting, we seek only structured points which  
 287 satisfy certain assumptions summarized in our definition of a *qualified point*.

288 **DEFINITION 3.1** (Qualified points). *A point  $\bar{x} \in \mathbb{R}^n$  is qualified relative to a*  
 289 *decomposition  $(g, c)$  of  $F$  and manifold  $\mathcal{M}^g$  if*

- 290 *i)  $g$  is prox-bounded and prox-regular at  $c(\bar{x})$ ;*
- 291 *ii)  $g$  is partly smooth at  $c(\bar{x})$  relative to  $\mathcal{M}^g$ ;*
- 292 *iii)  $c$  is transversal to  $\mathcal{M}^g$  at  $\bar{x}$ ;*
- 293 *iv)  $g$  satisfies [Properties 2.7](#) and [2.9](#) at point  $c(\bar{x})$ .*

294 Three of these assumptions constrain only the nonsmooth function  $g$  and are  
 295 easily verifiable in practice. Only the transversality condition limits the range of  
 296 acceptable smooth mappings; see *e.g.* [[18](#), Sec. 4]. For such *qualified* points, we get  
 297 two useful properties: first,  $F$  is partly smooth at  $\bar{x}$  relative to  $\mathcal{M} = c^{-1}(\mathcal{M}^g) \ni \bar{x}$   
 298 (by the chain rule of [[18](#), Th. 4.2]) and second, the operator  $\text{prox}_{\gamma g}$  is single-valued,  
 299 locally Lipschitz, and defined by its optimality condition near  $c(\bar{x})$ .



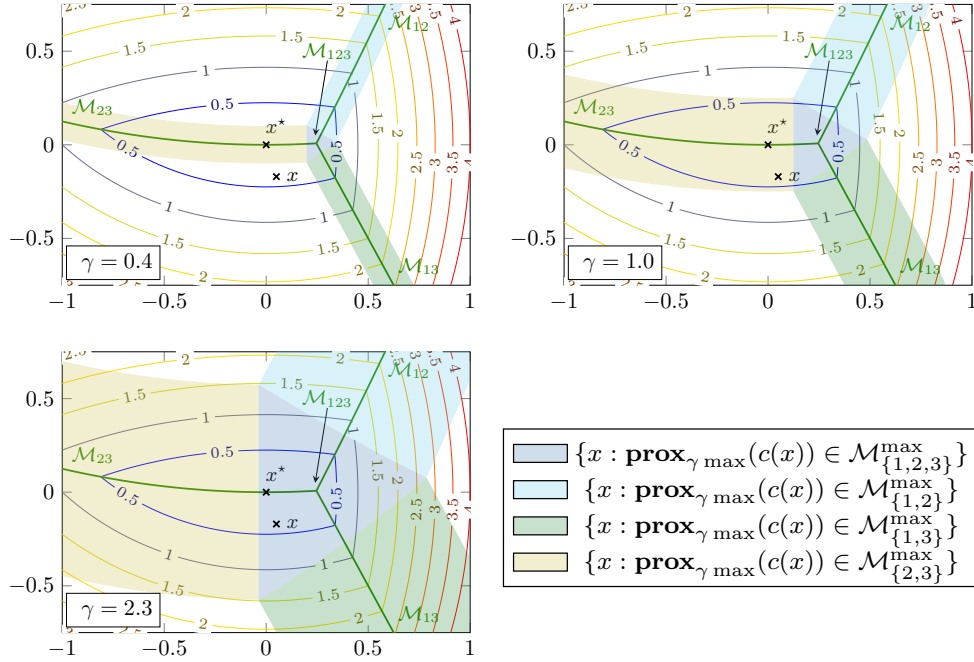


Fig. 2: Illustration of the main result on a maximum of three quadratic functions, with  $\bar{x} \in \mathcal{M}_{\{1,2\}}^{\max}$  and a point  $\tilde{x}$  near  $\bar{x}$ . The three figures show the areas where  $\mathbf{prox}_{\gamma g} \circ c$  detects manifolds for three stepsizes:  $\gamma = 0.4$  (upper left),  $\gamma = 1$  (upper right) and  $\gamma = 2.3$  (lower left). We see on the upper left fig. that  $\mathbf{prox}_{\gamma g} \circ c$  detects no structure from  $\tilde{x}$  because  $\gamma$  is too small, and in contrast, on the lower fig., that it wrongly detects too much structure ( $\mathcal{M}_{\{1,2,3\}}^{\max}$ ) because  $\gamma$  is too large. On the upper right fig., the optimal manifold is detected with  $\gamma$  chosen in the right interval.

300 **3.1. Main result.** We show in the following theorem that if  $x$  is near a qualified  
 301 point of  $F$  with structure  $\mathcal{M}$ , then  $\mathbf{prox}_{\gamma g}(c(x))$  will output a point on  $\mathcal{M}^g = c(\mathcal{M})$ ,  
 302 the structure manifold of  $g$  corresponding to  $\mathcal{M}$  (in the intermediate space). Our  
 303 theorem provides precise conditions on  $x$  and  $\gamma$  that guarantee this structure identifica-  
 304 tion and forms the main theoretical contribution of the paper. We illustrate this  
 305 behavior in Figure 2.

306 The position of this result with respect to the literature is discussed right after  
 307 in Remark 3.3, and the proof is given in the following Subsection 3.2, in a succession  
 308 of technical lemmas. We stress that we give guarantees on the *structure* to which the  
 309 point  $\mathbf{prox}_{\gamma g}(c(x))$  belongs, rather than on the point itself.

310 **THEOREM 3.2.** Consider a function  $F = g \circ c$  and a point  $\bar{x}$ . Assume that  $\bar{x}$  is  
 311 qualified relative to a manifold  $\mathcal{M}^g \subset \mathbb{R}^m$ . Then, there exists a neighborhood  $\mathcal{N}_{\bar{x}}$  of  
 312  $\bar{x}$  and a constant  $\Gamma$  such that, for all  $x \in \mathcal{N}_{\bar{x}}$ ,

$$313 \quad \mathbf{prox}_{\gamma g}(c(x)) \in \mathcal{M}^g \text{ for all } \gamma \in [\varphi(\text{dist}_{\mathcal{M}}(x)), \Gamma],$$

315 where  $\text{dist}_{\mathcal{M}}(x)$  denotes the distance from  $x$  to the manifold  $\mathcal{M}$  and  $\varphi$  is defined as

$$316 \quad \varphi(t) = \frac{c_{ri}}{2\tilde{L}} \left( 1 - \sqrt{1 - \frac{4\tilde{L}c_{map}t}{c_{ri}^2}} \right) = \frac{c_{map}}{c_{ri}}t + \frac{\tilde{L}c_{map}^2}{c_{ri}^3}t^2 + o(t^2),$$

317

318 with  $c_{ri}$ ,  $c_{map}$ , and  $\tilde{L}$  (of [Property 2.9](#)) positive constants.

319 In particular, there exists  $L > 0$ ,  $\epsilon > 0$  such that

$$320 \quad \|x - x^*\| \leq \epsilon \text{ and } L\|x - x^*\| \leq \gamma \leq \Gamma \implies \mathbf{prox}_{\gamma g}(c(x)) \in \mathcal{M}^g.$$

322 Note that [Property 2.9](#) is only used to compute explicitly an interval of  $\gamma$  guaranteed  
323 to provide the correct structure; the existence of that interval holds independently.

324 *Remark 3.3* (Relation with existing results). The difference between [Theorem 3.2](#)  
325 and existing results lies in two aspects. First, the identification properties of the  
326 proximal operator [8, Th. 28], the proximal-gradient operator [1, Th. 3.1], or even  
327 approximate prox-gradient operators [14] give structure information directly in the  
328 input space (even in abstract algorithmic frameworks [11, Th. 4]). In the composite  
329 case, the proximity operator reveals structure in the intermediate space only, and  
330 extra work is required to bring it back to the input space.

331 Second, most existing results investigate identification properties near minimizers,  
332 and not just arbitrary points (a notable exception is [1] in a different context). Here,  
333 we evaluate  $\mathbf{prox}_{\gamma g}$  near  $c(\bar{x})$ , a point without any specific properties even if  $\bar{x}$  is a  
334 local minimizer. This is why we need [Property 2.7](#) to guarantee identification in the  
335 intermediate space, and bring the structure information to the input space.

336 **3.2. Proof of [Theorem 3.2](#).** The proof consists in giving conditions on  $y$  and  
337  $\gamma$  so that  $p = \mathbf{prox}_{\gamma g}(y)$  lies on the manifold. We characterize this relation by the  
338 first-order optimality condition:

$$339 \quad y \in p + \gamma \partial g(p).$$

341 We first show that, for  $y$  near  $\bar{y}$  and  $\gamma$  small, there is a unique point  $p$  on the manifold  
342  $\mathcal{M}^g$  that satisfies the tangent component of this optimality condition:

$$343 \quad (3.1) \quad \text{proj}_{T_p \mathcal{M}^g}(y - p) = \gamma \text{grad } g(p),$$

345 where  $\text{grad } g(p) \triangleq \text{proj}_{T_p \mathcal{M}^g} \partial(g(p))$  is unique by the sharpness property of partial  
346 smoothness, and matches the Riemannian gradient of  $g$  on  $\mathcal{M}^g$  (see [6, Sec. 7.7]).  
347 Such points  $p$  are given by a smooth manifold-valued application  $e(y, \gamma)$ , the existence  
348 of which is guaranteed by the following lemma.

349 **LEMMA 3.4.** *Consider a function  $g : \mathbb{R}^m \rightarrow \mathbb{R} \cup \{+\infty\}$ , a point  $\bar{y} \in \mathbb{R}^m$ , and a*  
350 *manifold  $\mathcal{M}^g$  with  $g$  partly smooth at  $\bar{y}$  relative to  $\mathcal{M}^g$ . Then, there exists a smooth*  
351 *curve  $e : \mathcal{N}_{\bar{y}} \times \mathcal{N}_0 \rightarrow \mathcal{M}$  defined on a neighborhood of  $(\bar{y}, 0)$  in  $\mathbb{R}^m \times \mathbb{R}_+$  such that*

- 352 • for all  $y \in \mathcal{N}_{\bar{y}}$ ,  $e(y, 0) = \text{proj}_{\mathcal{M}^g}(y)$  and  $\frac{d}{d\gamma} e(y, \gamma)|_{\gamma=0} = -\text{grad } g(\text{proj}_{\mathcal{M}^g}(y))$ ;
- 353 • for all  $y \in \mathcal{N}_{\bar{y}}$ ,  $\gamma \in \mathcal{N}_0$ , *Eq. (3.1) is satisfied for  $p = e(y, \gamma)$ .*

354 *Proof.* We define the mapping  $\Phi : \mathbb{R}^m \times \mathbb{R} \times \mathcal{M}^g \rightarrow \cup_{x \in \mathcal{M}^g} T_x \mathcal{M}^g$  as

$$355 \quad \Phi(y, \gamma, p) = \gamma \text{grad } g(p) - \text{proj}_{T_p \mathcal{M}^g}(y - p)$$

356

357 and consider the equation  $\Phi(y, \gamma, p) = 0$  near the point  $(\bar{y}, 0, \bar{y})$ . Using the smoothness  
358 of  $g$  on  $\mathcal{M}^g$  given by partial smoothness, we have that this mapping is continuously

359 differentiable on a neighborhood of  $(\bar{y}, 0, \bar{y})$ . We see that its differential with respect  
 360 to  $p$  is  $D_p \Phi(\bar{y}, 0, \bar{y}) = I$ . Indeed, for  $\eta \in T_p \mathcal{M}^g$ ,

$$361 \quad D_p \Phi(y, \gamma, p)[\eta] = \gamma \text{Hess } g(p)[\eta] + \eta - D_{p'} \left( p' \mapsto \text{proj}_{T_{p'} \mathcal{M}^g}(y - p) \right) (p)[\eta].$$

363 At point  $(\bar{y}, 0, \bar{y})$ , the first term vanishes, and the third term writes

$$364 \quad D_{p'} \left( p' \mapsto \text{proj}_{T_{p'} \mathcal{M}^g}(0) \right) (\bar{y})[\eta]$$

366 and vanishes as well as the differential of the null function  $p' \mapsto \text{proj}_{T_{p'} \mathcal{M}^g}(0)$ . Thus  
 367  $D_p \Phi(\bar{y}, 0, \bar{y}) = I$  is invertible. The implicit functions theorem thus grants the existence  
 368 of neighborhoods  $\mathcal{N}_{\bar{y}}^1, \mathcal{N}_0^2, \mathcal{N}_{\bar{y}}^3$  of  $\bar{y}, 0, \bar{y}$  in  $\mathbb{R}^m, \mathbb{R}, \mathcal{M}^g$  and a continuously  
 369 differentiable function  $c : \mathcal{N}_{\bar{y}}^1 \times \mathcal{N}_0^2 \rightarrow \mathcal{N}_{\bar{y}}^3$  such that, for any  $(y, \gamma) \in \mathcal{N}_{\bar{y}}^1 \times \mathcal{N}_0^2$ , Equation  
 370 (3.1) is satisfied with  $p = e(y, \gamma)$ . For  $y \in \mathcal{N}_{\bar{y}}^1$ ,  $e(y, 0)$  satisfies  $y - e(y, 0) \in$   
 371  $N_{e(y, 0)} \mathcal{M}^g$ , which is the first-order optimality condition of  $e(y, 0) = \text{proj}_{\mathcal{M}^g}(y)$ . Possibly  
 372 reducing  $\mathcal{N}_{\bar{y}}^1$  so that, for all  $y \in \mathcal{N}_{\bar{y}}^1$ ,  $\text{proj}_{\mathcal{M}^g}(y)$  is well-defined and unique, the  
 373 previous optimality condition is equivalent to  $e(y, 0) = \text{proj}_{\mathcal{M}^g}(y)$ . Besides, differentiating  
 374  $\Phi(y, \gamma, e(y, \gamma)) = 0$  relative to  $\gamma$  at  $\gamma = 0$  yields

$$376 \quad D_\gamma e(y, 0) = -[D_p \Phi(y, 0, \text{proj}_{\mathcal{M}^g}(y))]^{-1} D_\gamma \Phi(y, 0, \text{proj}_{\mathcal{M}^g}(y)) = -\text{grad } g(\text{proj}_{\mathcal{M}^g}(y)),$$

377 which concludes the proof.  $\square$

378 The previous lemma shows that for every  $(y, \gamma)$  one can find a point  $e(y, \gamma)$  on the  
 379 manifold  $\mathcal{M}^g$  that solves the tangent part of the optimality condition (3.1). The next  
 380 lemma determines the values of  $y$  and  $\gamma$  for which the whole optimality condition

$$381 \quad (3.2) \quad y \in e(y, \gamma) + \gamma \text{ri } \partial g(e(y, \gamma))$$

383 holds, as illustrated in Figure 3a.

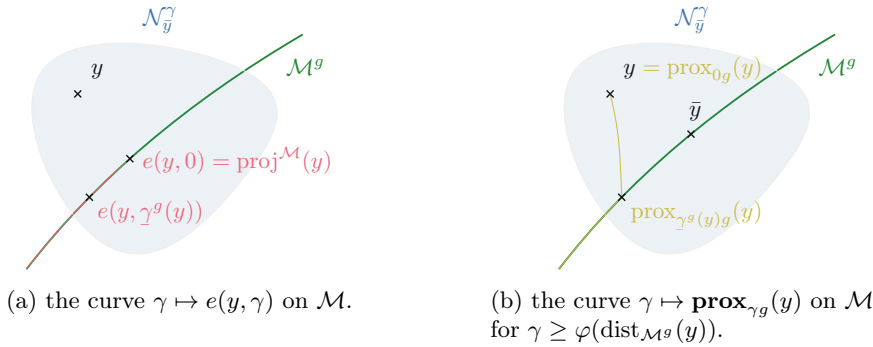


Fig. 3: Illustration of Lemma 3.5 and its consequences.

384 LEMMA 3.5. Consider a function  $g$ , a point  $\bar{y} \in \mathbb{R}^m$  and a manifold  $\mathcal{M}^g$  such  
 385 that  $g$  is partly smooth at  $\bar{y}$  relative to  $\mathcal{M}^g$  and that  $g$  satisfies Property 2.7 at  $\bar{y}$ . Let  
 386  $e$  denote a smooth  $\mathcal{M}$ -valued application defined on a neighborhood of  $(\bar{y}, 0)$  provided  
 387 by Lemma 3.4. Then, there exists  $C > 0$  such that:

- 388 i) for all  $\gamma \in [0, C]$ ,  $e(\bar{y}, \gamma)$  verifies (3.2) with  $y = \bar{y}$ ,  
 389 ii) for all  $\gamma \in [0, C]$ , there exists a neighborhood  $\mathcal{N}_{\bar{y}}^\gamma$  of  $\bar{y}$  such that, for all  
 390  $y \in \mathcal{N}_{\bar{y}}^\gamma$ ,  $e(y, \gamma)$  verifies (3.2),

391 Further assume that  $g$  satisfies [Property 2.9](#) at  $\bar{y}$  with constant  $\tilde{L}$ , then  
 392 *iii)* there exist  $\Gamma^g > 0$  and a neighborhood  $\mathcal{N}_{\bar{y}}$  of  $\bar{y}$  such that for all  $y \in \mathcal{N}_{\bar{y}}$

393  $e(y, \gamma)$  verifies [\(3.2\)](#) for all  $\gamma \in [\varphi^g(\text{dist}_{\mathcal{M}^g}(y)), \Gamma^g]$ ,

395 where  $c_{\text{ri}} \geq 0$  and  $\varphi^g(t) = \frac{c_{\text{ri}}}{2\tilde{L}} \left( 1 - \sqrt{1 - \frac{4\tilde{L}t}{c_{\text{ri}}^2}} \right) = \frac{1}{c_{\text{ri}}}t + \frac{\tilde{L}}{c_{\text{ri}}^3}t^2 + o(t^2)$ .

396 The proof consists in finding the points  $y, \gamma$  such that  $0 \in \text{ri } \Psi(y, \gamma)$ , where the  
 397 mapping  $\Psi : \mathbb{R}^m \times \mathbb{R} \rightarrow \cup_{x \in \mathcal{M}^g} N_x \mathcal{M}^g$  is defined as

$$398 \quad \Psi(y, \gamma) = \text{proj}_{N_{e(y, \gamma)} \mathcal{M}^g} \left( \frac{1}{\gamma} (e(y, \gamma) - y) + \partial g(e(y, \gamma)) \right).$$

399

400 Items *i)* and *ii)* are shown by extending the property  $0 \in \Psi(\bar{y}, 0)$  to a neighborhood of  
 401  $(\bar{y}, 0)$ , using the inner-semicontinuity properties of  $\Psi$ . We then derive explicit bounds  
 402 on the interval of steps such that  $0 \in \text{ri } \Psi(y, \gamma)$ : for a fixed  $y \in \mathcal{N}_{\bar{y}}$ , when  $\gamma$  decreases  
 403 past some value, say  $\underline{\gamma}(y)$ , the condition  $0 \in \text{ri } \Psi(y, \gamma)$  no longer holds. Precisely at  
 404  $\underline{\gamma}(y)$ ,  $0$  lies on the (relative) boundary of  $\Psi(y, \underline{\gamma}(y))$ : denoting  $\text{rbd } S \triangleq S \setminus \text{ri } S$  the  
 405 relative boundary of set  $S$ ,

$$406 \quad 0 \in \text{rbd } \text{proj}_{N_{e(y, \underline{\gamma}(y))} \mathcal{M}^g} \left( \frac{1}{\underline{\gamma}(y)} (e(y, \underline{\gamma}(y)) - y) + \partial g(e(y, \underline{\gamma}(y))) \right).$$

407

408 Denoting  $\partial^N g(p) \triangleq \text{proj}_{N_p \mathcal{M}^g}(\partial g(p))$  the projection of the subdifferential on the  
 409 normal space of its structure manifold and taking norms yields:

$$410 \quad \|\text{proj}_{N_{e(y, \underline{\gamma}(y))} \mathcal{M}^g}(y - e(y, \underline{\gamma}(y)))\| \geq \underline{\gamma}(y) \inf_{v_n \in \text{rbd } \partial^N g(e(y, \underline{\gamma}(y)))} \|v_n\|$$

$$411 \quad \geq \underline{\gamma}(y) \underbrace{\inf_{p \in \mathcal{N}_{\bar{y}}} \inf_{v_n \in \text{rbd } \partial^N g(p)} \|v_n\|}_{\triangleq c_{\text{ri}}}.$$

412

413 Since  $0 \in \text{ri } \text{proj}_{N_{\bar{y}} \mathcal{M}^g} \partial g(\bar{y})$  and  $\partial g$  is inner-semicontinuous, the former property  
 414 actually holds on a neighborhood of  $\bar{y}$  in  $\mathcal{M}^g$ , thus making the constant  $c_{\text{ri}}$  positive.  
 415 We note that this kind of quantity also appears as the *modulus of identifiability* in  
 416 the recent [\[16, Def. 2.3\]](#) where it has the same property: its positivity enables the  
 417 identification of the associated structure manifold.

418 Using [Property 2.9](#), the left-hand side is upper bounded by a simpler expression:

$$419 \quad \tilde{L}\underline{\gamma}(y)^2 + \text{dist}_{\mathcal{M}^g}(y) \geq c_{\text{ri}}\underline{\gamma}(y), \quad \text{that is} \quad \underline{\gamma}(y) \leq \frac{c_{\text{ri}}}{2\tilde{L}} \left( 1 - \sqrt{1 - \frac{4\tilde{L} \text{dist}_{\mathcal{M}^g}(y)}{c_{\text{ri}}^2}} \right),$$

420

421 which provides the expression for  $\varphi^g$  used in the lemma.

422 *Proof. Item i)* We first consider  $\Psi_{\bar{y}}(\cdot) = \Psi(\bar{y}, \cdot)$ . Since  $\bar{y} \in \mathcal{M}^g$ , [Lemma 3.4](#) tells  
 423 us that  $e(\bar{y}, \gamma) = \bar{y} - \gamma \text{grad } g(\bar{y}) + o(\gamma)$ , and thus

$$424 \quad \Psi_{\bar{y}}(0) = \text{proj}_{N_{\bar{y}} \mathcal{M}^g} (-\text{grad } g(\bar{y}) + \partial g(\bar{y})) = \text{proj}_{N_{\bar{y}} \mathcal{M}^g}(\partial g(\bar{y}))$$

425

426 where we used that  $\text{grad } g(\bar{y}) \in T_{\bar{y}} \mathcal{M}^g$  is orthogonal to  $N_{\bar{y}} \mathcal{M}^g$ .

427 **Property 2.7** provides that  $0 \in \text{ri } \Psi_{\bar{y}}(0)$ . We now turn to showing that there exists  
 428  $C'$  such that, for all  $\gamma \in [0, C']$ ,  $0 \in \text{ri } \Psi_{\bar{y}}(\gamma)$ .

429 By contradiction, assume there exist a sequence  $\gamma_k \rightarrow 0$  such that  $0 \notin \text{ri } \Psi_{\bar{y}}(\gamma_k)$ .  
 430 This means that there exists a sequence of unit norm vectors  $(s_k)$  such that for all  $k$ ,

$$431 \quad (3.3) \quad \langle s_k, z \rangle \leq 0 \text{ for all } z \in \Psi_{\bar{y}}(\gamma_k).$$

433 As a bounded sequence,  $s_k$  admits at least one limit point, say  $\bar{s}$ . Take  $\bar{z} \in \Psi_{\bar{y}}(0)$ . The  
 434 continuity of  $\partial g$  (by partial smoothness, item iv), of  $\gamma \mapsto (e(\bar{y}, \gamma) - \bar{y})/\gamma$  (by smooth-  
 435 ness of  $e$ ), and of  $\gamma \mapsto \text{proj}_{N_{e(\bar{y}, \gamma)} \mathcal{M}^g}$  (by smoothness of  $\mathcal{M}^g$ ) yield the continuity of  
 436  $\Psi_{\bar{y}}$  as a set-valued map. This mapping is thus inner-semicontinuous [26, Def. 5.4], so  
 437 there exists a sequence  $z_k \in \Psi_{\bar{y}}(\gamma_k)$  such that  $z_k$  converges to  $\bar{z}$ . Taking the correct  
 438 subsequence and renaming iterates, we can write  $s_k \rightarrow \bar{s}$  and  $z_k \rightarrow \bar{z}$ . Equation (3.3)  
 439 provides  $\langle s_k, z_k \rangle \leq 0$  for all  $k$ , which gives at the limit  $\langle \bar{s}, \bar{z} \rangle \leq 0$ . This actually holds  
 440 for all  $\bar{z} \in \Psi_{\bar{y}}(0)$ :  $\bar{s}$  separates 0 and  $\Psi(0)$ , which contradicts  $0 \in \text{ri } \Psi_{\bar{y}}(0)$ .

441 Finally, let us take the constant  $C$  such that  $[0, C]$  is included in  $[0, C']$  and the  
 442 neighborhood of 0 provided by **Lemma 3.4**. Then, for any  $\gamma \in [0, C]$ , adding the  
 443 two orthogonal inclusions  $0 \in \text{ri } \Psi_{\bar{y}}(\gamma)$  and  $0 = \Phi(y, \gamma, c(y, \gamma))$ , we obtain that  $e(\bar{y}, \gamma)$   
 444 verifies (3.2) with  $y = \bar{y}$ .

445 *Item ii)* Let  $\gamma \in [0, C]$ . We turn to show the existence of a neighborhood  $\mathcal{N}_{\bar{y}}^\gamma$  of  $\bar{y}$  such  
 446 that, for all  $y \in \mathcal{N}_{\bar{y}}^\gamma$ ,  $e(y, \gamma)$  verifies (3.2). By contradiction, assume that there exists  
 447 a sequence  $(y_k)$  that converges to  $\bar{y}$  such that (3.2) fails for  $(y_k, \gamma)$ . Since the tangent  
 448 component of (3.2) does hold, necessarily  $0 \notin \text{ri } \Psi(y_k, \gamma)$ . However, the mapping  
 449  $y \mapsto \Psi(y, \gamma)$  is inner-semicontinuous (from the same arguments as in the proof of *item*  
 450 *i)* and there holds  $0 \in \text{ri } \Psi(\bar{y}, \gamma)$ . A reasoning similar to that of *item i)* reveals the  
 451 contradiction.

452 *Item iii)* Define  $\mathcal{N}_{\bar{y}}$  a neighborhood of  $\bar{y}$  and  $\Gamma^g$  a positive constant such that **Prop-**  
 453 **erty 2.9** applies over  $\mathcal{N}_{\bar{y}} \times [0, \Gamma^g]$ ,  $\mathcal{N}_{\bar{y}}$  is contained in  $\cup_{\gamma \in [0, C]} \mathcal{N}_{\bar{y}}^\gamma \cap \mathcal{N}_{\bar{y}}^C$  and  $0 \in \text{ri } \Psi(y, \gamma)$   
 454 holds for all  $(y, \gamma) \in \mathcal{N}_{\bar{y}} \times [0, \Gamma^g]$ . This last condition can be met on a nontrivial neigh-  
 455 borhood of  $(\bar{y}, 0)$  since it holds at that point, and  $\Psi$  is inner-semicontinuous ( $e(y, \gamma)$   
 456 lies on  $\mathcal{M}^g$  and  $\partial g$  is inner-semicontinuous by partial smoothness of  $g$ ).

457 Let  $y \in \mathcal{N}_{\bar{y}}$  and  $\gamma > 0$  such that  $\varphi^g(\text{dist}_{\mathcal{M}^g}(y)) \leq \gamma \leq \Gamma^g$ . The lower bound on  
 458  $\gamma$  implies that  $\tilde{L}\gamma^2 + \text{dist}_{\mathcal{M}}(y) < c_{\text{ri}}\gamma$ . We have successively by **Property 2.9** and the  
 459 above bound that

$$460 \quad \|\text{proj}_{N_{e(y, \gamma)} \mathcal{M}^g}(y - e(y, \gamma))\| \leq \text{dist}_{\mathcal{M}}(y) + \tilde{L}\gamma^2 < c_{\text{ri}}\gamma$$

$$461 \quad \leq \gamma \inf\{\|n\|, n \in \text{rd } \partial^N g(e(y, \gamma))\}.$$

463 This means that  $\text{proj}_{N_{e(y, \gamma)} \mathcal{M}^g}(y - e(y, \gamma))$  belongs to the ball of center 0 and radius  
 464  $\gamma \inf\{\|n\|, n \in \text{rd } \partial^N g(e(y, \gamma))\}$  in  $N_{e(y, \gamma)} \mathcal{M}^g$ . In addition, this ball is included in  
 465  $\partial^N g(e(y, \gamma))$  since  $0 \in \partial^N g(e(y, \gamma))$  by definition of  $\mathcal{N}_{\bar{y}}$ . Therefore,  $0 \in \text{ri } \Psi(y, \gamma)$  for all  
 466  $y \in \mathcal{N}_{\bar{y}}$  and  $\gamma \in [\varphi^g(\text{dist}_{\mathcal{M}^g}(y)), \Gamma^g]$ .  $\square$

467 We can now proceed to the proof of **Theorem 3.2**. It consists in first identifying  
 468 the curve  $e(y, \gamma)$  to  $\text{prox}_{\gamma g}(y)$  and thus prove that it belongs to the sought manifold,  
 469 as illustrated in **Figure 3b**. Then, this intermediate identification result is brought  
 470 back to the input space using transversality.

471 *Proof.* The standing assumptions allow to call **Lemma 3.5** at point  $c(\bar{x})$  with  
 472 manifold  $\mathcal{M}^g$ . This yields the neighborhood  $\mathcal{N}_{c(\bar{x})}$ , constants  $\Gamma^g$  and  $C$ , a function

473  $\varphi^g$ , and a smooth mapping  $e : \mathcal{N}_{c(\bar{x})} \times [0, C] \rightarrow \mathcal{M}^g$  such that, for  $y \in \mathcal{N}_{c(\bar{x})}$  and  
 474  $\gamma \in [\varphi^g(\text{dist}_{\mathcal{M}^g}(y)), \Gamma^g]$ ,  $e(y, \gamma)$  verifies the optimality condition (3.2) of  $e(y, \gamma) =$   
 475  $\mathbf{prox}_{\gamma g}(y)$ . Besides, since  $g$  is prox-regular and prox-bounded at point  $c(\bar{x})$ , these  
 476 properties also hold on a neighborhood of that point. Under these conditions, **Propo-**  
 477 **sition 2.2** allows to recover the equality  $e(y, \gamma) = \mathbf{prox}_{\gamma g}(y)$ . Take  $\mathcal{N}_{\bar{x}} = c^{-1}(\mathcal{N}_{c(\bar{x})})$ ,  
 478 a neighborhood of  $\bar{x}$  as the preimage of a neighborhood of  $c(\bar{x})$  by the continuous  $c$ .  
 479 For all  $x \in \mathcal{N}_{\bar{x}}$ ,

$$480 \quad \mathbf{prox}_{\gamma g}(c(x)) \in \mathcal{M}^g \text{ for all } \gamma \in [\varphi^g(\text{dist}_{\mathcal{M}^g}(c(x))), \Gamma^g].$$

482 We turn to show that, for some constant  $c_{\text{map}} > 0$ , there holds  $\text{dist}_{\mathcal{M}^g}(c(x)) \leq$   
 483  $c_{\text{map}} \text{dist}_{\mathcal{M}}(x)$  for all  $x \in \mathcal{N}_{\bar{x}}$ . Let  $x \in \mathcal{N}_{\bar{x}}$  and  $x^{\mathcal{M}} = \text{proj}_{\mathcal{M}}(x)$ , so that  $\text{dist}_{\mathcal{M}}(x) =$   
 484  $\|x^{\mathcal{M}} - x\|$ . Using successively that  $c(x^{\mathcal{M}}) \in \mathcal{M}^g$  and smoothness of  $c$ , there holds for  
 485  $x$  near  $\bar{x}$

$$486 \quad \begin{aligned} \text{dist}_{\mathcal{M}^g}(c(x)) &\leq \|c(x) - c(x^{\mathcal{M}})\| \\ 487 \quad &\leq \|\text{Jac}_c(x^{\mathcal{M}}) \cdot (x - x^{\mathcal{M}})\| + \mathcal{O}(\|x - x^{\mathcal{M}}\|^2) \\ 488 \quad &\leq \left( \sup_{v_n \in N_{x^{\mathcal{M}}}\mathcal{M}, \|v_n\|=1} \|\text{Jac}_c(x^{\mathcal{M}}) \cdot v_n\| \right) \|x - x^{\mathcal{M}}\| + \mathcal{O}(\|x - x^{\mathcal{M}}\|^2) \\ 489 \quad &\leq \underbrace{\left( \sup_{u \in \mathcal{N}_{\bar{x}}} \sup_{v_n \in N_u\mathcal{M}, \|v_n\|=1} \|\text{Jac}_c(u) \cdot v_n\| \right)}_{C''} \|x - x^{\mathcal{M}}\| + \mathcal{O}(\|x - x^{\mathcal{M}}\|^2). \end{aligned}$$

490 Since  $c$  is transversal to  $\mathcal{M}^g$  at  $\bar{x}$ , its Jacobian restricted to the normal space is invertible:  
 491  $C''$  is positive. Therefore, for all  $x \in \mathcal{N}_{\bar{x}}$  and a constant  $c_{\text{map}} > C''$ , there holds  
 492  $\text{dist}_{\mathcal{M}^g}(c(x)) \leq c_{\text{map}} \text{dist}_{\mathcal{M}}(x)$ . Monotony of  $\varphi^g$  implies that  $\varphi^g(\text{dist}_{\mathcal{M}^g}(c(x))) \leq$   
 493  $\varphi^g(c_{\text{map}} \text{dist}_{\mathcal{M}}(x))$ . Hence the claimed bounds with

$$495 \quad \varphi(t) = \frac{c_{\text{ri}}}{2\bar{L}} \left( 1 - \sqrt{1 - \frac{4\tilde{L}c_{\text{map}}t}{c_{\text{ri}}^2}} \right) \quad \text{and} \quad \Gamma = \Gamma^g.$$

497 Finally, we show the existence of positive constants  $\epsilon, L$  such that

$$498 \quad \|x - \bar{x}\| \leq \epsilon \text{ and } L\|x - \bar{x}\| \leq \gamma \leq \Gamma \implies \mathbf{prox}_{\gamma g}(c(x)) \in \mathcal{M}^g.$$

500 Since  $\bar{x} \in \mathcal{M}$ ,  $\text{dist}_{\mathcal{M}}(\cdot) \leq \|\cdot - \bar{x}\|$ . By monotony and smoothness of  $\varphi$ , there exists  
 501  $L > 0$  such that  $\varphi(\text{dist}_{\mathcal{M}^g}(\cdot)) \leq L\|\cdot - x^*\|$  over  $\mathcal{B}(x^*, \epsilon)$ . Reducing  $\epsilon$  if necessary so  
 502 that  $L\epsilon < \Gamma$  yields the result.  $\square$

503 **4. Proposed method.** In this section, we use the results of [Section 3](#) to propose  
 504 an optimization method that locally identifies the structure of a minimizer and  
 505 converges quadratically to this point.

506 Recall the basic idea: if the optimal manifold  $\mathcal{M}^*$  corresponding to a minimizer  
 507  $x^*$  is known, the *nonsmooth* optimization problem turns into a *smooth constrained*  
 508 optimization problem. In turn, this problem can be solved using algorithms from  
 509 smooth constrained optimization such as Sequential Quadratic Programming.

510 Using this idea and the structure identification mechanism developed in the previ-  
 511 ous section, we propose a method which: i) uses the proximity operator of  $g$  to gather  
 512 structure in the intermediate space, ii) brings back this structure to the input space,  
 513 and iii) optimizes smoothly along the identified manifold. The resulting algorithm is  
 514 precisely described in [Subsection 4.1](#) and then analyzed in [Subsection 4.2](#).

515 **4.1. Description of the algorithm.** We proceed to describe the three steps  
 516 exposed above. The full algorithm is depicted in [Algorithm 4.1](#).

517 *Gathering structure.* We showed in [Theorem 3.2](#) that near a qualified point in  $\mathbb{R}^n$ ,  
 518 the operator  $\mathbf{prox}_{\gamma g}(c(\cdot))$  provides the optimal structure  $\mathcal{M}^{g^*}$  (in the intermediate  
 519 space  $\mathbb{R}^m$ ) for an explicit range of steps. We thus define from the current iterate  
 520  $x_k \in \mathbb{R}^n$  and stepsize  $\gamma_k$  the working manifold  $\mathcal{M}_k^g$  (in the intermediate space) as  
 521 the structure of  $\mathbf{prox}_{\gamma_k g}(c(x_k))$ . One technical point is to guarantee that, after some  
 522 time,  $\gamma_k \in [L\|x_k - x^*\|, \Gamma]$  so that the optimal manifold is identified; this is done by  
 523 decreasing  $\gamma_k$  linearly at each iteration.

524 *From the intermediate to the input space.* We now have a structure manifold  $\mathcal{M}_k^g$   
 525 in the intermediate space, and can define  $\tilde{g}_k$ , a smooth extension of  $g$  on  $\mathcal{M}_k^g$  to  $\mathbb{R}^m$ .  
 526 Using a local equation  $h_k^g$  of  $\mathcal{M}_k^g$ , we define the smooth map  $h_k = h_k^g \circ c : \mathbb{R}^n \rightarrow \mathbb{R}^{p_k}$ ,  
 527 which locally defines  $\mathcal{M}_k = c^{-1}(\mathcal{M}_k^g)$ . Similarly, a smooth extension of  $F$  on  $\mathcal{M}_k$  is  
 528 defined by  $\tilde{F}_k = \tilde{g}_k \circ c$ .

529 *Optimizing in the input space.* We can now take steps to minimize the smooth  
 530 extension  $\tilde{F}_k$  on the smooth set  $\mathcal{M}_k$  characterized by  $h_k(x) = 0$ :

$$531 \quad \min_{x \in \mathbb{R}^n} \tilde{F}_k(x) \quad \text{s.t.} \quad h_k(x) = 0.$$

533 We turn to an elementary version of the traditional second-order Sequential Quadratic  
 534 Programming methodology; see *e.g.* [5, Chap. 14]. At iteration  $k$ , the SQP direction  
 535  $d_k^{\text{SQP}}(x_k)$  at point  $x_k$  is defined as the solution of the following quadratic problem:

$$536 \quad (4.1) \quad d_k^{\text{SQP}}(x_k) = \operatorname{argmin}_{d \in \mathbb{R}^n} \quad \langle \nabla \tilde{F}_k(x_k), d \rangle + \frac{1}{2} \langle \nabla_{xx}^2 L_k(x_k, \lambda_k(x_k)) d, d \rangle$$

$$537 \quad \text{s.t.} \quad h_k(x_k) + D h_k(x_k) d = 0$$

538 where  $\nabla_{xx}^2 L_k$  denotes the Hessian of the Lagrangian  $L_k(x, \lambda) = \tilde{F}_k(x) + \langle \lambda, h_k(x) \rangle$ ,  
 539 and the multiplier  $\lambda_k(x_k)$  defined from the following least-squares problem:

$$540 \quad (4.2) \quad \lambda_k(x_k) = \operatorname{argmin}_{\lambda \in \mathbb{R}^{p_k}} \left\| \nabla \tilde{F}_k(x_k) + \sum_{i=1}^{p_k} \lambda_i \nabla h_{k,i}(x_k) \right\|^2.$$

542 Finally, we check that  $x_k + d_k^{\text{SQP}}(x_k)$  provides a functional decrease in order to  
 543 avoid degrading the iterate when the current structure is suboptimal. If the test is not  
 544 verified,  $x_k$  is not updated and  $\gamma_k$  is decreased until a satisfying structure is detected.

545 **4.2. Convergence of Algorithm 4.1.** We proceed to give the result guaran-  
 546 teeing identification and local quadratic convergence of [Algorithm 4.1](#).

547 In order to benefit from the quadratic rate of SQP, the elements of (4.1) should  
 548 have the minimal regularity typically required by smooth constrained Newton methods  
 549 (see *e.g.* [5, Th. 14.5]); we thus make the following assumption.

550 **ASSUMPTION 4.1** (Regularity of functions). *The smooth extension and the man-  
 551 ifold defining map are  $\mathcal{C}^2$  with Lipschitz second derivatives, and the Jacobian of the  
 552 constraints is full rank near the solution.*

553 In order to focus on the algorithmic originality of the method, we slightly simplify  
 554 the situation and make the two following algorithmic assumptions.

555 **ASSUMPTION 4.2** (Nonconvex stability). *The iterates of [Algorithm 4.1](#) remain  
 556 in the connex component of the sublevel set  $\{x : F(x) \leq F(x_0)\}$  that contains  $x^*$ .*



**Algorithm 4.1** General structure exploiting algorithm

**Require:** Pick  $x_0$  near a minimizer,  $\gamma_0$  large enough.

```

1: repeat
2:    $\gamma_k = \frac{\gamma_{k-1}}{2}$ 
3:   Compute  $\mathbf{prox}_{\gamma_k g}(c(x_k))$  and obtain  $\mathcal{M}_k^g$  locally defined by  $h_k^g$ 
4:    $h_k = h_k^g \circ c$  (local equation of  $\mathcal{M}_k$ ),  $\tilde{F}_k = \tilde{g}_k \circ c$  (smooth extension)
5:   Compute  $d_k^{\text{SQP}}(x_k)$  by solving (4.1)
6:   if  $F(x_k + d_k^{\text{SQP}}(x_k)) \leq F(x_k)$  then
7:      $x_{k+1} = x_k + d_k^{\text{SQP}}(x_k)$ 
8:   else
9:      $x_{k+1} = x_k$ 
10: until stopping criterion

```

557 This assumption ensures that an update that decreases the functional value remains  
558 in the neighborhood of the minimizer  $x^*$ . It is naturally satisfied when  $F$  is convex,  
559 or when  $x^*$  is a global minimizer of  $F$  and  $x_0$  is close enough to  $x^*$ .

560 ASSUMPTION 4.3 (No Maratos effect). *The iterates of Algorithm 4.1 are such*  
561 *that a step  $d$  that makes  $x+d$  quadratically closer to  $x$  yields descent:  $F(x+d) \leq F(x)$ .*

562 In smooth constrained optimization, getting closer (even at quadratic rate) to a min-  
563 imizer does not imply decrease of objective value and constraint violation (measured  
564 by a merit function). This so-called Maratos effect (see *e.g.* [5]) is one of the main  
565 difficulties in globalizing SQP schemes, which is out of the scope of the current paper.  
566 We thus assume this effect does not affect our algorithm in theory, and use in practice  
567 one of the successful refinements, as discussed in Subsection 5.2.

568 We are now ready for the main convergence result of Algorithm 4.1, which es-  
569 tablish that, after some finite time, the iterates identify exactly the optimal manifold  
570 and converge to the minimizer at a quadratic rate.

571 THEOREM 4.4 (Exact identification and quadratic convergence). *Consider a*  
572 *function  $F = g \circ c$  and  $x^*$  a strong minimizer,<sup>2</sup> qualified relative to the optimal*  
573 *manifold  $\mathcal{M}^*$ . Assume that the smooth extension  $\tilde{F}$  of  $F$  relative to  $\mathcal{M}^*$  and the*  
574 *corresponding manifold defining map  $h$  satisfy Assumption 4.1.*

575 *If  $x_0$  and  $F(x_0)$  are close enough to  $x^*$  and  $F(x^*)$ ,  $\gamma_0$  is large enough and the*  
576 *simplifying algorithmic Assumptions 4.2 and 4.3 hold, then there exists  $C > 0$  such*  
577 *that the iterates  $(x_k, \mathcal{M}_k)$  generated by Algorithm 4.1 verify:*

$$578 \quad \mathcal{M}_k = \mathcal{M}^* \quad \text{and} \quad \|x_{k+1} - x^*\| \leq C \|x_k - x^*\|^2 \quad \text{for all } k \text{ large enough.}$$

580 The proof of this result consists in two steps. We first show the existence of a  
581 neighborhood of initialization on which the proximity operator will eventually identify  
582 the optimal manifold, once the stepsize has been sufficiently decreased. From this  
583 point onward, we prove that the SQP step provides a quadratic improvement and  
584 that the stepsize policy makes the manifold identification stable.

585 *Proof. Local identification of the optimal structure.* By Theorem 3.2, there exists  
586 a ball centered around  $x^*$  of radius  $\epsilon_1 > 0$  and two positive constants  $L, \Gamma$  such that,  
587 for all  $x \in \mathcal{B}(x^*, \epsilon_1)$  and  $\gamma \in [L\|x - x^*\|, \Gamma]$ ,  $\mathbf{prox}_{\gamma g}(c(x))$  belongs to the optimal  
588 manifold  $\mathcal{M}^{g^*} = c(\mathcal{M}^*)$ .

<sup>2</sup>There exists  $\eta > 0, \epsilon > 0$  such that  $F(x) \geq F(x^*) + \eta\|x - x^*\|^2$  for all  $x \in \mathcal{B}(x^*, \epsilon)$ .

589 *Local quadratic convergence of SQP on the optimal structure.* Let us assume that the  
 590 optimal manifold has been identified. The least square multiplier  $\lambda$  is defined by the  
 591 optimality condition of (4.2):

$$592 \quad \lambda(x) = -[\text{Jac}_h(x) \text{Jac}_h(x)^\top]^{-1} \text{Jac}_h(x) \nabla \tilde{F}(x).$$

594 and since  $h$  is smooth and its Jacobian is full-rank near  $x^*$ ,  $\lambda$  is a Lipschitz continuous  
 595 function near  $x^*$ .

596 Since  $x^*$  is a strong minimizer of  $F$ , the Hessian of the Lagrangian restricted to  
 597 the tangent space is positive definite. Indeed, since  $x^*$  is a strong minimizer of  $F$   
 598 on  $\mathcal{M}^*$ , the Riemannian Hessian relative to the optimal manifold is positive definite.  
 599 With the choice of multiplier (4.2), the Riemannian Hessian is exactly the Hessian of  
 600 the Lagrangian restricted to the tangent space to  $\mathcal{M}^*$  at  $x^*$  (see [6, Sec. 7.7]), which  
 601 is thus itself positive definite.

602 Thus, using the local quadratic convergence of SQP [5, Th. 14.5], we get that there  
 603 exists a ball centered around  $x^*$  of radius  $\epsilon_2 > 0$  such that the SQP step computed at  
 604 a point  $x$  in that neighborhood relative to the optimal manifold provides a quadratic  
 605 improvement towards  $x^*$ . Reducing  $\epsilon_2$  if necessary, we can in addition have that the  
 606 convergence is at least linear with rate  $1/2$ .

607 *Initialization, identification, and quadratic convergence.* Let  $\epsilon = \min(\epsilon_1, \epsilon_2, \Gamma/(2L))$ .  
 608 We will now show that initializing with  $x_0 \in \{x : F(x) \leq F(x^*) + \eta\epsilon^2\}$  and  $\gamma_0 \geq \Gamma$   
 609 provides the claimed behavior.

610 First, the functional decrease test of the algorithm and Assumption 4.3 guarantee  
 611 that all iterates satisfy  $F(x_k) \leq F(x_0)$ . Using that  $x^*$  is a strong minimizer, we get  
 612 that  $\eta\|x_k - x^*\|^2 \leq F(x_k) - F(x^*) \leq F(x_0) - F(x^*) \leq \eta\epsilon^2$ , and thus that the iterates  
 613 remain in  $\mathcal{B}(x^*, \epsilon)$ .

614 Second, as  $L\|x - x^*\| \leq \Gamma/2$  for all  $x \in \mathcal{B}(x^*, \epsilon)$  by construction, the fact that  
 615  $\gamma_0 > \Gamma$  and  $(\gamma_k)$  decreases with geometric rate  $1/2$  implies that there exists  $K$  such  
 616 that  $L\|x_K - x^*\| \leq \gamma_K \leq \Gamma$ .

617 Now, assume that at iteration  $k \geq K$ ,  $L\|x_k - x^*\| \leq \gamma_k \leq \Gamma$ . Since  $x_k \in \mathcal{B}(x^*, \epsilon_1)$ ,  
 618 we have from above that  $\mathcal{M}^*$  is identified. Thus, the SQP step is performed relative  
 619 to the optimal manifold and  $x_k + d_k^{\text{SQP}}(x_k)$  brings a linear improvement of factor  $1/2$   
 620 at least. Assumption 4.2 ensures that  $F(x_k + d_k^{\text{SQP}}(x_k)) \leq F(x_k)$  so that  $x_{k+1} =$   
 621  $x_k + d_k^{\text{SQP}}(x_k)$  and thus

$$622 \quad L\|x_{k+1} - x^*\| \leq \frac{L}{2}\|x_k - x^*\| \leq \frac{\gamma_k}{2} = \gamma_{k+1}.$$

624 This shows that  $L\|x_{k+1} - x^*\| \leq \gamma_{k+1} \leq \Gamma$ , which completes the induction. We  
 625 get that  $\gamma_k \in [L\|x_k - x^*\|, \Gamma]$  for all  $k \geq K$ . Finally, we have that for all  $k \geq K$ ,  
 626  $\mathcal{M}_k = \mathcal{M}^*$  and  $x_{k+1}$  is quadratically closer to  $x^*$  than  $x_k$ .  $\square$

627 *Direct generalizations.* Theorem 4.4 actually holds for any decrease factor of  $\gamma_k$   
 628 in  $(0, 1)$  with the presented SQP update, or actually any superlinearly convergent  
 629 update (e.g. a quasi-Newton type update). The above result is also readily adapted  
 630 to an update that converges merely linearly, as long as its rate of convergence is faster  
 631 than that of  $\gamma_k$ . This opens the possibility of using SQP methods that rely only on  
 632 first-order information (see e.g. [4]).

633 **5. Numerical experiments.** In this section, we provide numerical illustrations  
 634 for our results. Our goal here is twofold:

- 635 i) to illustrate the identification of the optimal manifold by the proximity operator near a minimizer as provided by [Theorem 3.2](#);  
 636  
 637 ii) to demonstrate the applicability of [Algorithm 4.1](#) and observe the quadratic  
 638 rates predicted by [Theorem 4.4](#) on our running examples.

639 **5.1. Test problems.** We first consider the minimization of a pointwise maximum of smooth functions (1.2):

$$641 \quad \min_{x \in \mathbb{R}^n} \max_{i=1, \dots, m} (c_i(x)).$$

643 We take the celebrated `MaxQuad` instance, where  $n = 10$ ,  $m = 5$  and each  $c_i$  is  
 644 quadratic convex, making the whole function  $F$  convex [5, p. 153]. In this instance,  
 645 the optimal manifold is  $\mathcal{M}_F^{\max}$  with  $I = \{2, 3, 4, 5\}$ .

646 Second, we consider the minimization of the maximum eigenvalue of an affine  
 647 mapping (1.3):

$$648 \quad \min_{x \in \mathbb{R}^n} \lambda_{\max} \left( A_0 + \sum_{i=1}^n x_i A_i \right).$$

650 We take  $n = 25$  and we generate randomly  $n + 1$  symmetric matrices of size 50. In  
 651 this instance, the multiplicity of the maximum eigenvalue at the minimizer is  $r = 3$ .

652 **5.2. Numerical setup.** All the algorithms are implemented in Julia [2]; exper-  
 653 iments may be reproduced using the code available online<sup>3</sup>.

654 *Algorithm.* For the initialization of [Algorithm 4.1](#), we set  $\gamma_0$  as the smallest  $\gamma$  such  
 655 that  $\mathbf{prox}_{\gamma g}(c(x_0))$  has the most structure (*e.g.* if  $g = \max$ , we increase  $\gamma$  until the  
 656 output of the proximity operator sets all coordinates to the same value). We solve the  
 657 quadratic subproblem (4.1) providing the SQP step by the reduced system approach  
 658 presented in [5, p. 133]. Tangent vectors are expressed in an orthonormal basis of  
 659 the nullspace of the Jacobian of the constraints at the current iterate. At iterate  $x_k$ ,  
 660 a second-order correction step  $d^{\text{corr}}(x_k)$  is added to the SQP step  $d^{\text{SQP}}(x_k)$ . It is  
 661 obtained as  $d^{\text{corr}}(x_k) = \operatorname{argmin}_{d \in \mathbb{R}^n} \{\|h(x_k) + \operatorname{Jac}_h(x_k) d\|, \text{ s.t. } d \in \operatorname{Im} \operatorname{Jac}_h(x_k)^\top\}$ .  
 662 The full-step is thus  $x_k + d^{\text{SQP}}(x_k) + d^{\text{corr}}(x_k)$ .

663 *Baselines.* For the two nonsmooth problems, we compare with the nonsmooth  
 664 BFGS algorithm of [19] (nsBFGS) and the gradient sampling algorithm [7]. The  
 665 nsBFGS method is not covered by any theoretical guarantees; it is known to perform  
 666 relatively well in practice, often displaying a linear rate of convergence. In contrast,  
 667 the Gradient Sampling algorithm generates with probability one a sequence of iterates  
 668 for which all cluster points are Clarke stationary for  $F$  [7, Th. 3.1].<sup>4</sup>

669 *Oracles.* Traditional methods for nonsmooth optimization, and notably bundle  
 670 methods, require a first-order oracle:

$$671 \quad x \mapsto (F(x), v) \quad \text{where } v \in \partial F(x)$$

673 while Gradient Sampling and nsBFGS require additionally to know if  $F$  is differen-  
 674 tiable at point  $x$ . [Algorithm 4.1](#) requires rather different information oracles:

$$675 \quad x \mapsto F(x)$$

$$676 \quad x \mapsto \mathcal{M}^g \ni \mathbf{prox}_{\gamma g}(c(x))$$

$$677 \quad \mathcal{M}, x \mapsto h(x), \operatorname{Jac}_h(x), \nabla \tilde{F}(x), \nabla^2 L(x, \lambda).$$

<sup>3</sup><https://github.com/GillesBareilles/LocalCompositeNewton.jl>

<sup>4</sup>This holds when  $F$  is locally Lipschitz over  $\mathbb{R}^n$  and lower bounded, the algorithm iterates indefinitely and the sampling radius decreases to 0.

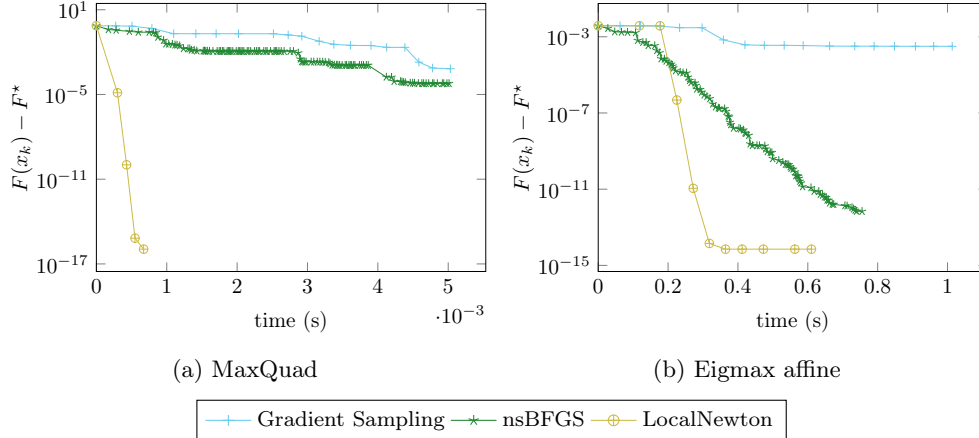


Fig. 4: Suboptimality vs time (s)

679 The second part of the oracle provides the candidate structure at point  $x$ . The last  
 680 part of the oracle, which requires a point *and a candidate structure*, provides the  
 681 second-order information of  $F$  required by the SQP step.

682 **5.3. Experiments.** Figure 4 reports the suboptimality of the considered meth-  
 683 ods in terms of CPU time and each marker corresponds to one iteration. All algorithms  
 684 are initialized at a point  $x_0$  obtained by running nsBFGS for several iterations.

685 Our algorithm compares favorably to nsBFGS and Gradient Sampling: it con-  
 686 verges in a handful of iterations and less time. Note that this happens even though  
 687 the iteration cost of our algorithm is higher than that of the other methods. Indeed,  
 688 the oracles of our method are more complex and a quadratic problem needs to be  
 689 solved, while the iteration cost of nsBFGS and Gradient Sampling is dominated by  
 690 the computation of function values and subgradients at each trials of the linesearch.

691 In terms of identification, our method finds the correct manifold at the first it-  
 692 eration for MaxQuad, and at the third iteration for Eigmax. From that point, the  
 693 iterates of Algorithm 4.1 reach machine precision in 3 iterations. This illustrates the  
 694 quadratic convergence, and supports the idea that, for nondifferentiable problems as  
 695 well, it is worth computing higher-order information to get fast local methods.

696 Figure 5 allows to observe the identification of the algorithm and the quality of  
 697 the bounds of Theorem 3.2. For each iterate  $x_k$  of Algorithm 4.1, we report the  
 698 current step  $\gamma_k$  along with the minimal and maximal steps  $\underline{\gamma}(x_k), \bar{\gamma}(x_k)$  such that  
 699  $\text{prox}_{\gamma_g}(c(x_k))$  belongs to the optimal manifold.<sup>5</sup> A first remark is that, as predicted  
 700 by Theorem 4.4, the pair  $x_k, \gamma_k$  satisfies the identification condition  $\gamma_k \in [L\|x_k - x^*, \Gamma]$   
 701 after a few iterations. We also observe that  $\bar{\gamma}(x_k)$  is near constant and that  $\underline{\gamma}(x_k)$   
 702 converges to zero linearly with  $\|x_k - x^*\|$ , as predicted by our result. Finally, we  
 703 note that even though the initial point is not structured and away from the minimizer  
 704 ( $\|x_0 - x^*\| \approx 10^{-2}$ ), the initialization of  $\gamma_0$  ensures a quick identification.

705 **6. Conclusions.** This paper studies the local structure of functions that write as  
 706 a composition of a nonsmooth function with a smooth mapping. When the proximity  
 707 operator of the nonsmooth function is explicitly available, we show that the structure

<sup>5</sup>To better illustrate the local behavior of our method, we also ran the algorithms with a high precision floating type. Details and corresponding experiments can be found in Appendix B.

708 of the minimizer can be detected. We further use this information to propose a  
 709 local Newton method to minimize the objective harnessing the detected structure.  
 710 This method is guaranteed to identify the structure of the minimizer and to converge  
 711 quadratically. We illustrate this behavior on two standard nonsmooth problems.

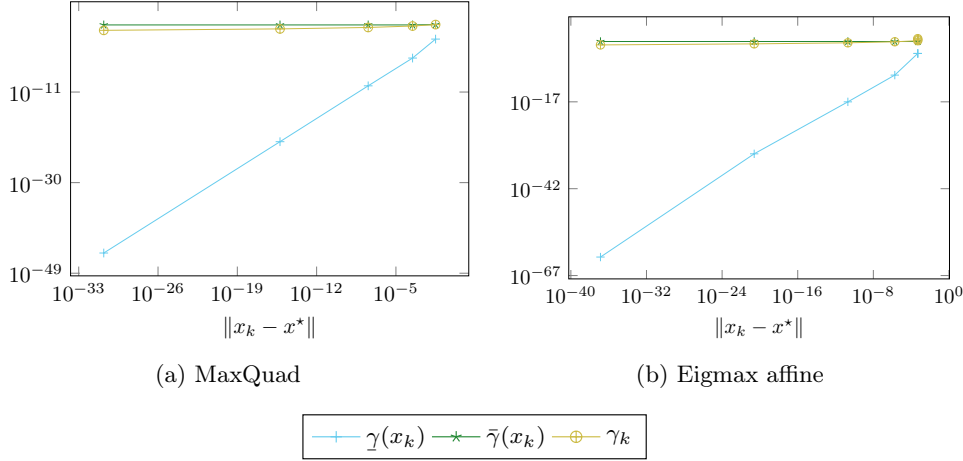


Fig. 5: Step size  $\gamma_k$  vs iteration

712 **Appendix A. The maximum and maximum eigenvalue satisfy the nor-**  
 713 **mal ascent and curve properties.** We show here that the maximum and the  
 714 maximum eigenvalue meet the normal ascent [Property 2.7](#) and curve properties [Prop-](#)  
 715 [erty 2.9](#). We begin with a lemma that simplifies verification of [Property 2.9](#).

716 **LEMMA A.1.** *Consider a function  $g$ , partly smooth at a point  $\bar{y}$  relative to a mani-*  
 717 *fold  $\mathcal{M}^g$ , and a smooth application  $e : \mathcal{N}_{\bar{y}} \times [0, T] \rightarrow \mathcal{M}^g$  defined for a neighborhood  $\mathcal{N}_{\bar{y}}$*   
 718 *of  $\bar{y}$  and  $T > 0$  such that  $e(y, 0) = \text{proj}_{\mathcal{M}^g}(y)$ ,  $\frac{d}{dt}e(y, t)|_{t=0} = -\text{grad}g(\text{proj}_{\mathcal{M}^g}(y))$ .*

719 *If  $D\left(t \mapsto \text{proj}_{N_{e(y,t)}\mathcal{M}^g}(\text{proj}_{\mathcal{M}}(y) - y)\right) = 0$  for all  $y \in \mathcal{N}_{\bar{y}}$ , then  $g$  satisfies Prop-*  
 720 *erty 2.9 at point  $\bar{y}$ .*

721 *Proof.* We denote  $\theta(y, t) = \text{proj}_{N_{e(y,t)}\mathcal{M}^g}(e(y, t) - y)$ . First,

$$722 \quad \frac{d}{dt}\theta(y, t)|_{t=0} = D\left(t \mapsto \text{proj}_{N_{e(y,t)}\mathcal{M}^g}(\text{proj}_{\mathcal{M}^g}(y) - y)\right)$$

$$723 \quad + \text{proj}_{N_{\text{proj}_{\mathcal{M}}(y)}\mathcal{M}^g}(D(t \mapsto (e(y, t) - y))(0)),$$

725 where the first term is null by assumption and the second is also null since it is the  
 726 normal projection of the tangent vector  $\text{grad}g(\text{proj}_{\mathcal{M}^g}(y))$ . Thus,  $\frac{d}{dt}\theta(y, t)|_{t=0} = 0$ .  
 727 Using this fact and smoothness of  $\theta$ , Taylor's theorem with Lagrange remainder yields,  
 728 for all  $y \in \mathcal{N}_{\bar{y}}$ , the existence of  $\bar{t} \in [0, T]$  such that, for all  $t \in [0, T]$ ,

$$729 \quad \theta(y, t) = \theta(y, 0) + \frac{t^2}{2} \frac{d^2}{dt^2}\theta(y, \bar{t}).$$

731 Therefore, for all  $y \in \mathcal{N}_{\bar{y}}$  and  $t \in [0, T]$ ,

$$732 \quad \|\theta(y, t)\| \leq \|\theta(y, 0)\| + \frac{t^2}{2} \sup_{\bar{t} \in [0, T]} \frac{d^2}{dt^2}\theta(y, \bar{t}) \leq \|\theta(y, 0)\| + t^2 \tilde{L},$$

733

734 where  $\tilde{L} = \sup_{y \in \mathcal{N}_{\bar{y}}} \sup_{\bar{t} \in [0, T]} \frac{d^2}{d\bar{t}^2} \theta(y, \bar{t})$ .  $\square$

735 We can now proceed with the proof of Lemma 2.10, divided into two parts cor-  
736 responding to the two cases of the result. The case  $g = \max$  comes easily, due to the  
737 polyhedral nature of the function.

738 LEMMA A.2. Consider  $g = \max$ , a point  $\bar{y} \in \mathbb{R}^m$  and the corresponding structure  
739 manifold  $\mathcal{M}_I^{\max}$  (of Example 2.5). Then Properties 2.7 and 2.9 hold at  $\bar{y}$ .

740 *Proof. Normal ascent* Take  $y \in \mathcal{M}_I^{\max}$  for some active indices  $I \subset \{1, \dots, m\}$ . A  
741 normal direction  $d \in N_y \mathcal{M}_I^{\max}$  is such that  $d_i = 0$  for  $i \notin I$  and  $\sum_{i \in I} d_i = 0$ . Thus  
742  $\max(y + td) = y_i + td_i$  with  $i = \operatorname{argmax}_i d_i$ , and  $D \max(y)[d] = \lim_{t \searrow 0} (\max(y + td) -$   
743  $\max(y))/t = d_i > 0$  for all  $d \neq 0$ .

744 *Curve assumption* Since the structure manifold of  $\max$  are affine subspaces, the  
745 normal spaces are equal at all points of the manifold. Therefore the derivative of the  
746 projection at a parametrized point is null and Lemma A.1 provides the result.  $\square$

747 The case  $g = \lambda_{\max}$  is not difficult *per se*, but requires a precise description of the  
748 geometry of the maximum eigenvalue function and its structure manifolds; we refer  
749 to [28, 25] for the derivation of these tools.

750 LEMMA A.3. Consider  $g = \lambda_{\max}$ , a point  $\bar{y} \in \mathbb{S}_m$  and the corresponding structure  
751 manifold  $\mathcal{M}_r^{\lambda_{\max}}$  (of Example 2.6). Then Properties 2.7 and 2.9 hold at  $\bar{y}$ .

752 *Proof. Normal ascent* Take  $y \in \mathcal{M}_r^{\lambda_{\max}}$ , let  $U \in \mathbb{R}^{m \times r}$  denote a basis of the first  
753 eigenspace of matrix  $y$  and  $d \in N_y \mathcal{M}_r^{\lambda_{\max}}$ . The normal space at  $y \in \mathcal{M}_r^{\lambda_{\max}}$  writes  
754 ([25, Th. 4.3, Cor. 4.8])

$$755 \quad N_y \mathcal{M}_r^{\lambda_{\max}} = \{U(y)ZU(y)^\top, Z \in \mathbb{S}_r, \operatorname{trace}(Z) = 0\}.$$

757 Therefore,  $d = UZU^\top$  for some  $Z \in \mathbb{S}_r$  such that  $\operatorname{trace}(Z) = 0$ . Let  $s = U(I/r +$   
758  $\alpha Z)U^\top$  where  $\alpha > 0$  is small enough so that  $s$  is positive definite. Since  $s$  has  
759 also unit trace, it is a subgradient of  $\lambda_{\max}$  at  $y$  [25, Th. 4.1]. Thus  $\lambda'_{\max}(y; d) =$   
760  $\sup_{v \in \partial \lambda_{\max}(y)} \langle v, d \rangle \geq \langle s, d \rangle = \langle I/r + \alpha Z, Z \rangle = \alpha \|Z\|^2$ . Hence  $\lambda'_{\max}(y; d) > 0$  for any  
761  $d \in N_y \mathcal{M}_r^{\lambda_{\max}} \setminus \{0\}$ .

762 *Curve assumption* Let  $\bar{y} \in \mathcal{M}_r^{\lambda_{\max}}$ . For any  $y \in \mathbb{S}_m$ , we denote by  $P(y)$  the  
763 orthogonal projection on the eigenspace corresponding to the  $r$  largest eigenvalues  
764 of  $y$  (counting multiplicities). This operator is smooth. We can define a mapping  
765  $U : \mathbb{S}_m \rightarrow \mathbb{R}^{m \times r}$  such that:  $U(y)^\top U(y) = I_r$ ,  $P(y) = U(y)U(y)^\top$ ,  $U$  is smooth near  
766 our reference point  $\bar{y}$  and its derivative at  $\bar{y}$  satisfies  $D U(\bar{y})^\top U(\bar{y}) = 0$ . The mapping  
767  $U$  defines a smooth orthonormal basis of the eigenspace corresponding to the  $r$  largest  
768 eigenvalues [28, p. 557]. Finally, for a point  $y' \in \mathcal{M}_r^{\lambda_{\max}}$ , the projection of  $d \in \mathbb{S}_m$  on  
769  $N_{y'} \mathcal{M}_r^{\lambda_{\max}}$  writes

$$770 \quad \operatorname{proj}_{N_{y'} \mathcal{M}_r^{\lambda_{\max}}}(d) = U(y') \left\{ U(y')^\top d U(y') - \frac{1}{r} \operatorname{trace}(U(y')^\top d U(y')) I_r \right\} U(y')^\top.$$

771

772 Now, fix  $y$  near  $\bar{y}$ , consider the eigenbasis  $U$  with reference point  $e(y, 0) =$   
773  $\operatorname{proj}_{\mathcal{M}_r^{\lambda_{\max}}}(y)$ . Following Lemma A.1, let  $\nu : t \mapsto \operatorname{proj}_{N_{e(y,t)} \mathcal{M}_r^{\lambda_{\max}}}(d)$  with  $d =$   
774  $\operatorname{proj}_{\mathcal{M}_r^{\lambda_{\max}}}(y) - y$ . We can now give an explicit expression of  $\nu(t)$  and show that

775  $\frac{d}{dt}\nu(0)$  is null. Denoting  $U(t) = U(e(y, t))$ , we have

$$776 \quad \nu(t) = U(t) \underbrace{\left\{ U(t)^\top dU(t) - \frac{1}{r} \text{trace}(U(t)^\top dU(t)) I_r \right\}}_{\triangleq \chi(t)} U(t)^\top.$$

777

778 First, as  $d$  is a normal vector to  $\mathcal{M}_r^{\lambda_{\max}}$  at point  $\text{proj}_{\mathcal{M}_r^{\lambda_{\max}}}(y)$ , there exists  $Z \in \mathbb{S}_r$   
779 such that  $d = U(0)ZU(0)^\top$ . Using that  $DU(0)^\top U(0) = 0$  yields

$$780 \quad DU(0)^\top dU(0) = DU(0)^\top U(0)ZU(0)^\top U(0) = 0.$$

782 Then, one readily checks that  $U(0)D\chi(0)U(0) = 0$ .

783 We turn to the term  $DU(0)\chi(0)U(0)^\top$ . A quick computation from the eigen  
784 decomposition of  $y$  shows that  $d$  writes  $U(0)ZU(0)^\top$ , where  $Z$  is actually diagonal.  
785 Therefore,  $\chi(0) = Z - (1/r)\text{trace}(Z)I_r$  is a diagonal matrix, so that

$$786 \quad DU(0)\chi(0)U(0)^\top = \sum_{i=1}^r \chi(0)_{ii} DU_i(0)U_i(0)^\top.$$

787

788 Following [28], the differential of  $t \mapsto U(e(y, t))$  at  $t = 0$  writes

$$789 \quad DU_i(0) = \sum_{k=r+1}^m \frac{1}{\lambda_1 - \lambda_k} U_k(0)U_k(0)^\top \eta U_i(0),$$

790

791 with  $\eta = \text{grad } \lambda_{\max}(\text{proj}_{\mathcal{M}_r^{\lambda_{\max}}}(y))$ . Using that  $\lambda_{\max}(y) = (1/r) \sum_{i=1}^r U_i(y)^\top y U_i(y)$ ,  
792 we compute the Riemannian gradient:  $\text{grad } \lambda_{\max}(y) = (1/r) \sum_{i=1}^r U_i(y)^\top U_i(y)$  (see [6,  
793 Sec. 7.7]). By orthogonality of the smooth basis of eigenvectors, the terms  $U_k(0)^\top U_i(0)$   
794 vanish for all  $i \in \{1, \dots, r\}$  and  $k \in \{r+1, \dots, m\}$ . We get that  $DU(0)\chi(0)U(0)^\top = 0$ ,  
795 and thus that  $D\nu(0) = 0$ . Thus, [Lemma A.1](#) applies and yields the result.  $\square$

796 **Appendix B. Numerical experiments in high precision.** We report in [Fig-](#)  
797 [ure 6](#) the evolution of suboptimality versus computing time, for the same problems  
798 and algorithms as in [section 5](#), but with a high precision floating type. Indeed, the  
799 flexibility of the Julia language allows to use the same implementation with the high  
800 precision `BigFloat` type, which precision is  $1.73 \cdot 10^{-72}$ , or the usual `Float64` type,  
801 which precision is  $2.22 \cdot 10^{-16}$ .

802 **Acknowledgments.** This work is funded by the ANR JCJC project STROLL  
803 (ANR-19-CE23-0008) and MIAI@Grenoble Alpes (ANR-19-P3IA-0003).

## 804 REFERENCES

- 805 [1] G. BAREILLES, F. IUTZELER, AND J. MALICK, *Newton acceleration on manifolds identified by*  
806 *proximal-gradient methods*, arXiv preprint arXiv:2012.12936, (2020).  
807 [2] J. BEZANSON, A. EDELMAN, S. KARPINSKI, AND V. B. SHAH, *Julia: A fresh approach to*  
808 *numerical computing*, SIAM review, 59 (2017), pp. 65–98.  
809 [3] J. BOLTE, Z. CHEN, AND E. PAUWELS, *The multiproximal linearization method for convex*  
810 *composite problems*, Mathematical Programming, 182 (2020), pp. 1–36, [https://doi.org/](https://doi.org/10.1007/s10107-019-01382-3)  
811 [10.1007/s10107-019-01382-3](https://doi.org/10.1007/s10107-019-01382-3).  
812 [4] J. BOLTE AND E. PAUWELS, *Majorization-Minimization Procedures and Convergence of SQP*  
813 *Methods for Semi-Algebraic and Tame Programs*, Mathematics of Operations Research, 41  
814 (2016), pp. 442–465, <https://doi.org/10.1287/moor.2015.0735>.



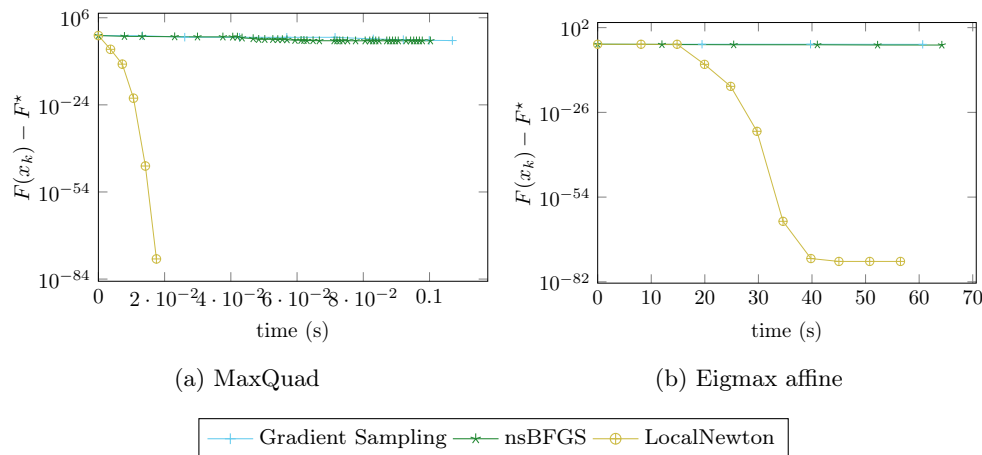


Fig. 6: Suboptimality vs time (s)

- 815 [5] J.-F. BONNANS, J. C. GILBERT, C. LEMARÉCHAL, AND C. A. SAGASTIZÁBAL, *Numerical*  
 816 *optimization: theoretical and practical aspects*, Springer Science & Business Media, 2006.
- 817 [6] N. BOUMAL, *An introduction to optimization on smooth manifolds*. To appear with Cambridge  
 818 University Press, Jan 2022, <http://www.nicolasboumal.net/book>.
- 819 [7] J. V. BURKE, F. E. CURTIS, A. S. LEWIS, M. L. OVERTON, AND L. E. SIMÕES, *Gradient*  
 820 *sampling methods for nonsmooth optimization*, in Numerical Nonsmooth Optimization,  
 821 Springer, 2020, pp. 201–225.
- 822 [8] A. DANILIDIS, W. HARE, AND J. MALICK, *Geometrical interpretation of the predictor-*  
 823 *corrector type algorithms in structured optimization problems*, Optimization, 55 (2006),  
 824 pp. 481–503.
- 825 [9] D. DRUSVYATSKIY, A. D. IOFFE, AND A. S. LEWIS, *Nonsmooth optimization using Taylor-like*  
 826 *models: Error bounds, convergence, and termination criteria*, Mathematical Programming,  
 827 185 (2021), pp. 357–383, <https://doi.org/10.1007/s10107-019-01432-w>.
- 828 [10] X. Y. HAN AND A. S. LEWIS, *Survey Descent: A Multipoint Generalization of Gradient*  
 829 *Descent for Nonsmooth Optimization*, (2021), p. 29.
- 830 [11] W. HARE AND A. S. LEWIS, *Identifying active constraints via partial smoothness and prox-*  
 831 *regularity*, Journal of Convex Analysis, 11 (2004), pp. 251–266.
- 832 [12] W. HARE AND C. SAGASTIZÁBAL, *Computing proximal points of nonconvex functions*, Math-  
 833 *ematical Programming*, 116 (2009), pp. 221–258.
- 834 [13] J.-B. HIRIART-URRUTY AND C. LEMARÉCHAL, *Convex Analysis and Minimization Algorithms*,  
 835 Springer Verlag, Heidelberg, 1993. Two volumes.
- 836 [14] C.-P. LEE, *Accelerating Inexact Successive Quadratic Approximation for Regularized Opti-*  
 837 *mization Through Manifold Identification*, arXiv:2012.02522 [math], (2021), <https://arxiv.org/abs/2012.02522>.
- 838 [15] J. M. LEE, *Introduction to Smooth Manifolds*, Graduate Texts in Mathematics, Springer-  
 839 *Verlag*, New York, 2003, <https://doi.org/10.1007/978-0-387-21752-9>.
- 840 [16] A. LEWIS AND T. TIAN, *Identifiability, the kl property in metric spaces, and subgradient*  
 841 *curves*, arXiv preprint arXiv:2205.02868, (2022).
- 842 [17] A. LEWIS AND C. WYLIE, *A simple Newton method for local nonsmooth optimization*,  
 843 arXiv:1907.11742 [cs, math], (2019), <https://arxiv.org/abs/1907.11742>.
- 844 [18] A. S. LEWIS, *Active sets, nonsmoothness, and sensitivity*, SIAM Journal on Optimization, 13  
 845 (2002), pp. 702–725.
- 846 [19] A. S. LEWIS AND M. L. OVERTON, *Nonsmooth optimization via quasi-Newton meth-*  
 847 *ods*, Mathematical Programming, 141 (2013), pp. 135–163, [https://doi.org/10.1007/](https://doi.org/10.1007/s10107-012-0514-2)  
 848 [s10107-012-0514-2](https://doi.org/10.1007/s10107-012-0514-2).
- 849 [20] A. S. LEWIS AND S. J. WRIGHT, *A proximal method for composite minimization*, Mathemat-  
 850 *ical Programming*, 158 (2016), pp. 501–546.
- 851 [21] R. MIFFLIN AND C. SAGASTIZÁBAL, *A  $\mathcal{VU}$ -algorithm for convex minimization*, Mathematical  
 852 *programming*, 104 (2005), pp. 583–608.
- 853 [22] S. A. MILLER AND J. MALICK, *Newton methods for nonsmooth convex minimization: con-*  
 854

- 855            *nections among-lagrangian, riemannian newton and sqp methods*, Mathematical program-  
856            ming, 104 (2005), pp. 609–633.
- 857 [23] J. NOCEDAL AND S. WRIGHT, *Numerical optimization*, Springer Science & Business Media,  
858            2006.
- 859 [24] D. NOLL AND P. APKARIAN, *Spectral bundle methods for non-convex maximum eigenvalue*  
860            *functions: second-order methods*, Mathematical Programming, 104 (2005), pp. 729–747.
- 861 [25] F. OUSTRY, *The U-Lagrangian of the Maximum Eigenvalue Function*, SIAM Journal on Opti-  
862            mization, 9 (1999), pp. 526–549, <https://doi.org/10.1137/S1052623496311776>.
- 863 [26] R. T. ROCKAFELLAR AND R. J.-B. WETS, *Variational analysis*, vol. 317, Springer Science &  
864            Business Media, 2009.
- 865 [27] A. SHAPIRO, *On a Class of Nonsmooth Composite Functions*, Mathematics of Operations  
866            Research, 28 (2003), pp. 677–692, <https://doi.org/10.1287/moor.28.4.677.20512>.
- 867 [28] A. SHAPIRO AND M. K. H. FAN, *On Eigenvalue Optimization*, SIAM Journal on Optimization,  
868            5 (1995), pp. 552–569, <https://doi.org/10.1137/0805028>.
- 869 [29] R. S. WOMERSLEY AND R. FLETCHER, *An algorithm for composite nonsmooth optimization*  
870            *problems*, Journal of Optimization Theory and Applications, 48 (1986), pp. 493–523, <https://doi.org/10.1007/BF00940574>.  
871