



Author Proof

## Proximal Methods Avoid Active Strict Saddles of Weakly Convex Functions

Damek Davis<sup>1</sup> · Dmitriy Drusvyatskiy<sup>2</sup>

Received: 22 January 2020 / Revised: 16 February 2021 / Accepted: 8 April 2021  
© SFoCM 2021

### Abstract

We introduce a geometrically transparent strict saddle property for nonsmooth functions. This property guarantees that simple proximal algorithms on weakly convex problems converge only to local minimizers, when randomly initialized. We argue that the strict saddle property may be a realistic assumption in applications, since it provably holds for generic semi-algebraic optimization problems.

**Keywords** Strict saddle · Proximal gradient · Proximal point · Center stable manifold theorem · Semi-algebraic

**Mathematics Subject Classification** 65K05 · 65K10 · 90C30

### 1 Introduction

Nonconvex optimization techniques are increasingly playing a major role in modern signal processing, high-dimensional statistics, and machine learning. A driving theme, fully supported by empirical evidence, is that simple algorithms often work well in highly nonconvex and even nonsmooth settings. Gradient descent, for example, often

---

Communicated by Michael Overton.

---

D. Drusvyatskiy: Research of Drusvyatskiy was supported by the NSF DMS 1651851 and CCF 1740551 awards.

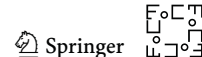
---

✉ Dmitriy Drusvyatskiy  
ddrusv@uw.edu  
<http://www.sites.math.washington.edu/ddrusv/>

Damek Davis  
<http://www.people.orie.cornell.edu/dsd95/>

<sup>1</sup> School of Operations Research and Information Engineering, Cornell University, Ithaca, NY 14850, USA

<sup>2</sup> Department of Mathematics, University of Washington, Seattle, WA 98195, USA



15 finds points with small objective value, despite existence of many highly suboptimal  
 16 critical points. A growing body of literature provides one compelling explanation for  
 17 this phenomenon. Namely, typical smooth objective functions provably satisfy the  
 18 *strict saddle property*, meaning each critical point is either a local minimizer or has  
 19 a direction of strictly negative curvature (e.g., [6,28,29,61,62]). For such functions,  
 20 randomly initialized gradient-type methods provably converge to local minimizers,  
 21 escaping all strict saddle points [35,51]. Moreover, stochastically perturbed gradient  
 22 methods escape strict saddles efficiently, indeed, in polynomial time [22,27,33].

23 Smoothness of the objective plays a crucial role in the existing literature on sad-  
 24 dle avoidance. Some extensions to constrained optimization do exist. The papers  
 25 [15,27,63] investigate saddle point avoidance for the problem of minimizing a smooth  
 26 functions over a smooth manifold. The works [30,44,49] propose algorithms for min-  
 27 imizing a smooth objective over a closed convex set. At each step of these algorithms,  
 28 one must minimize a nonconvex quadratic over a certain convex set (an NP hard prob-  
 29 lem in general). The work [4] proposes a polynomial time first-order algorithm for  
 30 minimizing a smooth objective over linear inequality constraints.<sup>1</sup> At each step of this  
 31 algorithm, one identifies the “active linear constraints” and then attempts to find a  
 32 “second-order stationary point” of the loss in the restricted subspace.

33 Though impressive in scope, existing work has yet to answer the following question:

34 Do simple algorithms on typical nonsmooth and nonconvex optimization prob-  
 35 lems converge only to local minimizers?

36 This question as stated is purposefully vague, since “simple algorithms” and “typical  
 37 optimization problems” can be interpreted in multiple ways. Let us try to formalize both  
 38 ideas. First, if one believes that gradient descent is a canonical first-order method for  
 39 smooth minimization, it is natural to focus on three concrete algorithms for nonsmooth  
 40 and nonconvex problems: the proximal point [42,43,46,55], proximal gradient [5,48],  
 41 and proximal linear [9,20,21,40,47] methods. This is the path we follow in the current  
 42 work.

43 The latter issue, identifying a typical optimization problem, is more subtle. To moti-  
 44 vate our approach, let us revisit the following question: why is the strict saddle property  
 45 a reasonable assumption for smooth minimization? A first compelling reason is that  
 46 the property holds in practice for specific problems of interest. There is, however, a  
 47 more classical justification, one rooted in stability to perturbations. Namely, consider  
 48 the task of minimizing a smooth function  $f$  on  $\mathbb{R}^d$ . Then, for a full measure set of per-  
 49 turbations  $v \in \mathbb{R}^d$ , the perturbed function  $x \mapsto f(x) - \langle v, x \rangle$  is guaranteed to satisfy  
 50 the strict saddle property—a direct consequence of Sard’s theorem. Consequently, in  
 51 a precise mathematical sense, the strict saddle property holds *generically* in smooth  
 52 optimization. This justification does not suggest one can omit verification of the strict  
 53 saddle property in concrete circumstances, but it does suggest that the strict saddle  
 54 property is widely prevalent.

55 Seeking to identify a similarly reasonable class of nonsmooth objectives on which  
 56 simple algorithms converge to local minimizers, the current paper accomplishes the  
 57 following.

<sup>1</sup> This work appeared concurrently with our manuscript.

58 We formulate natural geometric conditions, guaranteeing the proximal point,  
 59 proximal gradient, and proximal linear algorithms escape all saddle points.  
 60 Moreover, the proposed conditions are generic: they hold for almost all linear  
 61 perturbations of weakly convex and semi-algebraic problems.

62 The class of optimization problems we consider is broad. A function  $f$  is called  $\rho$ -  
 63 weakly convex if the assignment  $x \mapsto f(x) + \frac{\rho}{2}\|x\|^2$  is convex for some  $\rho > 0$ .<sup>2</sup>  
 64 Common examples include pointwise maxima of smooth functions and all compo-  
 65 sitions of Lipschitz convex functions with smooth maps. For detailed contemporary  
 66 examples, we refer the reader to [13,16,17,23,32]. The genericity guarantee applies  
 67 to semi-algebraic functions,<sup>3</sup> and more broadly, to those that are definable in an o-  
 68 minimal structure—a virtually exhaustive function class in applications.

### 69 1.1 The Role of Curvature

70 To motivate our core geometric conditions, we revisit the role that curvature plays in  
 71 saddle-point avoidance. Setting the stage for the rest of the paper, consider the task  
 72 of minimizing a weakly convex function  $f$  on  $\mathbb{R}^d$ . First-order optimality conditions  
 73 show that any local minimizer  $\bar{x}$  of  $f$  satisfies the *criticality condition*:

$$74 \quad df(\bar{x})(v) \geq 0 \quad \text{for all } v \in \mathbb{R}^d,$$

75 where  $df(\bar{x})(v)$  denotes the directional derivative of  $f$  at  $\bar{x}$  in direction  $v$  (see Defi-  
 76 nition 2.1). Conversely, sufficient conditions for local optimality at a critical point  $\bar{x}$   
 77 require a close look at the second-order variations of  $f$  along particular directions,  
 78 namely those where the directional derivative is zero. Mirroring the smooth setting,  
 79 one may naively call a critical point  $\bar{x}$  a strict saddle if there exists a direction  $v$  such  
 80 that  $df(\bar{x})(v) = 0$  and  $f$  decreases quadratically along  $v$ . This definition, however, is  
 81 insufficient for saddle avoidance: simple examples show that typical algorithms can  
 82 converge to such saddle points from a positive measure of initial conditions.

83 Negative curvature alone does not guarantee escape from saddle points.

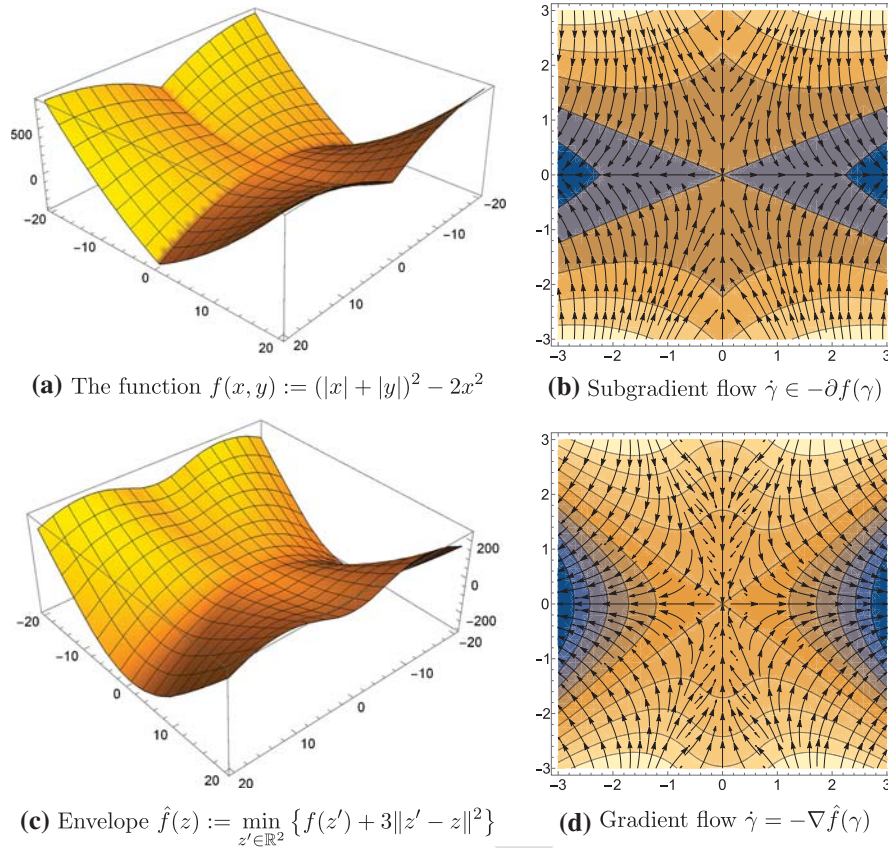
84 To illustrate what can go wrong, consider the example

$$85 \quad \min_{(x,y) \in \mathbb{R}^2} f(x,y) = (|x| + |y|)^2 - 2x^2, \quad (1.1)$$

86 the graph of which is shown in Fig. 1a. First, observe that the origin meets the con-  
 87 ditions of the candidate “strict saddle” definition. Indeed,  $f$  is differentiable at the  
 88 origin and the origin is a critical point. Moreover,  $f$  decreases quadratically along all  
 89 directions making a small angle with the  $x$ -axis. Next, we turn to algorithm dynam-  
 90 ics. Figure 1b depicts the subgradient flow  $-\dot{\gamma}(t) \in \partial f(\gamma(t))$ . From the picture we

<sup>2</sup> Weakly convex functions also go by other names such as lower- $C^2$ , uniformly prox-regularity, paraconvex, and semiconvex. We refer the reader to the seminal works on the topic [2,50,53,56,58].

<sup>3</sup> A function is called semi-algebraic if its graph decomposes into a finite union of sets, each defined by finitely many polynomial inequalities.



**Fig. 1** The function  $f$  in (1.1), its Moreau envelope, and their subgradient flows. **a** The function  $f(x, y) := (|x| + |y|)^2 - 2x^2$ . **b** Subgradient flow  $\dot{\gamma} \in -\partial f(\gamma)$ . **c** Envelope  $\hat{f}(z) := \min_{z' \in \mathbb{R}^2} \{f(z') + 3\|z' - z\|^2\}$ . **d** Gradient flow  $\dot{\gamma} = -\nabla \hat{f}(\gamma)$

91 find a positive measure cone, surrounding the  $y$ -axis and consisting of origin-attracted  
 92 initial conditions. Moreover, we show in Appendix B that a typical algorithm—the  
 93 proximal point method—initialized anywhere within this cone also converges to the  
 94 origin, illustrating the inadequacy of the definition. While this argument shows that  
 95 negative curvature is insufficient for nonsmooth optimization, it can be made even  
 96 more definitive by smoothing the problem at hand. Namely, an alternative view of  
 97 the proximal point method recognizes that the dynamics of the algorithm coincide  
 98 with the dynamics of gradient descent on a  $C^1$  smooth approximation of  $f$ , called  
 99 the *Moreau envelope* (see Sect. 2.3). The smooth envelope, whose graph and gradient  
 100 flow are shown in Fig. 1c, d, has the same cone of directions of second-order negative  
 101 curvature as  $f$ , but despite its smoothness and negative curvature, gradient descent  
 102 cannot escape the origin. The problem persists under a variety of different choices of  
 103 the step-size. Note that there is no contradiction with the saddle avoidance property  
 104 of gradient descent on smooth functions, since the envelope is not  $C^2$ , but merely  $C^1$ -

105 smooth around the origin. Although this example appears damning at first, it is highly  
 106 unstable, since small linear tilts of the function do not exhibit the same pathological  
 107 behavior around critical points—a direct consequence of the forthcoming results.

## 108 1.2 The Role of the Active Manifold

109 We have seen that negative curvature alone is insufficient for saddle avoidance. We  
 110 argue here that in addition to negative curvature, one must make a structural assumption  
 111 on the way nonsmoothness manifests. To illustrate and contrast with example (1.1)  
 112 above, consider:

$$113 \min_{(x,y) \in \mathbb{R}^2} g(x, y) = |x| - y^2. \quad (1.2)$$

114 The graph of  $g$  is shown in Fig. 2a, while its subgradient flow appears in Fig. 2b.  
 115 Looking at the figure, we see that the subgradient flow of  $g$  sharply contrasts with  
 116 that of the pathological example (1.1). Indeed, while both functions have directions  
 117 of negative curvature, the set of origin-attracted initial conditions of the flow  $-\partial g$  is  
 118 simply the  $x$ -axis—a measure zero set. This favorable behavior of  $g$  arises because its  
 119 nonsmoothness manifests in a structured way: its unique critical point  $\bar{z}$  (the origin)  
 120 lies on a smooth manifold  $\mathcal{M}$  (the  $y$ -axis). The function  $g$  then varies smoothly along  
 121  $\mathcal{M}$  and sharply normal to  $\mathcal{M}$  meaning:

$$122 \inf\{\|v\| : v \in \partial g(z), z \in U \setminus \mathcal{M}\} > 0,$$

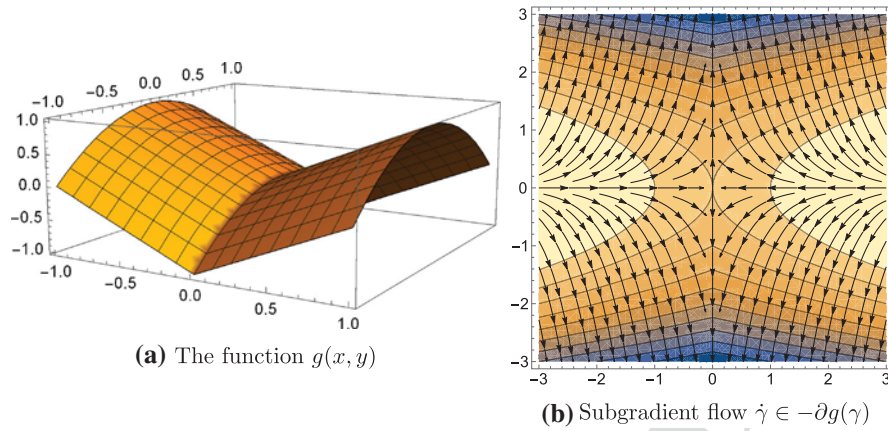
123 where  $U$  is some neighborhood of  $\bar{z}$ . Such “active manifolds” have classical roots in  
 124 optimization and serve as a far reaching generalization of “active sets” in nonlinear  
 125 programming. Important references include both the original works [1, 10–12, 24–26]  
 126 and the more recent work on identifiable surfaces [64],  $UV$ -decomposition [38], partial  
 127 smoothness [39], and cone reducibility [8]. Here, we most closely follow the  
 128 framework developed in [19].

## 129 1.3 Escape from Saddles by the Center Stable Manifold Theorem

130 Formalizing the favorable behavior of example (1.2), we will call a critical point  $\bar{x}$  of a  
 131 function  $g$  a *strict saddle* whenever (i)  $g$  admits an active manifold containing  $\bar{x}$ , and (ii)  
 132 the function  $g$  decreases quadratically along some direction  $v$  satisfying  $dg(\bar{x})(v) = 0$ .  
 133 A function  $g$  is said to have the *strict saddle property* if each of its critical points is  
 134 either a local minimizer or a strict saddle.<sup>4</sup> Though it may seem that this definition is  
 135 stringent at first, the strict saddle property is in a precise mathematical sense generic.  
 136 Namely, it follows from [18] that given any semi-algebraic weakly convex function  
 137  $g$ , the perturbed function  $g_v(x) = g(x) - \langle v, x \rangle$  has the strict saddle property for  
 138 almost all  $v \in \mathbb{R}^d$ .<sup>5</sup> In particular, almost all linear perturbations of the function  $f$  in

<sup>4</sup> Perhaps more appropriate would be the terms *active strict saddle* and the *active strict saddle property*. For brevity, we omit the word “active.”

<sup>5</sup> Weak convexity is not essential here, provided one modifies the definitions appropriately. Moreover, this guarantee holds more generally for functions definable in an o-minimal structure.



**Fig. 2** The function  $g(x, y) = |x| - y^2$  has an active manifold at the origin. **a** The function  $g(x, y)$ . **b** Subgradient flow  $\dot{\gamma} \in -\partial g(\gamma)$

**Table 1** The three algorithms with the update  $S(x) = \operatorname{argmin}_y f_x(y)$ ; we assume  $h$  is convex,  $r$  is weakly convex, and both  $g$  and  $F$  are smooth

Algorithm	Objective	Update function $f_x(y)$
Prox-point	$r(x)$	$r(y) + \frac{1}{2\mu} \ y - x\ ^2$
Prox-gradient	$g(x) + r(x)$	$g(x) + \langle \nabla g(x), y - x \rangle + r(y) + \frac{1}{2\mu} \ y - x\ ^2$
Prox-linear	$h(F(x)) + r(x)$	$h(F(x) + \nabla F(x)(y - x)) + r(y) + \frac{1}{2\mu} \ y - x\ ^2$

139 (1.1) do have the strict saddle property. That being said, it is important to note that  
 140 under more nuanced perturbations, the strict saddle property may fail. For example,  
 141 the classical NP-complete problem of checking copositivity of a matrix  $A \in \mathbb{R}^{d \times d}$   
 142 amounts to verifying if the origin  $\bar{x} = 0$  locally minimizes the quadratic  $x^T A x$  over  
 143 the nonnegative orthant  $\mathbb{R}_+^d$ . It is straightforward to see that this constrained problem  
 144 does not admit an active manifold at  $\bar{x}$  for any matrix  $A$ .

145 With the definition of a strict saddle at hand, we can now outline the main results  
 146 of the paper. As in the smooth setting, first explored in the seminal paper [35], our  
 147 arguments will be based on the center stable manifold theorem. Namely, we will  
 148 interpret the three simple minimization algorithms as fixed point iterations

149 
$$x_{k+1} = S(x_k) \quad \text{for some maps } S: \mathbb{R}^d \rightarrow \mathbb{R}^d.$$

150 Table 1 lists the maps  $S(\cdot)$  for the proximal point, proximal gradient, and proximal  
 151 linear algorithms. In each case, the fixed points of  $S(\cdot)$  are precisely the critical points  
 152 of the minimization problem.

153 To put our guarantees in context, it will be useful to recall the center stable manifold  
 154 theorem. To this end, suppose that the iteration map  $S(\cdot)$  is  $C^1$ -smooth on a neigh-  
 155 borhood of some fixed point  $\bar{x}$ . Then,  $\bar{x}$  is called an *unstable fixed point* of  $S$  if the  
 156 Jacobian  $\nabla S(\bar{x})$  has at least one eigenvalue whose magnitude is strictly greater than

157 one. The center stable-manifold theorem [60, Theorem III.7] guarantees the following:  
 158 if  $\bar{x}$  is an unstable fixed point of  $S$  and the Jacobian  $\nabla S(\bar{x})$  is invertible, then almost  
 159 all initializers  $x$  in a neighborhood  $U$  of  $\bar{x}$  generate iterates  $\{S^k(x)\}_{k \geq 1}$  that eventually  
 160 escape the neighborhood. More precisely, the theorem guarantees that the set of initial  
 161 conditions

$$162 \quad \left\{ x \in U : S^k(x) \in U \text{ for all } k \geq 1 \right\}$$

163 has zero Lebesgue measure. All that is needed to globalize this guarantee is to ensure  
 164 that the preimage  $S^{-1}(V)$  of any measure zero set  $V$  is itself measure zero. Then, for  
 165 almost all initial conditions  $x \in \mathbb{R}^d$ , the limit  $\lim_{k \rightarrow \infty} S^k(x)$ , when it exists, is not  
 166 an unstable fixed point of  $S$ . A straightforward way to ensure that the inverse  $S^{-1}$   
 167 respects null sets is by introducing the relaxation map:

$$168 \quad T(x) := (1 - \alpha)x + \alpha S(x). \quad (1.3)$$

169 Both  $T$  and  $S$  have the same fixed points, and any fixed point  $\bar{x}$  at which  $\nabla S(\bar{x})$  has a  
 170 real eigenvalue strictly greater than one is an unstable fixed point of  $T$ . Moreover, if  
 171 the map  $S$  is Lipschitz, then the inverse  $T^{-1}$  preserves null-sets for sufficiently small  
 172  $\alpha \in (0, 1)$ .

#### 173 1.4 The Main Results

174 We can now summarize our main results:

175 We show that around each strict saddle of the problem, each of the iterations  
 176 maps  $S(\cdot)$  in Table 1 is  $C^1$  smooth. Moreover, if  $\bar{x}$  is a strict saddle, then the  
 177 Jacobian  $\nabla S(\bar{x})$  has a real eigenvalue strictly greater than one.

178 From this result, the center stable manifold theorem guarantees that iteration (1.3)  
 179 locally escapes strict saddles. Seeking to globalize the guarantees, we compute the  
 180 global Lipschitz constants for the proximal point and proximal gradient methods.  
 181 We deduce that, when randomly initialized, the relaxed iterations (1.3) for both the  
 182 proximal point and proximal gradient methods converge to local minimizers of weakly  
 183 convex functions, provided they have the strict saddle property. On the other hand,  
 184 without placing further restrictions on the problem data, we are unable to compute  
 185 the global Lipschitz constant of the map  $S(\cdot)$  corresponding to the proximal linear  
 186 algorithm. We leave it as an intriguing open question to determine Lipschitz properties  
 187 of the proximal linear update.

188 The outlined results may seem surprising at first: the optimization problem is non-  
 189 smooth and yet we prove the iteration maps  $S(\cdot)$  are  $C^1$ -smooth around any strict  
 190 saddle. The reason is transparent and derives from the interplay between the active  
 191 manifold and weak convexity. Take the proximal point method, for example. The very  
 192 definition of the active manifold guarantees that the fixed point iteration  $S(\cdot)$  maps an  
 193 entire neighborhood  $\mathcal{X}$  around an strict saddle  $\bar{x}$  into the active manifold  $\mathcal{M}$ . Conse-  
 194 quently, for all  $x \in \mathcal{X}$ , the update  $S(x)$  can be realized as a minimizer of a smooth

195 function over the active manifold:

196 
$$S(x) = \operatorname{argmin}_{y \in \mathcal{M}} f(y) + \frac{1}{2\mu} \|y - x\|^2. \quad (1.4)$$

197 Weak convexity, in turn, ensures that  $S(\bar{x})$  satisfies a quadratic growth condition for  
 198 the problem (1.4), which by classical perturbation theory guarantees that  $S(\cdot)$  is  $C^1$ -  
 199 smooth on a neighborhood of  $\bar{x}$ . It only remains to argue that the negative curvature of  
 200 the objective function at  $\bar{x}$  implies that the Jacobian  $\nabla S(\bar{x})$  has at least one real eigen-  
 201 value greater than one. Though this computation is straightforward for the proximal  
 202 point method, it becomes more interesting (and surprising) for the proximal gradient  
 203 and proximal linear algorithms.

204 **Roadmap** The outline of the paper is as follows. Section 2 is a self-contained pre-  
 205 sentation of the necessary preliminaries for formalizing the ideas of the introduction.  
 206 Then, in Sects. 3, 4, and 5 we directly analyze the iteration maps for the proximal  
 207 point, proximal gradient, and proximal linear algorithms. Section 6 establishes iterate  
 208 convergence of the relaxed schemes (1.3) under the Kurdyka–Łojasiewicz property.

209 **2 Preliminaries**

210 Throughout, we follow standard notation in convex and variational analysis, as set  
 211 out, for example, in the monographs [14,45,54,57]. We consider a Euclidean space  
 212  $\mathbb{R}^d$  endowed with an inner product  $\langle \cdot, \cdot \rangle$  and the induced norm  $\|x\| = \sqrt{\langle x, x \rangle}$ . The  
 213 unit sphere in  $\mathbb{R}^d$  will be denoted by  $\mathbb{S}^{d-1}$ . For any function  $f: \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$ , the  
 214 domain and epigraph are the sets

215 
$$\operatorname{dom} f = \{x \in \mathbb{R}^d : f(x) < \infty\}, \quad \operatorname{epi} f = \{(x, r) \in \mathbb{R}^d \times \mathbb{R} : r \geq f(x)\},$$

216 respectively. The function  $f$  is called *closed* if  $\operatorname{epi} f$  is a closed set. For any set  
 217  $\mathcal{M} \subset \mathbb{R}^d$ , the indicator function  $\delta_{\mathcal{M}}$  evaluates to zero on  $\mathcal{M}$  and to  $+\infty$  off it.  
 218 For any function  $f: \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$  and a set  $\mathcal{M} \subset \mathbb{R}^d$ , we define the restriction  
 219  $f_{\mathcal{M}} := f + \delta_{\mathcal{M}}$ . Throughout, the symbol  $o(r)$  will denote any univariate function  
 220 satisfying  $o(r)/r \rightarrow 0$  as  $r \searrow 0$ .

221 Consider a differentiable mapping  $F(x) = (F_1(x), \dots, F_m(x))$  from  $\mathbb{R}^d$  to  $\mathbb{R}^m$ .  
 222 Throughout, the symbol  $\nabla F(x) \in \mathbb{R}^{m \times d}$  will denote the Jacobian matrix, whose  $ij$ 'th  
 223 entry is given by  $\frac{d}{dx_j} F_i(x)$ . Thus, row  $i$  of  $\nabla F(x)$  is the gradient of the coordinate  
 224 function  $F_i(x)$ . In the particular case  $m = 1$ , we will treat  $\nabla F(x)$  either as a column or  
 225 as a row vector, depending on context. For a  $C^2$ -smooth function  $g: \mathbb{R}^d \times \mathbb{R}^n \rightarrow \mathbb{R}^m$ ,  
 226 we partition the Hessian as follows:

227 
$$\nabla^2 g(x, y) = \begin{bmatrix} \nabla_{xx} g(x, y) & \nabla_{xy} g(x, y) \\ \nabla_{yx} g(x, y) & \nabla_{yy} g(x, y) \end{bmatrix}$$

Author Proof



228 **2.1 Subdifferentials and Subderivatives**

229 The following definition records the standard first- and second-order differential con-  
 230 structions, which we will use in the paper. After the definition, we will comment on  
 231 the role of each construction. For further details we refer the reader to [57, Definitions  
 232 8.1, 8.3, 13.59].

233 **Definition 2.1** (*Subdifferential and subderivatives*) Consider a function  $f: \mathbb{R}^d \rightarrow$   
 234  $\mathbb{R} \cup \{\infty\}$  and a point  $\bar{x}$  with  $f(\bar{x})$  finite. Then, the *subdifferential* of  $f$  at  $\bar{x}$ , denoted  
 235  $\partial f(\bar{x})$ , consists of all vectors  $v$  satisfying

236 
$$f(x) \geq f(\bar{x}) + \langle v, x - \bar{x} \rangle + o(\|x - \bar{x}\|) \quad \text{as } x \rightarrow \bar{x}.$$

237 The *subderivative* of  $f$  at  $\bar{x}$  in direction  $\bar{u} \in \mathbb{R}^d$  is

238 
$$df(\bar{x})(\bar{u}) := \liminf_{\substack{t \searrow 0 \\ u \rightarrow \bar{u}}} \frac{f(\bar{x} + tu) - f(\bar{x})}{t}.$$

239 The *critical cone* of  $f$  at  $\bar{x}$  for  $\bar{v} \in \mathbb{R}^d$  is

240 
$$C_f(\bar{x}, \bar{v}) := \{u \in \mathbb{R}^d : \langle \bar{v}, u \rangle = df(\bar{x})(u)\}.$$

241 The *parabolic subderivative* of  $f$  at  $\bar{x}$  for  $\bar{u} \in \text{dom } df(\bar{x})$  with respect to  $\bar{w}$  is

242 
$$d^2 f(\bar{x})(\bar{u}|\bar{w}) = \liminf_{\substack{t \searrow 0 \\ w \rightarrow \bar{w}}} \frac{f(\bar{x} + t\bar{u} + \frac{1}{2}t^2w) - f(\bar{x}) - df(\bar{x})(\bar{u})}{\frac{1}{2}t^2}.$$

243 We now comment on these definitions, in order. First, a vector  $v$  lies in the subdif-  
 244 ferential  $\partial f(\bar{x})$  precisely when the affine function  $x \mapsto f(\bar{x}) + \langle v, x - \bar{x} \rangle$  minorizes  
 245  $f$  up to first order near  $\bar{x}$ . The definition reduces to familiar objects in classical  
 246 circumstances. For example, differentiability of  $f$  at  $\bar{x}$  implies the set  $\partial f(\bar{x})$  is a  
 247 singleton, containing only the gradient  $\nabla f(\bar{x})$ . Convexity of  $f$  too entails a simplifi-  
 248 cation, wherein  $\partial f(\bar{x})$  reduces to the subdifferential of convex analysis.

249 While the subdifferential encodes the set of approximate affine minorants, the sub-  
 250 derivative measures the maximal instantaneous rate of decrease of  $f$  in direction  $\bar{u}$ .  
 251 Like the subdifferential, the subderivative reduces to familiar objects in classical cir-  
 252 cumstances. For example, if  $f$  is locally Lipschitz at  $\bar{x}$ , then one may set  $u = \bar{u}$  in  
 253 its defining expression. Simplifying further, if  $f$  is differentiable at  $\bar{x}$ , we recover the  
 254 directional derivative expression  $df(\bar{x})(\bar{u}) = \langle \nabla f(\bar{x}), \bar{u} \rangle$ . Finally, if  $f$  is convex, then  
 255 the subderivative reduces to the support function of the subdifferential

256 
$$df(\bar{x})(\bar{u}) = \sup\{\langle \bar{u}, v \rangle : v \in \partial f(\bar{x})\},$$

257 highlighting the dual roles of the subdifferential and subderivative constructions.

258 For smooth losses, necessary optimality conditions entail vanishing gradients, while  
 259 sufficient optimality conditions follow from second-order growth properties of  $f$ . Similar  
 260 characterizations persist in the nonsmooth setting. In particular, the subderivative  
 261 and the subdifferential feature in first-order necessary optimality conditions, where  
 262 the (dual) criticality condition  $0 \in \partial f(\bar{x})$  is equivalent to the (primal) nonnegativity  
 263 condition

$$264 \quad df(\bar{x})(u) \geq 0 \quad \text{for all } u \in \mathbb{R}^d. \quad (2.1)$$

265 A point  $\bar{x}$  satisfying these first-order necessary conditions (2.1) is thus called *critical*  
 266 for  $f$ . Sufficient optimality conditions, on the other hand, make use of second-order  
 267 variations of  $f$ . Namely, suppose that a point  $\bar{x}$  is critical for  $f$  and consider a direction  
 268  $\bar{u} \in \mathbb{R}^d$ . There are two possibilities to consider. On the one hand, if  $df(\bar{x})(\bar{u}) > 0$ ,  
 269 then  $f$  must locally increase in direction  $\bar{u}$ . On the other hand, if  $df(\bar{x})(\bar{u}) = 0$ , then  
 270 we must examine second-order variations of  $f$  to determine local optimality. Such  
 271 directions of ambiguity for the subderivative make up the critical cone  $C_f(\bar{x}, 0)$ . For  
 272 these directions, we must look to the parabolic derivative  $d^2 f(\bar{x})(\bar{u}|\bar{w})$ , a measurement  
 273 of the second-order variation of  $f$  along a parabolic arc with tangent direction  $\bar{u}$  and  
 274 second-order variation  $\bar{w}$ . This construction too simplifies when  $f$  is  $C^2$  smooth at  $\bar{x}$ ,  
 275 reducing to the familiar second-order variation:

$$276 \quad d^2 f(\bar{x})(\bar{u}|\bar{w}) = \langle \nabla^2 f(\bar{x})\bar{u}, \bar{u} \rangle.$$

277 This relation suggests second-order optimality conditions for nonsmooth problems.  
 278 Although we will not appeal to such conditions directly in this work, we record them  
 279 here for completeness. If  $\bar{x}$  is a local minimizer of  $f$ , then  $df(\bar{x})(u) \geq 0$  for all  $u \in \mathbb{R}^n$ ,  
 280 and moreover  $\inf_{w \in \mathbb{R}^n} d^2 f(\bar{x})(u|w) \geq 0$  for any nonzero  $u \in C_f(\bar{x}, 0)$ . Complementing  
 281 this necessary condition, a large class of functions, those that are *parabolically*  
 282 *regular*, may also be endowed with a sufficient optimality condition. Namely, if  
 283  $df(\bar{x})(u) \geq 0$  for all  $u \in \mathbb{R}^n$  and  $\inf_{w \in \mathbb{R}^n} d^2 f(\bar{x})(u|w) > 0$  for any nonzero  
 284  $u \in C_f(\bar{x}, 0)$ , then  $\bar{x}$  is a local minimizer of  $f$ . We refer the reader to [8] or [57,  
 285 Theorem 13.66] for details.

## 286 2.2 Smooth Minimization on a Manifold

287 The main results of this work exploit local smooth features of nonsmooth optimization  
 288 problems (c.f. Definition 2.6). In the presence of these features, the constructions of  
 289 Definition 2.1 locally simplify. Before moving to the general setting, we thus interpret  
 290 the various derivative constructions in the classical setting of minimizing a  $C^2$ -smooth  
 291 function  $f$  on a  $C^2$ -smooth manifold  $\mathcal{M}$ . To that end, we first recall the definition of  
 292 a manifold.

293 **Definition 2.2** (*Smooth manifold*) A subset  $\mathcal{M} \subset \mathbb{R}^n$  is a  $C^p$  manifold of dimension  
 294  $r$  around  $\bar{x} \in \mathcal{M}$  if there is an open neighborhood  $U$  around  $\bar{x}$  and a mapping  $G$  from  
 295  $\mathbb{R}^n$  to  $\mathbb{R}^{n-r}$  such that following hold:  $G$  is  $C^p$ -smooth, the derivative  $\nabla G(\bar{x})$  has full  
 296 rank, and equality holds:

$$297 \quad \mathcal{M} \cap U = \{x \in U : G(x) = 0\}.$$

298 We call  $G = 0$  the *local defining equations* for  $\mathcal{M}$  around  $\bar{x}$ . The *tangent space* to  
 299  $\mathcal{M}$  at  $\bar{x}$  is  $T_{\mathcal{M}}(\bar{x}) := \ker \nabla G(\bar{x})$  and the *normal space* to  $\mathcal{M}$  at  $\bar{x}$  is  $N_{\mathcal{M}}(\bar{x}) :=$   
 300  $\text{range } \nabla G(\bar{x})^*$ .

301 Turning to the classical setting, consider the optimization problem

$$302 \quad \min_{y \in \mathbb{R}^d} f(y) \quad \text{subject to } y \in \mathcal{M}. \quad (2.2)$$

303 Fix a point  $\bar{y} \in \mathcal{M}$  and suppose that both the function  $f$  is  $C^2$ -smooth around  $\bar{y}$  and  
 304  $\mathcal{M}$  is a  $C^2$ -smooth manifold around  $\bar{y}$ . Due to local smoothness, the subdifferential  
 305 admits the simple expression:

$$306 \quad \partial f_{\mathcal{M}}(\bar{y}) = \nabla f(\bar{y}) + N_{\mathcal{M}}(\bar{y}).$$

307 Recall that we use the shorthand  $f_{\mathcal{M}} := f + \delta_{\mathcal{M}}$ . From this expression, we see that a  
 308 point  $\bar{y} \in \mathcal{M}$  is first-order critical for the problem (2.2) precisely when the inclusion  
 309 holds:

$$310 \quad 0 \in \nabla f(\bar{y}) + N_{\mathcal{M}}(\bar{y}). \quad (2.3)$$

311 This inclusion can be equivalently stated in terms of the Lagrangian function. Namely,  
 312 let  $G = 0$  be the local defining equations for  $\mathcal{M}$  around  $\bar{y}$  and define the Lagrangian  
 313 function

$$314 \quad \mathcal{L}(y, \lambda) := f(y) + \langle G(y), \lambda \rangle.$$

315 Then, (2.3) amounts to existence of a (unique) multiplier vector  $\bar{\lambda} \in \mathbb{R}^m$  satisfying  
 316  $0 = \nabla_y \mathcal{L}(\bar{y}, \bar{\lambda})$ . Next, assuming  $\bar{y}$  is critical, second-order necessary conditions read

$$317 \quad \langle \nabla_{yy}^2 \mathcal{L}(\bar{y}, \bar{\lambda})u, u \rangle \geq 0 \quad \text{for all } u \in T_{\mathcal{M}}(\bar{y}). \quad (2.4)$$

318 Conversely, second-order sufficient conditions read

$$319 \quad \langle \nabla_{yy}^2 \mathcal{L}(\bar{y}, \bar{\lambda})u, u \rangle > 0 \quad \text{for all } 0 \neq u \in T_{\mathcal{M}}(\bar{y}). \quad (2.5)$$

320 It is well known that the sufficient condition (2.5) implies more than just local mini-  
 321 mality; namely, (2.5) holds if and only if there exists  $\alpha > 0$  such that

$$322 \quad f(y) - f(\bar{y}) \geq \alpha \|y - \bar{y}\|^2, \quad \text{for all } y \in \mathcal{M} \text{ near } \bar{y}. \quad (2.6)$$

323 Any point  $\bar{y}$  satisfying (2.6) is called a *strong local minimizer* of  $f$  on  $\mathcal{M}$ .

324 The Lagrangian conditions (2.4) and (2.5) may be succinctly expressed through  
 325 parabolic subderivatives of  $f_{\mathcal{M}}(y)$ , yielding a form independent of the choice of local  
 326 defining equations  $G = 0$ . In particular, a quick computation shows that for any  
 327  $u \in T_{\mathcal{M}}(\bar{y})$ , the function  $w \mapsto d^2 f(\bar{y})(u|w)$  is constant on its domain.<sup>6</sup> Dropping

<sup>6</sup> The domain of  $d^2 f_{\mathcal{M}}(\bar{y})(u|\cdot)$  consists of  $w$  satisfying  $(\langle \nabla^2 G_1(\bar{y})u, u \rangle, \dots, \langle \nabla^2 G_{n-r}(\bar{y})u, u \rangle) = -\nabla G(\bar{y})w$ , where  $G_i$  are the coordinate functions of  $G$ .

328 the dependence on  $w$ , the equation then holds:

329 
$$d^2 f_{\mathcal{M}}(\bar{y})(u) = \langle \nabla_{yy}^2 \mathcal{L}(\bar{y}, \bar{\lambda})u, u \rangle \quad \text{for all } u \in T_{\mathcal{M}}(\bar{y}).$$

330 The use of (2.5) goes far beyond verifying local optimality; indeed, this condition  
 331 plays a fundamental role in certifying solution stability under small perturbations. To  
 332 illustrate, consider the value function of the parametric family

333 
$$\varphi(x) = \inf_y \{f(x, y) : y \in \mathcal{M}\}, \quad (\mathcal{P}_x)$$
  
 334

335 where  $f$  is  $C^2$ -smooth and  $\mathcal{M} \subset \mathbb{R}^d$  is a closed set. Let  $\bar{y}$  be a minimizer of  $\mathcal{P}_{\bar{x}}$  for a  
 336 fixed parameter  $\bar{x}$ , and suppose that  $\mathcal{M}$  is a  $C^2$ -smooth manifold around  $\bar{y}$ . Let  $G = 0$   
 337 be the local defining equations for  $\mathcal{M}$  around  $\bar{y}$  and define the parametric Lagrangian  
 338 function

339 
$$\mathcal{L}(x, y, \lambda) = f(x, y) + \langle G(y), \lambda \rangle.$$

340 Since  $\bar{y}$  solves  $\mathcal{P}_{\bar{x}}$ , there is a multiplier vector  $\bar{\lambda}$  satisfying  $0 = \nabla_y \mathcal{L}(\bar{x}, \bar{y}, \bar{\lambda})$ .

341 The following perturbation result will form the core of our arguments. In short:  
 342 both the value function  $\varphi(x)$  and the minimizer of  $\mathcal{P}_x$  vary smoothly with  $x$ , provided  
 343 two mild conditions hold (level-boundedness and quadratic growth). Moreover, the  
 344 derivatives of both the value function and the solution maps can be computed explicitly.  
 345 For details and a much more general perturbation result, see [59, Theorem 3.1].

346 **Theorem 2.3** (Perturbation analysis) *Suppose that the following two properties hold.*

- 347 1. **(Level-boundedness)** *There exists a number  $\gamma > \varphi(\bar{x})$  and a neighborhood  $\mathcal{X}$  of*  
 348  $\bar{x}$  *such that the set*

349 
$$\bigcup_{x \in \mathcal{X}} \{y \in \mathcal{M} : f(x, y) \leq \gamma\} \quad \text{is bounded.}$$

- 350 2. **(Quadratic growth)** *The point  $\bar{y}$  is a strong local minimizer and a unique global*  
 351 *minimizer of  $\mathcal{P}_{\bar{x}}$ .*

352 *Define the partial Hessian matrices*

353 
$$H_{xx} = \nabla_{xx}^2 \mathcal{L}(\bar{x}, \bar{y}, \bar{\lambda}), \quad H_{xy} = \nabla_{xy}^2 \mathcal{L}(\bar{x}, \bar{y}, \bar{\lambda}), \quad H_{yy} = \nabla_{yy}^2 \mathcal{L}(\bar{x}, \bar{y}, \bar{\lambda}),$$

354 *and the quantities*

355 
$$\eta(h) = \min_{v \in T_{\mathcal{M}}(\bar{y})} \langle H_{xx}h, h \rangle + 2\langle H_{xy}v, h \rangle + \langle H_{yy}v, v \rangle,$$
  
 356 
$$\Phi(h) = \operatorname{argmin}_{v \in T_{\mathcal{M}}(\bar{y})} \langle H_{xx}h, h \rangle + 2\langle H_{xy}v, h \rangle + \langle H_{yy}v, v \rangle.$$
  
 357

358 Then, for every  $x$  near  $\bar{x}$ , the problem  $\mathcal{P}_x$  admits a unique solution  $y(x)$ , which varies  
 359  $C^1$ -smoothly and admits the first-order expansion

360 
$$\bar{y}(\bar{x} + h) = \bar{y} + \Phi(h) + o(\|h\|) \quad \text{as } h \rightarrow 0.$$

361 Moreover, the function  $\varphi$  is  $C^2$ -smooth around  $\bar{x}$  and admits the second-order expansion  
 362

363 
$$\varphi(\bar{x} + h) = \varphi(\bar{x}) + \langle \nabla_x f(\bar{x}, \bar{y}), h \rangle + \frac{1}{2} \eta(h) + o(\|h\|^2) \quad \text{as } h \rightarrow 0.$$

364 The two assumptions of the theorem play different roles. The level-boundedness  
 365 property ensures that the solutions of the perturbed problems  $\mathcal{P}_x$  lie in a compact  
 366 set around  $\bar{y}$ . The quadratic growth property in turn ensures smoothness of both the  
 367 solution map and the value function. In what follows, we will apply this result several  
 368 times. Both conditions will follow in all cases from the next simple lemma.

369 **Lemma 2.4** (Sufficient conditions for level boundedness) Consider a closed function  
 370  $\varphi: \mathbb{R}^d \times \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$  and fix a point  $\bar{x} \in \mathbb{R}^d$ . Suppose there exists  $\alpha > 0$  such  
 371 that for all  $x$  near  $\bar{x}$ , the function  $\varphi(x, \cdot)$  is  $\alpha$ -strongly convex and its minimizer  $y(x)$   
 372 varies continuously. Then,  $y(x)$  is a strong global minimizer of  $\varphi(x, \cdot)$  for all  $x$  near  $\bar{x}$ .  
 373 Moreover, there exists a neighborhood  $\mathcal{X}$  of  $\bar{x}$  such that for any real  $\gamma > \varphi(\bar{x}, y(\bar{x}))$ ,  
 374 the set

375 
$$\bigcup_{x \in \mathcal{X}} \{y \in \mathbb{R}^n : \varphi(x, y) \leq \gamma\} \quad \text{is bounded.}$$

376 **Proof** Strong convexity ensures there is a neighborhood  $\mathcal{X}$  of  $\bar{x}$  such that for any  
 377  $x \in \mathcal{X}$ , the estimate holds:

378 
$$\varphi(x, y(x)) + \frac{\alpha}{2} \|y - y(x)\|^2 \leq \varphi(x, y) \quad \forall y \in \mathbb{R}^n, \quad (2.7)$$

379 showing  $y(x)$  is a strong global minimizer of  $\varphi(x, \cdot)$ . Shrinking  $\mathcal{X}$  if necessary, we  
 380 may assume that  $y(\cdot)$  also varies continuously on  $\mathcal{X}$ . Choose any  $\delta > 0$ . Then, by  
 381 shrinking  $\mathcal{X}$  again and by leveraging both closedness of  $\varphi$  and continuity of  $y$ , we may  
 382 ensure that

383 
$$\|y(x) - y(\bar{x})\| \leq \delta \quad \text{and} \quad \varphi(x, y(x)) \geq \varphi(\bar{x}, y(\bar{x})) - \delta \quad \text{for all } x \in \mathcal{X}. \quad (2.8)$$

384 The proof will now follow quickly from the bound (2.8). Indeed, consider any points  
 385  $x \in \mathcal{X}$  and  $y \in \mathbb{R}^d$  satisfying  $\varphi(x, y) \leq \gamma$ . Then, (2.7) yields

386 
$$\|y - y(x)\| \leq \sqrt{\frac{2(\gamma - \varphi(x, y(x)))}{\alpha}}.$$

387 Applying (2.8) then gives the uniform bound

$$\|y - y(\bar{x})\| \leq \|y(x) - y(\bar{x})\| + \sqrt{\frac{2(\gamma - \varphi(x, y(x)))}{\alpha}} \leq \delta + \sqrt{\frac{2(\gamma + \delta - \varphi(\bar{x}, y(\bar{x})))}{\alpha}},$$

completing the proof. □

### 2.3 Weak Convexity and the Moreau Envelope

In general, the little- $o$  error term in the definition of  $\partial f(\bar{x})$  (Definition 2.1) may depend both on the base point  $\bar{x}$  and on the subgradient  $v$ . In this work, we focus on a particular class of functions for which the error in approximation is uniform. Namely, we focus on the class of  $\rho$ -weakly convex functions  $f: \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$ , meaning those for which the assignment  $x \mapsto f(x) + \frac{\rho}{2}\|x\|^2$  defines a convex function. Subgradients of a  $\rho$ -weakly convex function  $f$  automatically yield a uniform lower bound:

$$f(y) \geq f(x) + \langle v, y - x \rangle - \frac{\rho}{2}\|y - x\|^2, \quad \forall x, y \in \mathbb{R}^d, v \in \partial f(x). \quad (2.9)$$

A useful feature of weakly convex functions is that they admit a smooth approximation that preserves critical points. Setting the notation, fix a  $\rho$ -weakly convex function  $f: \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$  and a parameter  $\mu < \rho^{-1}$ . Define the *Moreau envelope* and the *proximal point map*, respectively:

$$f_\mu(x) = \inf_{y \in \mathbb{R}^d} \left\{ f(y) + \frac{1}{2\mu}\|y - x\|^2 \right\},$$

$$\text{prox}_{\mu f}(x) = \operatorname{argmin}_{y \in \mathbb{R}^d} \left\{ f(y) + \frac{1}{2\mu}\|y - x\|^2 \right\}.$$

We will use a few basic properties of these two constructions, summarized below.

**Lemma 2.5** (Moreau envelope and the proximal point map) *Consider a  $\rho$ -weakly convex function  $f: \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$  and fix a parameter  $\mu < \rho^{-1}$ . Then, the following are true.*

1. The envelope  $f_\mu$  is  $C^1$ -smooth with its gradient given by

$$\nabla f_\mu(x) = \mu^{-1}(x - \text{prox}_{\mu f}(x)). \quad (2.10)$$

2. The envelope  $f_\mu(\cdot)$  is  $\mu^{-1}$ -smooth and  $\frac{\rho}{1-\mu\rho}$ -weakly convex meaning:

$$-\frac{\rho}{2(1-\mu\rho)}\|x' - x\|^2 \leq f_\mu(x') - f_\mu(x) - \langle \nabla f_\mu(x), x' - x \rangle \leq \frac{1}{2\mu}\|x' - x\|^2, \quad (2.11)$$

for all  $x, x' \in \mathbb{R}^d$ .

3. The proximal map  $\text{prox}_{\mu f}(\cdot)$  is  $\frac{1}{1-\mu\rho}$ -Lipschitz continuous and the gradient map  $\nabla f_\mu$  is Lipschitz continuous with constant  $\max\{\mu^{-1}, \frac{\rho}{1-\mu\rho}\}$ .

416 4. The critical points of  $f$  and  $f_\mu$  coincide. In particular, they are exactly the fixed  
 417 points of the proximal map  $\text{prox}_{\mu f}$ .

418 **Proof** Claim 1 follows, for example, from [53, Theorem 4.4]. The left-hand side of  
 419 (2.11) is proved in [53, Theorem 5.2]. To see the right-hand side, observe

$$\begin{aligned}
 420 \quad f_\mu(x') &\leq f(\text{prox}_{\mu f}(x)) + \frac{1}{2\mu} \|\text{prox}_{\mu f}(x) - x'\|^2 \\
 421 \quad &= f_\mu(x) + \frac{1}{2\mu} \left( \|\text{prox}_{\mu f}(x) - x'\|^2 - \|x - \text{prox}_{\mu f}(x)\|^2 \right) \\
 422 \quad &= f_\mu(x) + \langle \mu^{-1}(x - \text{prox}_{\mu f}(x)), x' - x \rangle + \frac{1}{2\mu} \|x - x'\|^2. \\
 423
 \end{aligned}$$

424 Thus, claim 2 holds. The result [53, Theorem 4.4] shows that  $\text{prox}_{\mu f}(\cdot)$  is Lips-  
 425 chitz continuous with constant  $\frac{1}{1-\mu\rho}$ . Lipschitz continuity of  $\nabla f_\mu(\cdot)$  with constant  
 426  $\max\{\mu^{-1}, \frac{\rho}{1-\mu\rho}\}$  follows from (2.11) and Alexandrov’s theorem [57, Theorem 13.51].  
 427 Thus, claim 3 holds. Claim 4 is immediate from (2.10) and the observation that the  
 428 function  $y \mapsto f(y) + \frac{1}{2\mu} \|y - x\|^2$  is strongly convex for any  $x$ .  $\square$

429 **2.4 Active Manifolds**

430 The nonsmooth behavior of sets and functions arising in applications is typically far  
 431 from pathological and instead manifests in highly structured ways. Formalizing this  
 432 perspective we will assume that nonsmoothness, in a certain localized sense, only  
 433 occurs along an “active manifold.” This notion, introduced in [39] under the name of  
 434 partial smoothness and rooted in the earlier works [1,10–12,24–26,64], extends the  
 435 concept of *active sets* in nonlinear programming far beyond the classical setting. In  
 436 this work, we will take the related perspective developed in [19], since it will be most  
 437 expedient for our purpose.

438 Before giving the formal definition, we provide some intuition. Taking a geometric  
 439 view, we will assume that each critical point of a function  $f$  lies on a smooth manifold  
 440  $\mathcal{M}$ , and that the objective varies smoothly along the manifold, but sharply off of it.  
 441 For example, consider Fig. 2a: there, the function  $f(x, y) = |x| - y^2$  admits the  
 442 active manifold  $\mathcal{M} = \{0\} \times \mathbb{R}$  around its unique critical point (the origin). From an  
 443 algorithmic point of view, active manifolds are the sets that typical algorithms (e.g.,  
 444 proximal point, proximal gradient [31], and dual averaging [37]) identify in finite  
 445 time. Active manifolds also play a central role for sensitivity analysis, providing a  
 446 path to reduce such questions to the smooth setting. In particular, reasonable conditions  
 447 guarantee that the active manifold is smoothly traced out by critical points of slight  
 448 perturbations of the problem. We are now ready to state the formal definition.<sup>7</sup>

449 **Definition 2.6** (*Active manifold*) Consider a closed weakly convex function  $f: \mathbb{R}^d \rightarrow$   
 450  $\mathbb{R} \cup \{\infty\}$  and fix a set  $\mathcal{M} \subseteq \mathbb{R}^d$  containing a critical point  $\bar{x}$  of  $f$ . Then,  $\mathcal{M}$  is called

<sup>7</sup> What we call an *active manifold* here is called an *identifiable manifold* in [19]—the reference we most closely follow. The term active is more evocative in the context of the current work.

451 an active  $C^p$ -manifold around  $\bar{x}$  if there exist a neighborhood  $U$  around  $\bar{x}$  satisfying  
 452 the following.

- 453 • **(smoothness)** The set  $\mathcal{M} \cap U$  is a  $C^p$ -smooth manifold and the restriction of  $f$   
 454 to  $\mathcal{M} \cap U$  is  $C^p$ -smooth.
- 455 • **(sharpness)** The lower bound holds:

$$456 \quad \inf\{\|v\| : v \in \partial f(x), x \in U \setminus \mathcal{M}\} > 0.$$

457 If  $f$  admits an active manifold around a critical point  $\bar{x}$ , then it must be locally  
 458 unique: any two active manifolds at  $\bar{x}$  must coincide on a neighborhood of  $\bar{x}$  [19,  
 459 Proposition 2.4, Proposition 10.10].<sup>8</sup> Moreover, the critical cone  $C_f(\bar{x}, 0)$  coincides  
 460 with the tangent space  $T_{\mathcal{M}}(\bar{x})$  [19, Proposition 10.8]. With the definition of the active  
 461 manifold in hand, we can now introduce the strict saddle property for nonsmooth  
 462 functions.<sup>9</sup>

463 **Definition 2.7** (*Strict saddles*) Consider a weakly convex function  $f: \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$ .  
 464 Then, we say that a critical point  $\bar{x}$  is a *strict saddle* of  $f$  if there exists a  $C^2$ -active  
 465 manifold  $\mathcal{M}$  of  $f$  at  $\bar{x}$  and the inequality  $d^2 f_{\mathcal{M}}(\bar{x})(u) < 0$  holds for some vector  
 466  $u \in T_{\mathcal{M}}(\bar{x})$ . If every critical point of  $f$  is either a local minimizer or a strict saddle,  
 467 then we say that  $f$  satisfies the *strict saddle property*.

468 Looking at Fig. 2a, we see that the function  $f(x, y) = |x| - y^2$  indeed has the strict  
 469 saddle property: the restriction of  $f$  to the axis  $\mathcal{M} = \{0\} \times \mathbb{R}$ , namely  $f_{\mathcal{M}}(0, t) = -t^2$ ,  
 470 certainly has directions of negative curvature. Figure 2b depicts the subgradient flow  
 471 generated by this function. Notice that the set of initial conditions attracted to the origin  
 472 has measure zero. This observation suggests that typical algorithms are also unlikely to  
 473 stall at the strict saddle point, an observation made precise by the forthcoming results.

474 The curvature condition in the definition of the strict saddle pertains only to negative  
 475 curvature of the restriction of  $f$  to  $\mathcal{M}$ . One may instead ask whether existence of  
 476 directions of negative curvature for  $f$  alone suffice. The answer turns out to be yes.

477 **Theorem 2.8** ([18, Corollary 4.15]) Consider a closed weakly convex function  
 478  $f: \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$  that admits a  $C^3$ -active manifold  $\mathcal{M}$  around a critical point  
 479  $\bar{x}$ . Then, it holds:

$$480 \quad d^2 f(\bar{x})(u \mid w) \geq d^2 f_{\mathcal{M}}(\bar{x})(u) \quad \text{for all } u \in C_f(\bar{x}, 0), w \in \mathbb{R}^d.$$

481 A natural question is whether we expect the strict saddle property to hold typically.  
 482 One supporting piece of evidence is that the property holds under generic linear per-  
 483 turbations of semialgebraic problems.<sup>10</sup> This is almost immediate from guarantees

<sup>8</sup> Note that due to the convention  $\inf_{\emptyset} = +\infty$ , the entire space  $\mathcal{M} = \mathbb{R}^d$  is the active manifold for a globally  $C^p$ -smooth function  $f$  around any of its critical points.

<sup>9</sup> Better terminology would be the terms *active strict saddle* and the *active strict saddle property*. To streamline the notation, we omit the word active, as it should be clearly understood from context.

<sup>10</sup> A function is semi-algebraic if its graph can be written as a finite union of sets each cut out by finitely many polynomial inequalities.



484 in [18, Theorem 4.16], though this conclusion is not explicitly stated in the theorem  
 485 statement. We state this guarantee below and provide a quick proof in Sect. A for  
 486 completeness.

487 **Theorem 2.9** (Strict saddle property is generic). *Consider a closed, weakly convex,*  
 488 *semi-algebraic function  $f: \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$ . Then, for a full Lebesgue measure set of*  
 489 *perturbations  $v \in \mathbb{R}^d$ , the perturbed function  $x \mapsto f(x) - \langle v, x \rangle$  has the strict saddle*  
 490 *property.*

## 491 2.5 The Center Stable Manifold Theorem

492 In this work, we will show that a variety of simple algorithms escape strict saddle  
 493 points. To prove results of this type, we will interpret algorithms as fixed point iterations  
 494 of a nonlinear map  $T: \mathbb{R}^d \rightarrow \mathbb{R}^d$ , having certain favorable properties. As in the smooth  
 495 setting of [35], the core of our arguments will be based on the center stable manifold  
 496 theorem.

497 **Theorem 2.10** (The Center Stable Manifold Theorem [60, Theorem III.7]) *Let the*  
 498 *origin be a fixed point of the  $C^1$  local diffeomorphism  $T: U \rightarrow \mathbb{R}^d$  where  $U$  is a*  
 499 *neighborhood of the origin in  $\mathbb{R}^d$ . Let  $E^s \oplus E^c \oplus E^u$  be the invariant splitting of  $\mathbb{R}^d$*   
 500 *into the generalized eigenspaces of the Jacobian  $\nabla T(0)$  corresponding to eigenvalues*  
 501 *of absolute value less than one, equal to one, and greater than one. Then, there exists a*  
 502 *local  $T$  invariant  $C^1$  embedded disk  $W_{\text{loc}}^{\text{cs}}$ , tangent to  $E^s \oplus E^c$  at 0 and a neighborhood*  
 503  *$B$  around zero such that  $T(W_{\text{loc}}^{\text{cs}}) \cap B \subseteq W_{\text{loc}}^{\text{cs}}$ . In addition, if  $T^k(x) \in B$  for all  $k \geq 0$ ,*  
 504 *then  $x \in W_{\text{loc}}^{\text{cs}}$ .*

505 An immediate consequence of this theorem is the following: if  $\nabla T(0)$  is invertible  
 506 and has at least one eigenvalue of magnitude greater than one, then there exists a  
 507 neighborhood  $B$  of the origin such that the set

$$508 \quad \{x \in B : T^k(x) \in B \text{ for all } k \geq 0\},$$

509 has measure zero. This fact motivates the following key definition.

510 **Definition 2.11** (*Unstable fixed points*) A fixed point  $\bar{x}$  of a map  $T: \mathbb{R}^d \rightarrow \mathbb{R}^d$  is  
 511 called *unstable* if  $T$  is  $C^1$ -smooth around  $\bar{x}$  and the Jacobian  $\nabla T(\bar{x})$  has an eigenvalue  
 512 of magnitude strictly greater than one.

513 To globalize the guarantees of the center stable manifold theorem, we will need to  
 514 impose global regularity properties on  $T$ . In this work, we will require the map  $T$  to be  
 515 a *lipeomorphism*, namely, we require that  $T$  is globally Lipschitz and its inverse  $T^{-1}$   
 516 is a well-defined globally Lipschitz map. The following corollary is now immediate.  
 517 Its proof closely follows the presentation in [34, Theorem 2].

518 **Corollary 2.12** *Let  $T: \mathbb{R}^d \rightarrow \mathbb{R}^d$  be a lipeomorphism and let  $\mathcal{U}_T$  consist of all unstable*  
 519 *fixed points  $x$  of  $T$  at which the Jacobian  $\nabla T(x)$  is invertible. Then, the set of initial*

520 conditions attracted by such fixed points

Author Proof

521 
$$W := \left\{ x \in \mathbb{R}^d : \lim_{k \rightarrow \infty} T^k(x) \in \mathcal{U}_T \right\}$$

522 has zero Lebesgue measure.

523 **Proof** For every  $\bar{x} \in \mathcal{U}_T$  there exists a neighborhood  $U$  of  $\bar{x}$  such that  $T : U \rightarrow \mathbb{R}^d$  is  
 524 a local diffeomorphism. Thus, the center stable manifold theorem shows there exists  
 525 an open neighborhood  $B_{\bar{x}}$  of  $\bar{x}$  so that  $S_{\bar{x}} := \bigcap_{k=0}^{\infty} T^{-k}(B_{\bar{x}})$  is contained in a measure  
 526 zero set. In particular,  $S_{\bar{x}}$  itself is measure zero.

527 Now observe that  $\mathcal{U}_T \subseteq \bigcup_{\bar{x} \in \mathcal{U}_T} B_{\bar{x}}$  is an open cover of  $\mathcal{U}_T$ . Since  $\mathbb{R}^d$  is second  
 528 countable, this cover has a countable subcover  $\mathcal{U}_T \subseteq \bigcup_{i=1}^{\infty} B_{\bar{x}_i}$ . Observe the inclusion  
 529  $W \subseteq \bigcup_{i=1}^{\infty} \bigcup_{j=0}^{\infty} T^{-j}(S_{\bar{x}_i})$ . Since  $T$  is a lipeomorphism, the right-hand side is a  
 530 countable union of measure zero sets, and therefore,  $W$  has measure zero.  $\square$

531 To verify that a map  $T$  is a lipeomorphism, we will appeal to the following standard  
 532 sufficient condition. We provide a quick proof for completeness.

533 **Lemma 2.13** Let  $H : \mathbb{R}^d \rightarrow \mathbb{R}^d$  be a Lipschitz continuous map with constant  $\lambda < 1$ .  
 534 Then,  $I + H$  is invertible and  $(I + H)^{-1} : \mathbb{R}^d \rightarrow \mathbb{R}^d$  is Lipschitz continuous with  
 535 constant  $(1 - \lambda)^{-1}$ .

536 **Proof** To show that  $(I + H)$  is invertible, we must show that for every  $u \in \mathbb{R}^d$ , the  
 537 equation  $u = H(x) + x$  has a unique solution  $x(u) \in \mathbb{R}^d$ . Equivalently, we must show  
 538 that for every  $u \in \mathbb{R}^d$ , the mapping

539 
$$\zeta_u(x) := u - H(x)$$

540 has a unique fixed point. This is immediate from Banach's fixed point theorem since  
 541  $\zeta_u(\cdot)$  is strictly contractive.

542 To show that  $(I + H)^{-1}$  is Lipschitz, choose arbitrary  $u, v \in \mathbb{R}^d$  and define  $x :=$   
 543  $(I + H)^{-1}(u)$  and  $y := (I + H)^{-1}(v)$ . We then compute

544 
$$\begin{aligned} \|u - v\| &= \|(I + H)(x) - (I + H)(y)\| \geq \|x - y\| - \|H(x) - H(y)\| \\ &\geq (1 - \lambda)\|x - y\|, \end{aligned}$$

545  
546

547 where we have used the reverse triangle inequality and Lipschitz continuity of  $H$ .  
 548 Rearranging completes the proof.  $\square$

549 While the iteration mappings  $S$  of Sect. 1.3 can be Lipschitz, they are usually not  
 550 invertible. Thus, to ensure Lipschitz invertibility, we will consider damped fixed point  
 551 iterations, as summarized in the following elementary lemma. We provide a quick  
 552 proof for completeness.

553 **Lemma 2.14** (Damped fixed point iterations). Consider a map  $S : \mathbb{R}^d \rightarrow \mathbb{R}^d$  and fix  
 554 a damping parameter  $\alpha \in (0, 1)$ . Define the map

555 
$$T(x) = (1 - \alpha)x + \alpha \cdot S(x).$$

556 Then, the following are true.

- 557 1. The fixed points of  $T$  and  $S$  coincide.
- 558 2. If  $S$  is differentiable at  $\bar{x}$  and the Jacobian  $\nabla S(\bar{x})$  has a real eigenvalue strictly
- 559 greater than one, then  $\bar{x}$  is an unstable fixed point of  $T$ .
- 560 3. If the map  $S$  is continuous and the iterates generated by the process  $x_{k+1} = T(x_k)$
- 561 converge to some point  $\bar{x}$ , then  $\bar{x}$  must be a fixed point of  $S$ .
- 562 4. If the map  $I - S$  is  $L$ -Lipschitz, then  $T$  is a lipeomorphism for any  $\alpha \in (0, L^{-1})$ .

563 **Proof** Claims 1 and 2 follow directly from algebraic manipulations. Claim 4 follows  
 564 immediately from Lemma 2.13 by writing  $T = I + H$  with  $H = \alpha(S - I)$ . To see  
 565 claim 3, suppose that  $T$  is continuous and that  $x_k$  converge to some point  $\bar{x}$ . Then, we  
 566 deduce

$$567 \quad T(\bar{x}) = T\left(\lim_{k \rightarrow \infty} x_k\right) = \lim_{k \rightarrow \infty} T(x_k) = \lim_{k \rightarrow \infty} x_{k+1} = \bar{x}.$$

568 Therefore,  $\bar{x}$  is a fixed point of  $T$ . Using claim 1, we deduce that  $\bar{x}$  is a fixed point of  
 569  $S$ . □

### 570 3 The Proximal Point Method

571 We now turn to the saddle escape properties of the proximal-point method. Fixing the  
 572 problem at hand, we consider

$$573 \quad \min_{x \in \mathbb{R}^d} f(x),$$

574 where  $f: \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$  is a  $\rho$ -weakly convex function that is bounded from below.  
 575 For a fixed  $\mu < \rho^{-1}$ , the classical proximal-point method is precisely the fixed point  
 576 iteration

$$577 \quad x_{t+1} = \text{prox}_{\mu f}(x_t).$$

578 Key to our analysis is the equivalence between this algorithm and gradient descent on  
 579 the Moreau envelope. This equivalence follows from (2.10), which quickly yields the  
 580 description

$$581 \quad x_{k+1} = x_k - \mu \cdot \nabla f_\mu(x_k).$$

582 The saddle escape properties of the proximal point method thus flow from the strict  
 583 saddle properties of the Moreau envelope. Indeed, the following theorem shows that  
 584 when  $f$  admits a  $C^2$  active manifold around a critical point  $\bar{x}$ , the envelope  $f_\mu$  is  
 585 automatically  $C^2$ -smooth near  $\bar{x}$ . Moreover, if  $\bar{x}$  is a strict saddle of  $f$ , then it is also  
 586 a strict saddle of  $f_\mu$ . Consequently, any strict saddle point of  $f$  is an unstable fixed  
 587 point of the proximal map  $\text{prox}_{\mu f}(\cdot)$ .

588 **Theorem 3.1** (Saddle points of the Moreau envelope). *Let  $f: \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$  be a*  
 589 *closed and  $\rho$ -weakly convex function and let  $\bar{x}$  be any critical point of  $f$ . Suppose that*  
 590  *$f$  admits a  $C^2$  active manifold  $\mathcal{M}$  at  $\bar{x}$ . Then, for any  $\mu < \rho^{-1}$ , the Moreau envelope*  
 591  *$f_\mu$  is  $C^2$ -smooth around  $\bar{x}$  and its Hessian satisfies*

$$592 \quad \min_{h \in \mathbb{S}^{d-1} \cap T_{\mathcal{M}}(\bar{x})} \langle \nabla^2 f_\mu(\bar{x})h, h \rangle \leq \min_{h \in \mathbb{S}^{d-1} \cap T_{\mathcal{M}}(\bar{x})} d^2 f_{\mathcal{M}}(\bar{x})(h). \quad (3.1)$$

593 *Consequently, if  $\bar{x}$  is a strict saddle point of  $f$ , then  $\bar{x}$  is both a strict saddle point of*  
 594  *$f_\mu$  and an unstable fixed point of the proximal map  $\text{prox}_{\mu f}(\cdot)$ . Moreover,  $\nabla \text{prox}_{\mu f}(\bar{x})$*   
 595 *has a real eigenvalue that is strictly greater than one.*

596 **Proof** It is well known (for example, from [31]) that for all  $x$  near  $\bar{x}$ , the inclusion  
 597  $\text{prox}_{\mu f}(x) \in \mathcal{M}$  holds. From this inclusion, we will be able to view the proximal  
 598 subproblem through the lens of the perturbation result in Theorem 2.3. For the sake  
 599 of completeness, however, let us first quickly verify the claim. Consider a sequence  
 600  $x_i \rightarrow \bar{x}$  and observe the inclusion  $\nabla f_\mu(x_i) \in \partial f(\text{prox}_{\mu f}(x_i))$ . Since the gradient  $\nabla f_\mu$   
 601 is continuous, we deduce the limits  $\text{prox}_{\mu f}(x_i) \rightarrow \bar{x}$  and  $\nabla f_\mu(x_i) \rightarrow 0$ . Therefore,  
 602 by definition of the active manifold, we have  $\text{prox}_{\mu f}(x_i) \in \mathcal{M}$  for all sufficiency large  
 603 indices  $i$ , proving the claim.

604 Turning to the perturbation result, let  $F: \mathbb{R}^d \rightarrow \mathbb{R}$  be any  $C^2$ -smooth function  
 605 agreeing with  $f$  on a neighborhood of  $\bar{x}$  in  $\mathcal{M}$ .<sup>11</sup> Applying the claim, we find that the  
 606 equality

$$607 \quad f_\mu(x) = \min_{y \in \mathcal{M}} \left\{ F(y) + \frac{1}{2\mu} \|y - x\|^2 \right\},$$

608 holds for all  $x$  near  $\bar{x}$ . Our goal is to apply the perturbation result (Theorem 2.3) with  
 609  $f(x, y) := F(y) + \frac{1}{2\mu} \|y - x\|^2$  and  $\varphi(x) := f_\mu(x)$ . To that end, we now verify  
 610 the assumptions of Theorem 2.3. First, we verify the quadratic growth condition:  
 611 since we have chosen  $\mu < \rho^{-1}$ , it follows that for every  $x \in \mathbb{R}^d$  the function  $y \mapsto$   
 612  $f(x) + \frac{1}{2\mu} \|y - x\|^2$  is strongly convex with constant  $\mu^{-1} - \rho$ . Next, we verify the  
 613 level boundedness condition: since the minimizer  $y(x) := \text{prox}_{\mu f}(x)$  of this function  
 614 varies continuously and satisfies  $y(\bar{x}) = \bar{x}$ , the conditions of Lemma 2.4 are satisfied.  
 615 Therefore, the assumptions of Theorem 2.3 are valid.

616 We now apply Theorem 2.3. To that end, let  $G = 0$  be the defining equation of  $\mathcal{M}$   
 617 around  $\bar{x}$  and define the parametric Lagrangian function

$$618 \quad \mathcal{L}(x, y, \lambda) := F(y) + \frac{1}{2\mu} \|y - x\|^2 + \langle G(y), \lambda \rangle.$$

619 Since  $\bar{x}$  is critical for  $f$ , the equality  $\bar{x} = \text{prox}_{\mu f}(\bar{x})$  holds. Consequently,  $y(\bar{x}) =$   
 620  $\bar{x}$  minimizes the function  $y \mapsto F(y) + \frac{1}{2\mu} \|y - \bar{x}\|^2$  on  $\mathcal{M}$ . Therefore, first-order

<sup>11</sup> For example, let  $F$  be a  $C^2$  function defined on a neighborhood  $U$  of  $\bar{x}$  that agrees with  $f$  on  $U \cap \mathcal{M}$ . Using a partition of unity (e.g., [36, Lemma 2.26]), one can extend  $F$  from a slightly smaller neighborhood to be  $C^2$  on all of  $\mathbb{R}^d$ .

621 optimality conditions guarantee there exists a multiplier vector  $\bar{\lambda}$  satisfying

$$622 \quad 0 = \nabla_y \mathcal{L}(\bar{x}, \bar{x}, \bar{\lambda}) = \nabla F(\bar{x}) + \sum_{i \geq 1} \bar{\lambda}_i G_i(\bar{x}),$$

623 where  $G_i(\cdot)$  are the coordinate functions of  $G(\cdot)$ . Appealing to Theorem 2.3, we learn  
624 both that  $f_\mu$  is  $C^2$ -smooth around  $\bar{x}$  and that its Hessian satisfies

$$625 \quad \langle \nabla^2 f_\mu(\bar{x})h, h \rangle = \min_{u \in T_{\mathcal{M}}(\bar{x})} \langle H_{xx}h, h \rangle + 2\langle H_{xy}u, h \rangle + \langle H_{yy}u, u \rangle, \quad (3.2)$$

626 where the Hessian matrices are given by

$$627 \quad H_{xx} = \mu^{-1}I, \quad H_{xy} = -\mu^{-1}I, \quad H_{yy} = \nabla^2 F(\bar{x}) + \sum_{i=1}^m \bar{\lambda}_i \nabla^2 G_i(\bar{x}) + \mu^{-1}I.$$

628 Thus, rearranging (3.2) and setting  $D := \nabla^2 F(\bar{x}) + \sum_{i=1}^m \bar{\lambda}_i \nabla^2 G_i(\bar{x})$ , we have

$$629 \quad \langle \nabla^2 f_\mu(\bar{x})h, h \rangle = \min_{u \in T_{\mathcal{M}}(\bar{x})} \left\{ \langle Du, u \rangle + \mu^{-1} \|h - u\|^2 \right\}.$$

630 Therefore, we arrive at the estimate

$$631 \quad \min_{h \in \mathbb{S}^{d-1} \cap T_{\mathcal{M}}(\bar{x})} \langle \nabla^2 f_\mu(\bar{x})h, h \rangle = \min_{u \in T_{\mathcal{M}}(\bar{x})} \min_{h \in \mathbb{S}^{d-1} \cap T_{\mathcal{M}}(\bar{x})} \left\{ \langle Du, u \rangle + \mu^{-1} \|h - u\|^2 \right\}$$

$$632 \quad \leq \min_{h \in \mathbb{S}^{d-1} \cap T_{\mathcal{M}}(y_0)} \langle Dh, h \rangle = \min_{h \in \mathbb{S}^{d-1} \cap T_{\mathcal{M}}(\bar{x})} d^2 f_{\mathcal{M}}(\bar{x})(h),$$

633

634 thereby verifying (3.1). If  $\bar{x}$  is a strict saddle point of  $f$ , then (3.1) implies that  $\nabla^2 f_\mu(\bar{x})$   
635 has a strictly negative eigenvalue. From the expression  $\text{prox}_{\mu f} = I - \mu \nabla f_\mu$ , we  
636 therefore deduce that the Jacobian of  $\text{prox}_{\mu f}$  at  $\bar{x}$  has at least one real eigenvalue that  
637 is strictly greater than one. Consequently,  $\bar{x}$  is an unstable fixed point of  $\text{prox}_{\mu f}$ .  $\square$

638 Even if the proximal mapping has an unstable fixed-point, it often fails to meet the  
639 conditions of the center stable manifold theorem (Theorem 2.10). Indeed, the proximal  
640 mapping is generally not injective, even near critical points. To remedy this issue, we  
641 instead analyze a slightly damped version of the proximal point method

$$642 \quad x_{k+1} = (1 - \alpha)x_k + \alpha \cdot \text{prox}_{\mu f}(x_k),$$

643 where  $\alpha \in (0, 1)$  is a fixed constant. Reinterpreting this algorithm in terms of the  
644 Moreau envelope, we arrive at the recurrence

$$645 \quad x_{k+1} = x_k - (\alpha\mu) \cdot \nabla f_\mu(x_k). \quad (3.3)$$

646 Thus, the role of damping is clear: it still induces gradient descent on the Moreau  
647 envelope, but with a stepsize slightly below the “theoretically optimal” step  $\mu$ . This is

Author Proof

648 entirely in line with the saddle point escape guarantees for gradient descent in smooth  
649 minimization [35].

650 **Theorem 3.2** (Proximal point method: global escape). *Let  $f: \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$  be a*  
651 *closed and  $\rho$ -weakly convex function satisfying the strict saddle property. Choose a*  
652 *constant  $\mu < \rho^{-1}$  and a damping parameter  $\alpha \in (0, \min\{1, (\mu\rho)^{-1} - 1\})$ . With these*  
653 *choices, consider the algorithm*

$$654 \quad x_{k+1} = (1 - \alpha)x_k + \alpha \cdot \text{prox}_{\mu f}(x_k). \quad (3.4)$$

655 *Then, for almost all initializers  $x_0$ , the following holds: if the limit of  $\{x_k\}_{k \geq 0}$  exists,*  
656 *it must be a local minimizer of  $f$ .*

657 **Proof** Define the map  $S := \text{prox}_{\mu f}(x_k)$ . Lemma 2.5 guarantees that the map  $I - S =$   
658  $\mu \nabla f_\mu$  is Lipschitz continuous with constant  $\max\{1, \frac{\mu\rho}{1-\mu\rho}\}$ . Taking into account the  
659 range of  $\alpha$  and applying Lemma 2.14 and Theorem 3.1, we may deduce the following  
660 three properties: (1)  $T$  is a lipeomorphism, (2) the limit of the sequence  $x_k$ , if it exists,  
661 must be a critical point of  $f$ , and (3) if a critical point of  $f$  is not a local minimum,  
662 then it is an unstable fixed point of  $T$ . An application of Corollary 2.12 completes the  
663 proof.  $\square$

#### 664 4 The Proximal Gradient Method

665 We now turn to the saddle escape properties of the proximal gradient method. Fixing  
666 the problem at hand, we consider

$$667 \quad \min_{x \in \mathbb{R}^d} f(x) = g(x) + r(x), \quad (4.1)$$

668 where  $g: \mathbb{R}^d \rightarrow \mathbb{R}$  is a  $C^2$ -smooth function with  $\beta$ -Lipschitz gradient and  $r: \mathbb{R}^d \rightarrow$   
669  $\mathbb{R} \cup \{+\infty\}$  is a closed and  $\rho$ -weakly convex function. We assume throughout that  $f$  is  
670 bounded from below. For this problem, the proximal gradient method takes the form

$$671 \quad x_{k+1} = \text{prox}_{\mu r}(x_k - \mu \nabla g(x_k)).$$

672 Unlike the proximal point algorithm, the proximal gradient algorithm may not corre-  
673 spond to gradient descent on a smooth envelope of the problem. Still, as the following  
674 theorem shows, the iteration mapping is  $C^1$  smooth near  $\bar{x}$  whenever  $f$  admits a  $C^2$   
675 active manifold around a critical point  $\bar{x}$ . Moreover, if  $\bar{x}$  is a strict saddle point of  $f$ ,  
676 then  $\bar{x}$  is an unstable fixed point of the iteration mapping

677 **Theorem 4.1** (Unstable fixed points of the prox-gradient map). *Consider the optimiza-*  
678 *tion problem (4.1) and let  $\bar{x}$  be any critical point of  $f$ . Suppose that  $f$  admits a  $C^2$*   
679 *active manifold  $\mathcal{M}$  at  $\bar{x}$ . Then, for any  $\mu \in (0, \rho^{-1})$ , the proximal-gradient map*

$$680 \quad S(x) := \text{prox}_{\mu r}(x - \mu \nabla g(x))$$

681 is  $C^1$ -smooth on a neighborhood of  $\bar{x}$ . Moreover, if  $\bar{x}$  is a strict saddle point of  $f$ , then  
 682  $\nabla S(\bar{x})$  has a real eigenvalue that is strictly greater than one.

683 **Proof** It is well known (for example, from [31]) that for all  $x$  near  $\bar{x}$ , the point  $S(x)$  lies  
 684 in  $\mathcal{M}$ . From this inclusion, we will be able to view the proximal subproblem through the  
 685 lens of the perturbation result in Theorem 2.3. For the sake of completeness, however,  
 686 we provide a quick proof. Indeed, consider a sequence  $x_i \rightarrow \bar{x}$  and set  $y_i = S(x_i)$ .  
 687 Then, by definition of the proximal gradient map, we have  $0 \in \nabla g(x_i) + \mu^{-1}(y_i -$   
 688  $x_i) + \partial r(y_i)$ , and therefore

$$689 \quad \text{dist}(0, \partial f(y_i)) = \text{dist}(-\nabla g(y_i), \partial r(y_i)) \leq \text{dist}(-\nabla g(x_i), \partial r(y_i)) + \beta \|y_i - x_i\|$$

$$690 \quad \leq (\mu^{-1} + \beta) \|y_i - x_i\|.$$

692 Since  $S(\cdot)$  is continuous and  $S(\bar{x}) = \bar{x}$ , we deduce  $y_i \rightarrow \bar{x}$  and therefore  
 693  $\text{dist}(0, \partial f(y_i)) \rightarrow 0$ . Therefore, the points  $y_i$  lie in  $\mathcal{M}$  for all sufficiently large indices  
 694  $i$ , proving the claim.

695 Turning to the perturbation result, let  $R: \mathbb{R}^d \rightarrow \mathbb{R}$  be any  $C^2$ -smooth function  
 696 agreeing with  $r$  on a neighborhood of  $\bar{x}$  in  $\mathcal{M}$ . Applying the claim, we find that for  $x$   
 697 near  $\bar{x}$ , the point  $S(x)$  uniquely minimizes problem

$$698 \quad \min_{y \in \mathcal{M}} \left\{ g(x) + \langle \nabla g(x), y - x \rangle + R(y) + \frac{1}{2\mu} \|y - x\|^2 \right\}. \quad (\mathcal{P}_x)$$

700 Our goal is to apply the perturbation result (Theorem 2.3) with  $f(x, y) := g(x) +$   
 701  $\langle \nabla g(x), y - x \rangle + R(y) + \frac{1}{2\mu} \|y - x\|^2$ . To that end, we now verify the assumptions of  
 702 Theorem 2.3. First, we verify the quadratic growth condition: since we have chosen  
 703  $\mu < \rho^{-1}$ , it follows that for every  $x \in \mathbb{R}^d$  the function  $y \mapsto f(x, y)$  is strongly convex  
 704 with the constant  $\mu^{-1} - \rho$ . Next, we verify the level-boundedness condition: since the  
 705 minimizer  $S(x)$  clearly varies continuously and satisfies  $S(\bar{x}) = \bar{x}$ , the conditions of  
 706 Lemma 2.4 are satisfied. Therefore, the assumptions of Theorem 2.3 are valid.

707 We now apply Theorem 2.3. To that end, let  $G = 0$  be the defining equation of  $\mathcal{M}$   
 708 around  $\bar{x}$  and define the parametric Lagrangian function

$$709 \quad \mathcal{L}(x, y, \lambda) = g(x) + \langle \nabla g(x), y - x \rangle + R(y) + \frac{1}{2\mu} \|y - x\|^2 + \sum_{i \geq 1} \lambda_i G_i(y),$$

710 where  $G_i(\cdot)$  are the coordinate functions of  $G$ . Clearly  $y(\bar{x}) = \bar{x}$  minimizes  $f(\bar{x}, \cdot)$   
 711 on  $\mathcal{M}$ . Therefore, first-order optimality conditions guarantee there exists a multiplier  
 712 vector  $\bar{\lambda}$  satisfying

$$713 \quad 0 = \nabla_y \mathcal{L}(\bar{x}, \bar{x}, \bar{\lambda}) = \nabla g(\bar{x}) + \nabla R(\bar{x}) + \sum_{i \geq 1} \bar{\lambda}_i G_i(\bar{x}).$$

714 Appealing to Theorem 2.3, we learn that the solution map  $S(\cdot)$  is  $C^1$ -smooth around  
 715  $\bar{x}$  with

$$716 \quad \nabla S(\bar{x})h = \underset{v \in T_{\mathcal{M}}(\bar{x})}{\text{argmin}} \quad 2\langle H_{xy}v, h \rangle + \langle H_{yy}v, v \rangle, \quad (4.2)$$

717 where the Hessian matrices are given by

718 
$$H_{xy} = \nabla^2 g(\bar{x}) - \mu^{-1}I, \quad H_{yy} = \nabla^2 R(\bar{x}) + \mu^{-1}I + \sum_{i=1}^p \bar{\lambda}_i \nabla^2 G_i(\bar{x}).$$

719 We now simplify the expression (4.2). To that end, let  $W$  be the orthogonal projection  
 720 onto  $T_{\mathcal{M}}(\bar{x})$  and define the linear maps  $\overline{H_{yy}}: T_{\mathcal{M}}(\bar{x}) \rightarrow T_{\mathcal{M}}(\bar{x})$  and  $\overline{H_{xy}}: T_{\mathcal{M}}(\bar{x}) \rightarrow$   
 721  $T_{\mathcal{M}}(\bar{x})$  by setting  $\overline{H_{yy}} = WH_{yy}W$  and  $\overline{H_{xy}} = WH_{xy}W$ , respectively. Since  $\bar{x}$  is a  
 722 strong local minimizer of  $\mathcal{P}_{\bar{x}}$ , the map  $\overline{H_{yy}}$  is positive definite, and hence invertible.  
 723 Solving (4.2) then yields the expression

724 
$$\nabla S(\bar{x})h = -\overline{H_{yy}}^{-1}\overline{H_{xy}}^\top h \quad \text{for all } h \in T_{\mathcal{M}}(\bar{x}).$$

725 Note that  $\overline{H_{xy}}^\top$  is a symmetric matrix, so we drop the “ $\top$ ” throughout.

726 Let us now verify that if  $\bar{x}$  is a strict saddle of  $f$ , then  $\nabla S(\bar{x})$  has a real eigenvalue  
 727 that is greater than one. To this end, observe that  $\gamma \in \mathbb{R}$  is a real eigenvalue of  $\nabla S(\bar{x})$   
 728 with an associated eigenvector  $v \in T_{\mathcal{M}}(\bar{x})$  if and only if

729 
$$\nabla S(\bar{x})v = \gamma v \iff -\overline{H_{yy}}^{-1}\overline{H_{xy}}v = \gamma v \iff (\gamma\overline{H_{yy}} + \overline{H_{xy}})v = 0.$$

730 In particular, if the matrix  $\gamma\overline{H_{yy}} + \overline{H_{xy}}$  is singular, then  $\gamma$  is an eigenvalue of  $\nabla S(\bar{x})$ .  
 731 To prove such a  $\gamma$  exists, we will examine the following ray of symmetric matrices

732 
$$\{\gamma\overline{H_{yy}} + \overline{H_{xy}} : \gamma \geq 1\}.$$

733 Beginning with the end point, the strict saddle property shows that

734 
$$\overline{H_{yy}} + \overline{H_{xy}} = W \left( \nabla^2 g(\bar{x}) + \nabla^2 R(\bar{y}) + \sum_i \bar{\lambda}_i \nabla^2 G_i(\bar{x}) \right) W.$$

735 has a strictly negative eigenvalue. On the other hand, the direction of the ray  $\overline{H_{yy}}$  is a  
 736 positive definite matrix. Therefore, by continuity of the minimal eigenvalue function,  
 737 there exists some  $\gamma > 1$  such that the matrix  $\gamma\overline{H_{yy}} + \overline{H_{xy}}$  is singular, as claimed.  $\square$

738 Similar to the proximal point method, the proximal gradient mapping fails to meet  
 739 the conditions of the center stable manifold theorem (Theorem 2.10), since it generally  
 740 lacks invertibility. Therefore, as before we will analyze a slightly damped version of  
 741 the process, and prove the following theorem.

742 **Theorem 4.2** (Proximal gradient method: global escape). *Consider the optimization*  
 743 *problem (4.1) and suppose that  $f$  has the strict saddle property. Choose any constant*  
 744  *$\mu \in (0, \rho^{-1})$  and a damping parameter  $\alpha \in (0, 1)$  satisfying*

745 
$$\alpha \cdot \left( \mu\beta + (1 + \mu\beta) \max \left\{ 1, \frac{\mu\rho}{1-\mu\rho} \right\} \right) < 1.$$



746 Consider the algorithm

$$747 \quad x_{k+1} = (1 - \alpha)x_k + \alpha \cdot \text{prox}_{\mu r}(x_k - \mu \nabla g(x_k)). \quad (4.3)$$

748 Then, for almost all initializers  $x_0$ , the following holds: if the limit of  $\{x_k\}_{k \geq 0}$  exists,  
749 it must be a local minimizer of  $f$ .

750 **Proof** Define the maps  $S = \text{prox}_{\mu r}(I - \mu \nabla g)$ . We successively rewrite

$$751 \quad \begin{aligned} I - S &= (I - \mu \nabla g) - \text{prox}_{\mu r}(I - \mu \nabla g) + \mu \nabla g \\ 752 \quad &= \mu \cdot \nabla r_{\mu} \circ (I - \mu \nabla g) + \mu \nabla g. \end{aligned}$$

754 Lemma 2.5 implies that the map  $I - S$  is Lipschitz continuous with constant  $\mu\beta + (1 +$   
755  $\mu\beta) \max\left\{1, \frac{\mu\rho}{1-\mu\rho}\right\}$ . Taking into account the range of  $\alpha$  and applying Lemma 2.14 and  
756 Theorem 4.1, we may deduce the following three properties: (1)  $T$  is a lipeomorphism,  
757 (2) the limit of the sequence  $x_k$ , if it exists, must be a critical point for  $f$ , and (3) if a  
758 critical point of  $f$  is not a local minimum, then it is an unstable fixed point of  $T$ . An  
759 application of Corollary 2.12 then completes the proof.  $\square$

## 760 5 The Proximal Linear Method

761 We now turn to the saddle escape properties of the proximal linear method, a gen-  
762 eralization of the proximal point and proximal gradient methods. Setting the stage,  
763 consider the composite optimization problem

$$764 \quad \min_x f(x) = h(F(x)) + r(x), \quad (5.1)$$

765 where  $F: \mathbb{R}^d \rightarrow \mathbb{R}^m$  is a  $C^2$ -smooth map,  $h: \mathbb{R}^d \rightarrow \mathbb{R}$  is convex, and  $r: \mathbb{R}^d \rightarrow$   
766  $\mathbb{R} \cup \{\infty\}$  is  $\rho$ -weakly convex. As is standard in the literature, we will assume that there  
767 exists a constant  $\beta > 0$  satisfying

$$768 \quad |h(F(y)) - h(F(x) + \nabla F(x)(y - x))| \leq \frac{\beta}{2} \|y - x\|^2, \quad \forall x, y \in \mathbb{R}^d. \quad (5.2)$$

769 These assumptions then easily imply that  $f$  is weakly convex with constant  $\beta + \rho$ .

770 With the stage set, we now slightly refine the notion of a strict saddle, adapting it  
771 to the compositional nature of the problem. This refinement intuitively asks that the  
772 active manifold for  $f$  at a critical point  $\bar{x}$  is induced by active manifolds of  $h$  and  $r$ .  
773 Similar conditions have appeared elsewhere, for example, in [19,39,40]. To describe  
774 the condition formally, we will also revise the definition of an active manifold, allowing  
775 us to discuss active manifolds of  $h(\cdot)$  and  $r(\cdot)$  at noncritical points. The revision is  
776 intuitive, requiring just a linear tilt of the functions:

- 777 • Consider a set  $\mathcal{R} \subset \mathbb{R}^d$ , a point  $x \in \mathcal{R}$ , and a subgradient  $v \in \partial r(x)$ . We will say  
778 that  $\mathcal{R}$  is a  $C^2$  active manifold of  $r$  at  $x$  for  $v$  if  $\mathcal{R}$  is a  $C^2$  active manifold of the  
779 tilted function  $r - \langle v, \cdot \rangle$  at  $x$  in the sense of Definition 2.6.

780 We may likewise define the active manifold of  $h$  at  $z$  for  $w \in \partial h(z)$ , based on a tilting  
 781 of  $h$  by  $w$ . Coupling these definitions, we arrive at the active manifold concept for the  
 782 composite problem (5.1).

783 **Definition 5.1** (*Composite active manifold*) Consider the compositional problem (5.1)  
 784 and let  $\bar{x}$  be a critical point of  $f$ . Fix arbitrary vectors  $\bar{w} \in \partial h(F(\bar{x}))$  and  $\bar{v} \in \partial r(\bar{x})$   
 785 satisfying

$$0 \in \nabla F(\bar{x})^* \bar{w} + \bar{v}. \tag{5.3}$$

787 Suppose the following hold.

- 788 1. There exist  $C^2$ -smooth manifolds  $\mathcal{R} \subset \mathbb{R}^d$  and  $\mathcal{H} \subset \mathbb{R}^m$  containing  $\bar{x}$  and  $F(\bar{x})$ ,  
 789 respectively, and satisfying the transversality condition:

$$\nabla F(\bar{x}) [T_{\mathcal{R}}(\bar{x})] + T_{\mathcal{H}}(F(\bar{x})) = \mathbb{R}^m. \tag{5.4}$$

- 791 2.  $\mathcal{R}$  is an active manifold of  $r$  at  $\bar{x}$  for  $\bar{v}$  and  $\mathcal{H}$  is an active manifold of  $h$  at  $F(\bar{x})$   
 792 for  $\bar{w}$ .

793 Then, we will call  $\mathcal{M} := \mathcal{R} \cap F^{-1}(\mathcal{H})$  a *composite  $C^2$  active manifold* for the  
 794 problem (5.1) at  $\bar{x}$ . If in addition the inequality  $d^2 f_{\mathcal{M}}(\bar{x})(u) < 0$  holds for some  
 795 vector  $u \in T_{\mathcal{M}}(\bar{x})$ , then we will call  $\bar{x}$  a *composite strict saddle point*.

796 This definition has several important subtleties. First, the set  $\mathcal{M} := \mathcal{R} \cap F^{-1}(\mathcal{H})$   
 797 is indeed a  $C^2$ -smooth manifold around  $\bar{x}$ , due to the classical transversality condition  
 798 (5.4), a central fact in differential geometry [36, Theorem 6.30]. Next, the vectors  $\bar{v}$   
 799 and  $\bar{w}$  do exist. This follows since  $\bar{x}$  is first-order critical for  $f$ :

$$0 \in \nabla F(\bar{x})^* \partial h(F(\bar{x})) + \partial r(\bar{x}).$$

801 Beyond existence, the vectors  $\bar{v} \in \partial r(\bar{x})$  and  $\bar{w} \in \partial h(F(\bar{x}))$  are in fact the unique  
 802 elements satisfying (5.3), a second consequence of transversality. To see this, we state  
 803 (5.4) in dual terms as

$$(\nabla F(\bar{x})^*)^{-1} N_{\mathcal{R}}(\bar{x}) \cap N_{\mathcal{H}}(F(\bar{x})) = \{0\}. \tag{5.5}$$

805 Considering another pair  $v \in \partial r(\bar{x})$  and  $w \in \partial h(F(\bar{x}))$  satisfying (5.3), we deduce

$$0 = \nabla F^*(\bar{x})(\bar{w} - w) + (\bar{v} - v).$$

807 To conclude  $v = \bar{v}$  and  $w = \bar{w}$ , we use (5.5) and simply recall that  $\text{span } \partial h(F(\bar{x})) =$   
 808  $N_{\mathcal{H}}(F(\bar{x}))$  and  $\text{span } \partial r(\bar{x}) = N_{\mathcal{R}}(\bar{x})$ , as shown in [19, Proposition 10.12]. Finally,  
 809 collecting these facts together, it follows from the chain rule [19, Proposition 5.1] that  
 810  $\mathcal{M}$  is an active manifold of  $f$  at  $\bar{x}$  in the sense of Definition 2.7.

811 A natural question is whether we expect the composite strict saddle property to hold  
 812 typically. One supporting piece, of evidence, analogous to Theorem 2.9, is that the  
 813 property holds under generic linear perturbations of semialgebraic composite problems.  
 814 This result quickly follows from [18, Theorem 5.2]. We provide a proof sketch  
 815 in Sect. A.

816 **Theorem 5.2** (Strict saddle property is generic). *Consider the composite problem (5.1),*  
 817 *where  $h$ ,  $r$ , and  $F$  are in addition semi-algebraic. Then, for a full Lebesgue measure*  
 818 *set of perturbations  $(y, v) \in \mathbb{R}^m \times \mathbb{R}^d$ , the problem*

$$819 \quad \min_x h(F(x) + y) + r(x) - \langle v, x \rangle$$

820 *has the composite strict saddle property.*

821 Turning to our central task, we aim to analyze the saddle escape properties of the  
 822 proximal linear method:

$$823 \quad x_{k+1} = \operatorname{argmin}_y h(F(x_k) + \nabla F(x_k)(y - x_k)) + r(y) + \frac{1}{2\mu} \|y - x_k\|^2.$$

824 To analyze this method, we prove the following theorem, showing that any strict  
 825 saddle point of the composite problem (5.1) is an unstable fixed point of proximal  
 826 linear update.

827 **Theorem 5.3** (Unstable fixed points of the proximal linear map). *Consider the com-*  
 828 *posite problem (5.1) and let  $\bar{x}$  be any critical point of  $f$ . Suppose the problem admits*  
 829 *a composite  $C^2$  active manifold  $\mathcal{M}$  at  $\bar{x}$ . Then, for any  $\mu \in (0, \rho^{-1})$ , the proximal*  
 830 *linear map*

$$831 \quad S(x) := \operatorname{argmin}_y h(F(x) + \nabla F(x)(y - x)) + r(y) + \frac{1}{2\mu} \|y - x\|^2. \quad (5.6)$$

832 *is  $C^1$ -smooth on a neighborhood of  $\bar{x}$ . Moreover, if  $\bar{x}$  is a composite strict saddle point,*  
 833 *then the Jacobian  $\nabla S(\bar{x})$  has a real eigenvalue strictly greater than one.*

834 In most ways, the proof mirrors that of Theorem 3.2. There is, however, an impor-  
 835 tant complication: we must move beyond the perturbation result of Theorem 2.3 and  
 836 instead analyze a parametric family of optimization problems where both the objective  
 837 and the constraints depend on a perturbation parameter. Therefore, we will rely on  
 838 the following generalization of Theorem 2.3. For details and a much more general  
 839 perturbation result, see [59, Theorem 4.2].

840 **Theorem 5.4** (Perturbation analysis). *Consider the family of optimization problems*

$$841 \quad \min_y f(x, y) \quad \text{subject to} \quad G(x, y) = 0 \quad (\mathcal{Q}_x)$$

842 *Fix a point  $\bar{x}$  and a minimizer  $\bar{y}$  of  $\mathcal{Q}_{\bar{x}}$ , and suppose the following hold.*

- 844 1. **(Nondegeneracy)** *The function  $f(\cdot, \cdot)$  and the map  $G(\cdot, \cdot)$  are  $C^2$ -smooth near*  
 845  *$(\bar{x}, \bar{y})$ , and the Jacobian  $\nabla_y G(\bar{x}, \bar{y})$  is surjective.*
- 846 2. **(Level-boundedness)** *There exists a neighborhood  $\mathcal{X}$  of  $\bar{x}$  and a number  $\gamma$  greater*  
 847 *than the minimal value of  $\mathcal{Q}_{\bar{x}}$  such that the set*

$$848 \quad \bigcup_{x \in \mathcal{X}} \{y \in Y(x) : f(x, y) \leq \gamma\} \quad \text{is bounded,}$$

849 where  $Y(x) := \{y : G(x, y) = 0\}$  denotes the set of feasible points for  $\mathcal{Q}_x$ .  
 850 3. **(Quadratic growth)** The point  $\bar{y}$  is a strong local minimizer and a unique global  
 851 minimizer of  $\mathcal{Q}_{\bar{x}}$ .

852 Define the parametric Lagrangian function

$$853 \quad \mathcal{L}(x, y, \lambda) = f(x, y) + \langle G(x, y), \lambda \rangle.$$

854 Fix the multiplier vector  $\bar{\lambda}$  satisfying  $0 = \nabla_y \mathcal{L}(\bar{x}, \bar{y}, \bar{\lambda})$  and define the Hessian matrices  
 855

$$856 \quad H_{xx} = \nabla_{xx}^2 \mathcal{L}(\bar{x}, \bar{y}, \bar{\lambda}), \quad H_{xy} = \nabla_{xy}^2 \mathcal{L}(\bar{x}, \bar{y}, \bar{\lambda}), \quad H_{yy} = \nabla_{yy}^2 \mathcal{L}(\bar{x}, \bar{y}, \bar{\lambda}).$$

857 Then, for every  $x$  near  $\bar{x}$ , the problem  $\mathcal{Q}_x$  admits a unique solution  $y(x)$ , which varies  
 858  $C^1$ -smoothly. Moreover, its directional derivative in direction  $h$  given by

$$859 \quad \begin{aligned} \nabla y(\bar{x})h &= \underset{v}{\operatorname{argmin}} 2 \langle H_{xy}v, h \rangle + \langle H_{yy}v, v \rangle \\ &\text{s.t.} \quad \nabla_x G(\bar{x}, \bar{y})h + \nabla_y G(\bar{x}, \bar{y})v = 0. \end{aligned} \quad (5.7)$$

860 With these tools in hand, we now prove Theorem 5.3.

861 **Proof of Theorem 5.3** Let  $\bar{v}, \bar{w}, \mathcal{H}, \mathcal{R}$ , and  $\mathcal{M}$  be the vectors and manifolds specified in  
 862 Definition 5.1. It is known from [40, Theorem 4.11] that for all  $x$  near  $\bar{x}$ , the inclusions  
 863 hold:

$$864 \quad S(x) \in \mathcal{M} \quad \text{and} \quad F(x) + \nabla F(x)(S(x) - x) \in \mathcal{H}.$$

865 From this inclusion, we will be able to view the proximal subproblem through the lens  
 866 of the perturbation result in Theorem 5.4. For the sake of completeness, however, we  
 867 provide a quick proof. Indeed, consider a sequence  $x_i \rightarrow \bar{x}$  and define  $z_i = F(x_i) +$   
 868  $\nabla F(x_i)(S(x_i) - x_i)$ . Then, appealing to the optimality conditions of the proximal linear  
 869 subproblem, we deduce that there exist vectors  $v_i \in \partial r(x_i)$  and  $w_i \in \partial h(z_i)$  satisfying  
 870  $\frac{1}{\mu}(x_i - S(x_i)) = \nabla F(x_i)^* w_i + v_i$ . Since  $S(\cdot)$  is continuous and  $h$  is Lipschitz, the  
 871 vectors  $w_i$  and  $v_i$  are bounded. Passing to a subsequence, we may assume that  $w_i$  and  
 872  $v_i$  converge to some  $w \in \partial h(F(\bar{x}))$  and  $v \in \partial r(\bar{x})$ , respectively, and moreover, that

$$873 \quad 0 \in \nabla F(\bar{x})^* w + v.$$

874 We therefore deduce  $w = \bar{w}$  and  $v = \bar{v}$ . Taking into account that  $\mathcal{R}$  is a  $C^2$ -active  
 875 manifold at  $\bar{x}$  for  $\bar{v}$  and  $\mathcal{H}$  is a  $C^2$ -active manifold at  $F(\bar{x})$  for  $\bar{w}$ , we deduce  $S(x_i) \in \mathcal{R}$   
 876 and  $z_i \in \mathcal{H}$  for all large indices  $i$ , proving the claim.

877 Turning to the perturbation result, let  $\hat{h}: \mathbb{R}^m \rightarrow \mathbb{R}$  be any  $C^2$ -smooth function  
 878 agreeing with  $h$  on a neighborhood of  $F(\bar{x})$  in  $\mathcal{H}$ , and let  $\hat{r}: \mathbb{R}^d \rightarrow \mathbb{R}$  be any  $C^2$ -  
 879 smooth function agreeing with  $r$  on a neighborhood of  $\bar{x}$  in  $\mathcal{R}$ . Applying the claim,  
 880 we find that for  $x$  near  $\bar{x}$ , we may write

$$\begin{aligned}
 881 \quad S(x) &= \operatorname{argmin}_y \hat{h}(F(x) + \nabla F(x)(y - x)) + \hat{r}(y) + \frac{1}{2\mu} \|y - x\|^2 \\
 &\text{s.t. } F(x) + \nabla F(x)(y - x) \in \mathcal{H} \quad \text{and } y \in \mathcal{R}
 \end{aligned} \tag{5.8}$$

882 Our goal is to apply the perturbation result (Theorem 5.4) to the parametric family  
 883 (5.8). To this end, let  $\omega = 0$  be the local defining equations of  $\mathcal{H}$  around  $F(\bar{x})$  and  
 884 let  $\eta = 0$  be the local defining equation of  $\mathcal{R}$  around  $\bar{x}$ . We can now place (5.8) in the  
 885 setting of Theorem 5.4 by setting

$$886 \quad f(x, y) = \hat{h}(F(x) + \nabla F(x)(y - x)) + \hat{r}(y) + \frac{1}{2\mu} \|y - x\|^2$$

887 and

$$888 \quad G(x, y) := (G^{\mathcal{H}}(x, y), G^{\mathcal{R}}(x, y)) := (\omega(F(x) + \nabla F(x)(y - x)), \eta(y)).$$

889 For these functions, we now verify the assumptions of Theorem 2.3. First, the nonde-  
 890 generacy property follows from the transversality condition (5.4). Second, we verify  
 891 the quadratic growth condition: since we have chosen  $\mu < \rho^{-1}$ , it follows that for  
 892 every  $x \in \mathbb{R}^d$  the function  $y \mapsto f(x, y)$  is strongly convex with the constant  $\mu^{-1} - \rho$ .  
 893 Finally, we verify the level-boundedness condition: since the minimizer  $S(x)$  clearly  
 894 varies continuously and satisfies  $S(\bar{x}) = \bar{x}$ , the conditions of Lemma 2.4 are satisfied.  
 895 Therefore, the assumptions of Theorem 2.3 are valid. In particular, we learn that the  
 896 solution map  $S(\cdot)$  is  $C^1$ -smooth around  $\bar{x}$ .

897 Computing the Jacobian of the solution mapping will occupy the remainder of the  
 898 proof. To that end, define the parametric Lagrangian

$$899 \quad \mathcal{L}(x, y, \lambda) = f(x, y) + \langle G(x, y), \lambda \rangle.$$

900 Localizing, the identification properties then entail that  $y = \bar{x}$  is a minimizer of the  
 901 problem (5.8) corresponding to  $x = \bar{x}$ . We conclude there exists a Lagrange multiplier  
 902 vector  $\bar{\lambda} = (\bar{\lambda}^{\mathcal{H}}, \bar{\lambda}^{\mathcal{R}})$  satisfying  $0 = \nabla_y \mathcal{L}(\bar{x}, \bar{x}, \bar{\lambda})$ , a fact we will return to after a few  
 903 calculations.

904 We now compute the first-order variations of  $f$  and  $G$ . To simplify notation, we  
 905 adopt two conventions. First, we align the notation of gradients and Jacobians, viewing  
 906 every gradient as a row vector. Second, we let the symbol  $\nabla^2 F[x; v]$  denote the  $m \times d$   
 907 matrix whose  $i$ th row equals  $v^{\top} \nabla^2 F_i(x)$ . Then, defining the map

$$908 \quad \zeta(x, y) = F(x) + \nabla F(x)(y - x),$$

909 a quick computation shows

$$910 \quad \nabla_y \zeta(x, y) = \nabla F(x) \quad \text{and} \quad \nabla_x \zeta(x, y) = \nabla^2 F[x, y - x].$$

911 Therefore, using the chain rule, we compute the first-order variations

$$\begin{aligned}
 912 \quad \nabla_x G^{\mathcal{H}}(x, y) &= \nabla \omega(\zeta(x, y)) \cdot \nabla^2 F[x, y - x] \\
 913 \quad \nabla_y G^{\mathcal{H}}(x, y) &= \nabla \omega(\zeta(x, y)) \cdot \nabla F(x) \\
 914 \quad \nabla_x G^{\mathcal{R}}(x, y) &= 0 \\
 915 \quad \nabla_y G^{\mathcal{R}}(x, y) &= \nabla \eta(y) \\
 916 \quad \nabla_x f(x, y) &= \nabla \hat{h}(\zeta(x, y)) \cdot \nabla^2 F[x, y - x] + \mu^{-1}(x - y)^\top \\
 917 \quad \nabla_y f(x, y) &= \nabla \hat{h}(\zeta(x, y)) \cdot \nabla F(x) + \nabla \hat{r}(y) + \mu^{-1}(y - x)^\top.
 \end{aligned}$$

919 From these variations we deduce  $\nabla_x G(\bar{x}, \bar{x}) = 0$  and therefore the constraint in (5.7)  
 920 simply amounts to the inclusion

$$\begin{aligned}
 921 \quad v \in \ker \nabla_y G(\bar{x}, \bar{x}) &= \left( \ker \nabla \eta(\bar{x}) \right) \cap \left( \ker(\nabla \omega(F(\bar{x})) \cdot \nabla F(\bar{x})) \right) \\
 &= T_{\mathcal{R}}(\bar{x}) \cap \nabla F(\bar{x})^{-1} T_{\mathcal{H}}(F(\bar{x})) = T_{\mathcal{M}}(\bar{x}).
 \end{aligned} \tag{5.9}$$

922 In particular, formula (5.7) reduces to

$$923 \quad \nabla S(\bar{x})h = \operatorname{argmin}_{v \in T_{\mathcal{M}}(\bar{x})} 2\langle H_{xy}v, h \rangle + \langle H_{yy}v, v \rangle, \tag{5.10}$$

924 To find an explicit solution, we mirror the analysis of the proximal gradient method.  
 925 We let  $W$  be the orthogonal projection onto  $T_{\mathcal{M}}(\bar{x})$  and define the linear maps  
 926  $\overline{H}_{yy}: T_{\mathcal{M}}(\bar{x}) \rightarrow T_{\mathcal{M}}(\bar{x})$  and  $\overline{H}_{xy}: T_{\mathcal{M}}(\bar{x}) \rightarrow T_{\mathcal{M}}(\bar{x})$  by setting  $\overline{H}_{yy} = WH_{yy}W$   
 927 and  $\overline{H}_{xy} = WH_{xy}W$ , respectively. Since  $\bar{x}$  is a strong local minimizer of (5.7), the  
 928 map  $\overline{H}_{yy}$  is positive definite and invertible. Solving (5.7) then yields the expression

$$929 \quad \nabla S(\bar{x})h = -\overline{H}_{yy}^{-1} \overline{H}_{xy}^\top h \quad \text{for all } h \in T_{\mathcal{M}}(\bar{x}).$$

930 Let us now verify that if  $\bar{x}$  is a composite strict saddle of  $f$ , then  $\nabla S(\bar{x})$  has a real  
 931 eigenvalue that is greater than one. To this end, observe that  $\gamma \in \mathbb{R}$  is an eigenvalue  
 932 of  $\nabla S(\bar{x})$  with an associated eigenvector  $v \in T_{\mathcal{M}}(\bar{x})$  if and only if

$$933 \quad \nabla S(\bar{x})v = \gamma v \iff -\overline{H}_{yy}^{-1} \overline{H}_{xy}^\top v = \gamma v \iff (\gamma \overline{H}_{yy} + \overline{H}_{xy}^\top)v = 0.$$

934 In particular, if the matrix  $\gamma \overline{H}_{yy} + \overline{H}_{xy}^\top$  is singular, then  $\gamma$  is an eigenvalue of  $\nabla S(\bar{x})$ .  
 935 To prove such a  $\gamma \geq 1$  exists, we will show that  $\overline{H}_{xy}$  is self-adjoint, and then, we will  
 936 examine the following ray of symmetric matrices

$$937 \quad \{\gamma \overline{H}_{yy} + \overline{H}_{xy}^\top : \gamma \geq 1\}.$$

938 Beginning with the end point, we will show that the matrix  $\overline{H}_{yy} + \overline{H}_{xy}^\top$  has a strictly  
 939 negative eigenvalue. On the other hand, we already know the direction of the ray

940  $\overline{H}_{yy}$  is a positive definite matrix. Therefore, by continuity of the minimal eigenvalue  
 941 function, there will exist some  $\gamma > 1$  such that the matrix  $\gamma \overline{H}_{yy} + \overline{H}_{xy}$  is singular, as  
 942 claimed.

943 To this end, we now compute the second-order variations.

$$\begin{aligned}
 944 \quad \nabla_{xy} G_i^{\mathcal{H}}(x, y)v &= \nabla^2 F[x; v]^{\top} \nabla \omega_i(\zeta(x, y))^{\top} \\
 &\quad + \nabla^2 F[x; y-x]^{\top} \nabla^2 \omega_i(\zeta(x, y)) \nabla F(x)v \\
 945 \quad \nabla_{yy} G_i^{\mathcal{H}}(x, y)v &= \nabla F(x)^{\top} \nabla^2 \omega_i(\zeta(x, y)) \nabla F(x)v \\
 946 \quad \nabla_{xy} f(x, y)v &= \nabla^2 F[x; v]^{\top} \nabla \hat{h}(\zeta(x, y))^{\top} \\
 &\quad + \nabla^2 F[x; y-x] \nabla^2 \hat{h}(\zeta(x, y)) \nabla F(x)v - \mu^{-1}v \\
 947 \quad \nabla_{yy} f(x, y)v &= \nabla F(x)^{\top} \nabla^2 \hat{h}(\zeta(x, y)) \nabla F(x)v + \nabla^2 \hat{r}(y)v + \mu^{-1}v.
 \end{aligned}$$

951 A quick computation then shows that  $\nabla_{xy} f(\bar{x}, \bar{x})$  and  $\nabla_{xy} G_i^{\mathcal{H}}(\bar{x}, \bar{x})$  are self-adjoint  
 952 operators. Consequently, we obtain  $H_{xy} = H_{xy}^{\top}$  and the expression

$$\begin{aligned}
 953 \quad (H_{yy} + H_{xy}^{\top})v &= \nabla F(\bar{x})^{\top} \nabla^2 \hat{h}(F(\bar{x})) \nabla F(\bar{x})v + \nabla^2 \hat{r}(\bar{x})v + \nabla^2 F[\bar{x}; v]^{\top} \nabla \hat{h}(F(\bar{x}))^{\top} \\
 &\quad + \sum_{i \geq 1} \bar{\lambda}_i^{\mathcal{H}} \left( \nabla F(\bar{x})^{\top} \nabla^2 \omega_i(F(\bar{x})) \nabla F(\bar{x})v + \nabla^2 F[\bar{x}; v]^{\top} \nabla \omega_i(F(\bar{x}))^{\top} \right) \\
 954 \quad &\quad + \sum_{i \geq 1} \bar{\lambda}_i^{\mathcal{R}} \nabla^2 \eta_i(y)v.
 \end{aligned}$$

957 To prove that  $H_{yy} + H_{xy}^{\top}$  has a strictly negative eigenvalue, we will show that it  
 958 coincides with the Hessian of the Lagrangian of the problem:

$$959 \quad \min_x \hat{h}(F(x)) + \hat{r}(x) \quad \text{subject to} \quad \omega(F(x)) = 0, \eta(x) = 0.$$

960 Indeed, define the Lagrangian function

$$961 \quad \mathcal{L}_0(x, \lambda) = \hat{h}(F(x)) + \hat{r}(x) + \sum_{i \geq 1} \lambda_i^{\mathcal{H}} \omega(F(x)) + \sum_{i \geq 1} \lambda_i^{\mathcal{R}} \eta(x).$$

962 A quick computation shows

$$\begin{aligned}
 963 \quad \nabla^2(\hat{h} \circ F)(x)v &= \nabla F(x)^{\top} \nabla^2 \hat{h}(F(x)) \nabla F(x)v + \nabla^2 F[x, v]^{\top} \nabla \hat{h}(F(x))^{\top} \\
 964 \quad \nabla^2(\omega_i \circ F)(x)v &= \nabla F(x)^{\top} \nabla^2 \omega_i(F(x)) \nabla F(x)v + \nabla^2 F[x, v]^{\top} \nabla \omega_i(F(x))^{\top}
 \end{aligned}$$

965 and therefore the equality

$$967 \quad \nabla^2 \mathcal{L}_0(\bar{x}, \bar{\lambda}) = H_{yy} + H_{xy}^{\top}.$$

968 The composite strict saddle property guarantees that the matrix  $\nabla^2 \mathcal{L}_0(\bar{x}, \bar{\lambda})$  has a  
 969 strictly negative eigenvalue, completing the proof.  $\square$

970 In line with the previous sections, one could ask whether a damped and randomly  
 971 initialized proximal linear method almost surely escapes all composite strict saddle  
 972 points. An immediate obstacle is that the global Lipschitz constant of the proximal  
 973 linear map  $S(\cdot)$  defined in (5.6) seems unclear, and therefore, we are unable to find  
 974 an appropriate damping parameter. Instead we will settle for a local escape guarantee  
 975 supplied by the center stable manifold theorem. We leave it as an intriguing open  
 976 question to obtain global escape guarantees for the damped proximal linear algorithm.

977 A first difficulty in applying the center stable manifold theorem is that the Jacobian  
 978  $\nabla S(\bar{x})$  at the saddle point  $\bar{x}$  may not be invertible. Consequently, we will damp the  
 979 proximal linear method, forcing the update to be a local diffeomorphism. To compute  
 980 an appropriate threshold for the damping parameter, we will need to estimate the  
 981 operator norm of  $\nabla S(\bar{x})$ . This is the content of the following lemma.

982 **Lemma 5.5** (The slope at the critical points). *Consider the composite optimization*  
 983 *problem (5.1) and choose any  $\mu \in (0, (\rho + 2\beta)^{-1})$ . Then, for all points  $x \in \mathbb{R}^d$  and*  
 984 *all critical points  $\bar{x} \in \mathbb{R}^d$ , the proximal linear map  $S(\cdot)$  defined in (5.6) satisfies*

$$985 \quad \|S(x) - \bar{x}\| \leq \left(1 + \sqrt{\frac{2\beta\mu}{1 - \mu\beta - \mu\rho}}\right) \cdot \max\left\{1, \frac{\mu\rho + \mu\beta}{1 - \mu\rho - 2\mu\beta}\right\} \cdot \|x - \bar{x}\|.$$

986 **Proof** To simplify notation, define the map

$$987 \quad \zeta(x, y) = F(x) + \nabla F(x)(y - x).$$

988 Set  $\gamma := \mu^{-1} - \beta$ , fix an arbitrary point  $x \in \mathbb{R}^d$ , and define

$$989 \quad x^+ := S(x) \quad \text{and} \quad \hat{x} := \text{prox}_{f/\gamma}(x).$$

991 Using strong convexity of the prox-linear and proximal subproblems and the estimate  
 992 (5.2), we successively compute

$$\begin{aligned} 993 \quad h(\hat{x}) + r(\hat{x}) + \frac{\gamma}{2} \|\hat{x} - x\|^2 &\leq h(x^+) + r(x^+) + \frac{\gamma}{2} \|x^+ - x\|^2 - \frac{\gamma - \rho - \beta}{2} \|x^+ - \hat{x}\|^2 \\ 994 \quad &\leq h(\zeta(x, x^+)) + r(x^+) + \frac{\gamma + \beta}{2} \|x^+ - x\|^2 - \frac{\gamma - \rho - \beta}{2} \|x^+ - \hat{x}\|^2 \\ 995 \quad &\leq h(\zeta(x, \hat{x})) + r(\hat{x}) + \frac{\gamma + \beta}{2} \|\hat{x} - x\|^2 - (\gamma - \rho) \|x^+ - \hat{x}\|^2 \\ 996 \quad &\leq h(\hat{x}) + r(\hat{x}) + \frac{\gamma + 2\beta}{2} \|\hat{x} - x\|^2 - (\gamma - \rho) \|x^+ - \hat{x}\|^2. \end{aligned}$$

998 Rearranging yields the estimate

$$999 \quad (\gamma - \rho) \|x^+ - \hat{x}\|^2 \leq 2\beta \|\hat{x} - x\|^2 = 2\beta\gamma^{-2} \|\nabla f_{1/\gamma}(x)\|^2.$$



1001 Therefore, using Lipschitz continuity of the gradient  $\nabla f_{1/\gamma}$  (Lemma 2.5) and the  
 1002 triangle inequality yields

$$\begin{aligned}
 1003 \quad \|x^+ - \bar{x}\| &\leq \left( \gamma^{-1} + \sqrt{\frac{2\beta\gamma^{-2}}{\gamma - \rho}} \right) \cdot \max \left\{ \gamma, \frac{\rho + \beta}{1 - \gamma^{-1}(\rho + \beta)} \right\} \cdot \|x - \bar{x}\| \\
 1004 \quad &= \left( 1 + \sqrt{\frac{2\beta\mu}{1 - \mu\beta - \mu\rho}} \right) \cdot \max \left\{ 1, \frac{\mu\rho + \mu\beta}{1 - \mu\rho - 2\mu\beta} \right\} \cdot \|x - \bar{x}\|, \\
 1005
 \end{aligned}$$

1006 as claimed.  $\square$

1007 We are now ready to deduce that the damped proximal linear method almost locally  
 1008 escapes any composite strict saddle point.

1009 **Theorem 5.6** (Proximal linear method: local escape). *Consider the composite prob-*  
 1010 *lem (5.1) and let  $\bar{x}$  be any composite strict saddle point. Choose any constant*  
 1011  *$\mu \in (0, (\rho + 2\beta)^{-1})$  and a damping parameter  $\alpha \in (0, 1)$  satisfying*

$$1012 \quad \alpha \cdot \left( 1 + \left( \left( 1 + \sqrt{\frac{2\beta\mu}{1 - \mu\beta - \mu\rho}} \right) \cdot \max \left\{ 1, \frac{\mu\rho + \mu\beta}{1 - \mu\rho - 2\mu\beta} \right\} \right) \right) < 1.$$

1013 Define the damped proximal linear update

$$1014 \quad T(x) = (1 - \alpha)x + \alpha S(x),$$

1015 where  $S(\cdot)$  is the proximal linear map defined in (5.6). Then, there exists a neighbor-  
 1016 hood  $U$  of  $\bar{x}$  such that the set of initial conditions

$$1017 \quad \{x \in U : S^k(x) \in U \text{ for all } k \geq 0\}$$

1018 has zero Lebesgue measure.

1019 **Proof** First, using Theorem 5.3 and Lemma 2.14, we deduce that  $\bar{x}$  is an unstable fixed  
 1020 point of  $\bar{x}$ . Let us next verify that  $T$  is a local diffeomorphism around  $\bar{x}$ . To see this,  
 1021 observe

$$1022 \quad \nabla T(\bar{x}) = I - \alpha(I - \nabla S(\bar{x})).$$

1023 Using Theorem 5.5, we deduce  $\alpha\|I - \nabla S(\bar{x})\|_{\text{op}} < 1$  and therefore  $T$  is invertible.  
 1024 An application of the center stable manifold theorem (Theorem 2.10) completes the  
 1025 proof.  $\square$

1026 **6 Convergence of Relaxed Descent Methods**

1027 Thus far, all of our escape theorems made an assumption that the iterate sequence  
 1028 generated by the algorithms converges. In this section, we verify this assumption for  
 1029 the damped proximal point, proximal gradient, and proximal linear methods. Taking a  
 1030 general view, we see that the iterative methods of this paper can be understood within  
 1031 a broad family of damped model-based algorithms for minimizing a function  $f$ . These  
 1032 algorithms construct iterates  $x_0, x_1 \dots$  by repeatedly minimizing a local model  $f_x(\cdot)$   
 1033 of the function and moving in the direction of its minimizer. More specifically, in the  
 1034 section we suppose that there exist constant  $\rho, \eta, \beta > 0$  such that the the following  
 1035 properties hold:

- 1036 (A1) The function  $f: \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$  is closed and  $\rho$ -weakly convex.  
 1037 (A2) For all  $x \in \mathbb{R}^d$  there exists a closed  $\eta$ -weakly convex function  $f_x: \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$   
 1038 satisfying

1039 
$$|f(y) - f_x(y)| \leq \frac{\beta}{2} \|y - x\|^2 \quad \text{for all } y \in \mathbb{R}^d.$$

1040 Under these assumptions, we will study how the following algorithm behaves: given  
 1041 iterates  $x_0, \dots, x_t$  define

1042 
$$y_t = \operatorname{argmin}_{y \in \mathbb{R}^d} \left\{ f_{x_t}(y) + \frac{\tau}{2} \|y - x_t\|^2 \right\} \tag{MBA}$$
  
 1043 
$$x_{t+1} = (1 - \alpha)x_t + \alpha y_t,$$

1044 where  $\tau > 0$  and  $\alpha > 0$  are fixed constants, determined below.

1045 To analyze this algorithm, we rely on the seminal paper [3]. There, the authors  
 1046 identified three conditions, guaranteeing global convergence of a sequence  $\{z_t\}$  of  
 1047 “algorithm iterates” to a critical point of a closed function  $g: \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$ . Namely,  
 1048 they assume there exist  $a, b > 0$  such that the following holds:

- 1049 (B1) **(Sufficient Decrease.)** For each  $t \in \mathbb{N}$ , we have

1050 
$$g(z_{t+1}) + a \|z_{t+1} - z_t\|^2 \leq g(z_t)$$

- 1051 (B2) **(Relative Error Conditions.)** For each  $t \in \mathbb{N}$  there exists  $w_{t+1} \in \partial g(z_{t+1})$  such  
 1052 that

1053 
$$\|w_{t+1}\| \leq b \|z_{t+1} - z_t\|$$

- 1054 (B3) **(Continuity Condition.)** There exists a subsequence  $\{z_{t_j}\}$  and  $\tilde{z}$  such that

1055 
$$z_{t_j} \rightarrow \tilde{z} \text{ and } g(z_{t_j}) \rightarrow g(\tilde{z}), \quad \text{as } j \rightarrow \infty.$$

1056 The above assumptions alone may not guarantee that  $z_t$  converges to a critical  
 1057 point of  $g$ . Instead, the authors of [3] restrict their focus to the broad class of functions  
 1058 satisfying the *Kurdyka–Lojasiewicz property*.

1059 **Definition 6.1** (KŁ Function) Let  $g: \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$  be a closed function. We say that  
 1060  $g$  has the Kurdyka–Lojasiewicz (KL) property at a point  $\bar{x}$ , where  $\partial g(\bar{x})$  is nonempty,  
 1061 if there exists  $\varepsilon \in (0, +\infty]$ , a neighborhood  $U$  of  $\bar{x}$ , and a continuous convex function  
 1062  $\varphi: [0, \varepsilon) \rightarrow \mathbb{R}_+$  satisfying

- 1063 1.  $\varphi(0) = 0$ ,
- 1064 2.  $\varphi$  is  $C^1$  on  $(0, \varepsilon)$  with  $\varphi' > 0$ , and
- 1065 3. the KŁ inequality

$$1066 \quad \text{dist}(0, \partial g(x)) \geq \frac{1}{\varphi'(g(x) - g(\bar{x}))},$$

1067 holds for all  $x \in U$  satisfying  $g(\bar{x}) < g(x) < g(\bar{x}) + \varepsilon$ .

1068 If  $g$  satisfies the KŁ property at each point  $x$ , with  $\partial g(x) \neq \emptyset$ , then  $g$  is called a *KL*  
 1069 *function*.

1070 The class of KŁ functions is broad, containing all closed semialgebraic functions  
 1071 and more broadly any functions definable in an o-minimal structure, as shown in the  
 1072 pioneering work [7]. Under these assumptions we have the following theorem from [3,  
 1073 Theorem 2.9].

1074 **Theorem 6.2** Let  $g: \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$  be a closed function. Consider a sequence  $x_t$  that  
 1075 satisfies (B1), (B2), and (B3). If  $g$  satisfies the KŁ property at some cluster point  $\tilde{x}$ ,  
 1076 then  $\tilde{x}$  is a critical point of  $g$ , the entire sequence  $x_k$  converges to  $\tilde{x}$ , and the sequence  
 1077  $x_t$  has finite length

$$1078 \quad \sum_{t=0}^{\infty} \|x_{t+1} - x_t\| < +\infty.$$

1079 In the remainder of this section, we will verify assumption (B1), (B2), and (B3)  
 1080 for the sequence  $\{z_t\} = \{x_t\}$  and the Moreau envelope  $g := f_{1/\hat{\rho}}$ , where  $\hat{\rho}$  will be  
 1081 chosen in a moment. Since the critical points of  $f$  and  $f_{1/\hat{\rho}}$  agree, the result will imply  
 1082 convergence to critical points of  $f$ . To do so, we employ one final assumption.

1083 (A3) For every  $\hat{\rho} > 0$ , the Moreau envelope  $f_{1/\hat{\rho}}$  is a KŁ function.

1084 Although assumption (A3) may appear hard to verify, it holds whenever  $f$  is semi-  
 1085 algebraic since in this case  $f_{1/\hat{\rho}}$  is also semialgebraic. More generally, the analogous  
 1086 statement holds if  $f$  is definable in an o-minimal structure. The following is the main  
 1087 result of this section.

1088 **Theorem 6.3** (Convergence of relaxed model-based methods). Suppose that  $\alpha \in$   
 1089  $(0, 1]$ , that  $\tau > \max\{\eta, 2\rho, \frac{4\beta + \rho + \eta}{2}\}$ , and that assumptions (A1) and (A2) hold. Then,  
 1090 for all  $T \geq 0$ , we have

$$1091 \quad \min_{t=0, \dots, T} \|\nabla f_{1/\hat{\rho}}(x_t)\| \leq \sqrt{\frac{f_{1/\hat{\rho}}(x_0) - \inf f}{\frac{\alpha(2\hat{\rho} - \rho - \eta - \beta)}{2\hat{\rho}(\hat{\rho} + \tau - \rho - \eta)}(T + 1)}}.$$

1092 where  $\hat{\rho} = (1/2)\tau + (1/4)(\rho + \eta)$ . Moreover, if (A3) also holds and the sequence  $\{x_t\}$   
 1093 has a cluster point  $\bar{x}$ , then  $\bar{x}$  is critical for  $f$  and the entire sequence  $\{x_t\}$  converges  
 1094 to  $\bar{x}$ . Moreover, the sequence  $\{x_t\}$  has finite length.

1095 
$$\sum_{t=0}^{\infty} \|x_{t+1} - x_t\| < +\infty.$$

1096 This result is new and may be of independent interest. In particular, the conclusion  
 1097 of the theorem extends the convergence guarantees for the proximal linear method  
 1098 developed in [52] to all relaxed model-based algorithms.

1099 **6.1 Proof of Theorem 6.3**

1100 We are free to choose the parameter  $\hat{\rho}$  defining the Moreau envelope. To this end, we  
 1101 will need the existence of a parameter  $\hat{\rho}$ , satisfying the following inequalities.

1102 **Lemma 6.4** *Under the assumptions of Theorem 6.3, it holds that  $\hat{\rho} > \rho$  and*

- 1103 1.  $\tau - \hat{\rho} - \beta > 0,$   
 1104 2.  $2\hat{\rho} - \rho - \eta - \beta > 0,$   
 1105 3.  $\hat{\rho} + \tau - \rho - \eta > 0,$   
 1106 4.  $1 - \frac{2\hat{\rho} - \rho - \eta - \beta}{\hat{\rho} + \tau - \rho - \eta} > 0.$

1107 **Proof** Note that  $\hat{\rho} > \tau/2 > \rho > 0$  and that  $\hat{\rho} = \tau - \beta - \varepsilon/2$  for  $\varepsilon = (2\tau - 4\beta - \rho -$   
 1108  $\eta)/2 > 0$ . To prove the first inequality, notice that  $\tau - \hat{\rho} - \beta = \varepsilon/2 > 0$ . To prove  
 1109 the second inequality, notice that

1110 
$$2\hat{\rho} - \rho - \eta - \beta = 2\tau - 4\beta - \rho - \eta - \varepsilon = \varepsilon > 0.$$

1111 To prove the third inequality, observe

1112 
$$\hat{\rho} + \tau - \rho - \eta > \hat{\rho} + (\hat{\rho} + \beta) - \rho - \eta \geq 2\hat{\rho} - \rho - \eta - \beta > 0,$$

1113 where the first and second inequalities follow from items 6.4 and 6.4, respectively. To  
 1114 prove the fourth inequality, we compute

1115 
$$1 - \frac{2\hat{\rho} - \rho - \eta - \beta}{\hat{\rho} + \tau - \rho - \eta} = \frac{\beta + \tau - \hat{\rho}}{\hat{\rho} + \tau - \rho - \eta} = \frac{2\beta + \varepsilon/2}{\hat{\rho} + \tau - \rho - \eta} > 0,$$

1116 as desired. □

1117 Throughout the rest of this section, we fix a constant  $\hat{\rho}$  satisfying the conditions  
 1118 of Lemma 6.4. Critical to our proof is the following lemma, comparing the proximal  
 1119 point

1120 
$$\hat{x}_t := \text{prox}_{f/\hat{\rho}}(x_t)$$

1121 to the “approximately proximal point”  $y_t$ . A closely related estimate appeared in [16,  
1122 Lemma 4.2], driving the convergence analysis of that paper.

1123 **Lemma 6.5** *It holds that*

$$1124 \quad \|\hat{x}_t - y_t\|^2 \leq \|\hat{x}_t - x_t\|^2 - \frac{2\hat{\rho} - \rho - \eta - \beta}{\hat{\rho} + \tau - \rho - \eta} \|\hat{x}_t - x_t\|^2 - \frac{\tau - \hat{\rho} - \beta}{\hat{\rho} + \tau - \rho - \eta} \|x_t - y_t\|^2.$$

1125 **Proof** Since the function  $y \mapsto f(y) + \frac{\hat{\rho}}{2} \|y - x_t\|^2$  is  $(\hat{\rho} - \rho)$ -strongly convex and  $\hat{x}_t$   
1126 is its minimizer, we have

$$1127 \quad \frac{\hat{\rho} - \rho}{2} \|\hat{x}_t - y_t\|^2 \leq \left( f(y_t) + \frac{\hat{\rho}}{2} \|y_t - x_t\|^2 \right) - \left( f(\hat{x}_t) + \frac{\hat{\rho}}{2} \|\hat{x}_t - x_t\|^2 \right).$$

1129 Consequently, using the double-sided model property (A2), we find

$$1130 \quad \frac{\hat{\rho} - \rho}{2} \|\hat{x}_t - y_t\|^2 \leq f_{x_t}(y_t) - f_{x_t}(\hat{x}_t) + \frac{\hat{\rho} + \beta}{2} \|x_t - y_t\|^2 - \frac{\hat{\rho} - \beta}{2} \|\hat{x}_t - x_t\|^2. \quad (6.1)$$

1131 Since the function  $y \mapsto f_{x_t}(y) + \frac{\tau}{2} \|y - x_t\|^2$  is  $(\tau - \eta)$ -strongly convex and  $y_t$  is its  
1132 minimizer, we have

$$1133 \quad f_{x_t}(y_t) - f_{x_t}(\hat{x}_t) \leq \frac{\tau}{2} \|\hat{x}_t - x_t\|^2 - \frac{\tau}{2} \|y_t - x_t\|^2 - \frac{\tau - \eta}{2} \|y_t - \hat{x}_t\|^2.$$

1134 Combining this estimate with (6.1), we compute

$$1135 \quad \begin{aligned} \frac{\hat{\rho} - \rho}{2} \|\hat{x}_t - y_t\|^2 &\leq \frac{\tau}{2} \|\hat{x}_t - x_t\|^2 - \frac{\tau}{2} \|y_t - x_t\|^2 - \frac{\tau - \eta}{2} \|y_t - \hat{x}_t\|^2 \\ 1136 \quad &\quad + \frac{\hat{\rho} + \beta}{2} \|x_t - y_t\|^2 - \frac{\hat{\rho} - \beta}{2} \|\hat{x}_t - x_t\|^2 \\ 1137 \quad &= \frac{\beta + \tau - \hat{\rho}}{2} \|\hat{x}_t - x_t\|^2 + \frac{\hat{\rho} + \beta - \tau}{2} \|x_t - y_t\|^2 - \frac{\tau - \eta}{2} \|y_t - \hat{x}_t\|^2. \end{aligned}$$

1139 Rearranging, we conclude

$$1140 \quad \frac{\hat{\rho} + \tau - \rho - \eta}{2} \|\hat{x}_t - y_t\|^2 \leq \frac{\beta + \tau - \hat{\rho}}{2} \|\hat{x}_t - x_t\|^2 + \frac{\hat{\rho} + \beta - \tau}{2} \|x_t - y_t\|^2.$$

1141 Dividing both sides by  $\frac{\hat{\rho} + \tau - \rho - \eta}{2}$ , we achieve the result:

$$1142 \quad \begin{aligned} \|\hat{x}_t - y_t\|^2 &\leq \frac{\beta + \tau - \hat{\rho}}{\hat{\rho} + \tau - \rho - \eta} \|\hat{x}_t - x_t\|^2 + \frac{\hat{\rho} + \beta - \tau}{\hat{\rho} + \tau - \rho - \eta} \|x_t - y_t\|^2 \\ 1143 \quad &= \|\hat{x}_t - x_t\|^2 - \left( 1 - \frac{\beta + \tau - \hat{\rho}}{\hat{\rho} + \tau - \rho - \eta} \right) \|\hat{x}_t - x_t\|^2 + \frac{\hat{\rho} + \beta - \tau}{\hat{\rho} + \tau - \rho - \eta} \|x_t - y_t\|^2 \\ 1144 \quad &= \|\hat{x}_t - x_t\|^2 - \frac{2\hat{\rho} - \rho - \eta - \beta}{\hat{\rho} + \tau - \rho - \eta} \|\hat{x}_t - x_t\|^2 - \frac{\tau - \hat{\rho} - \beta}{\hat{\rho} + \tau - \rho - \eta} \|x_t - y_t\|^2. \end{aligned}$$

1146 This completes the proof of the lemma. □

1147 The following lemma verifies the Assumption (B1).

1148 **Lemma 6.6** (Sufficient Decrease) *We have*

1149 
$$f_{1/\hat{\rho}}(x_{t+1}) \leq f_{1/\hat{\rho}}(x_t) - \frac{\hat{\rho}(\tau - \hat{\rho} - \beta)}{2\alpha(\hat{\rho} + \tau - \rho - \eta)} \|x_{t+1} - x_t\|^2$$

1150 
$$- \frac{\alpha(2\hat{\rho} - \rho - \eta - \beta)}{2\hat{\rho}(\hat{\rho} + \tau - \rho - \eta)} \|\nabla f_{1/\hat{\rho}}(x_t)\|^2.$$

1151 *In particular,  $f_{1/\hat{\rho}}$  and  $\{x_t\}$  satisfy (B1). Moreover, for all  $T \geq 0$ , we have*

1152 
$$\min_{t=0, \dots, T} \|\nabla f_{1/\hat{\rho}}(x_t)\|^2 \leq \frac{1}{T+1} \sum_{t=0}^T \|\nabla f_{1/\hat{\rho}}(x_t)\|^2 \leq \frac{f_{1/\hat{\rho}}(x_0) - \inf f}{\frac{\alpha(2\hat{\rho} - \rho - \eta - \beta)}{2\hat{\rho}(\hat{\rho} + \tau - \rho - \eta)}(T+1)}$$

1153 **Proof** We successively compute

1154 
$$f_{1/\hat{\rho}}(x_{t+1}) = f(\hat{x}_{t+1}) + \frac{\hat{\rho}}{2} \|\hat{x}_{t+1} - x_{t+1}\|^2$$

1155 
$$\leq f(\hat{x}_t) + \frac{\hat{\rho}}{2} \|\hat{x}_t - x_{t+1}\|^2$$

1156 
$$= f(\hat{x}_t) + \frac{\hat{\rho}}{2} \|(1 - \alpha)(\hat{x}_t - x_t) + \alpha(\hat{x}_t - y_t)\|^2$$

1157 
$$\leq f(\hat{x}_t) + \frac{\hat{\rho}(1 - \alpha)}{2} \|\hat{x}_t - x_t\|^2 + \frac{\hat{\rho}\alpha}{2} \|\hat{x}_t - y_t\|^2$$

1158 
$$\leq f(\hat{x}_t) + \frac{\hat{\rho}}{2} \|\hat{x}_t - x_t\|^2$$

1159 
$$- \frac{\hat{\rho}\alpha}{2} \left( \frac{2\hat{\rho} - \rho - \eta - \beta}{\hat{\rho} + \tau - \rho - \eta} \|\hat{x}_t - x_t\|^2 + \frac{\tau - \hat{\rho} - \beta}{\hat{\rho} + \tau - \rho - \eta} \|x_t - y_t\|^2 \right)$$

1160 
$$\leq f_{1/\hat{\rho}}(x_t) - \frac{\hat{\rho}\alpha(\tau - \hat{\rho} - \beta)}{2(\hat{\rho} + \tau - \rho - \eta)} \|x_t - y_t\|^2 - \frac{\alpha(2\hat{\rho} - \rho - \eta - \beta)}{2\hat{\rho}(\hat{\rho} + \tau - \rho - \eta)} \|\nabla f_{1/\hat{\rho}}(x_t)\|^2,$$

1161 (6.2)

1162 where (6.2) follows from Lemma 6.5, and the final inequality follows since  $\hat{\rho}(x_t - \hat{x}_t) =$   
 1163  $\nabla f_{1/\hat{\rho}}(x_t)$ . To get the descent inequality, it remains to write  $x_t - y_t = (x_{t+1} - x_t)/\alpha$ .  
 1164 Finally, the bound on the average gradient norm follows by induction.  $\square$

1165 The following lemma verifies the Assumption (B2).

1166 **Lemma 6.7** (Relative Error). *It holds*

1167 
$$\|\nabla f_{1/\hat{\rho}}(x_{t+1})\| \leq \left( \max \left\{ \hat{\rho}, \frac{\rho}{1 - \rho/\hat{\rho}} \right\} + \frac{\hat{\rho}}{\alpha} \frac{1}{1 - \sqrt{\left(1 - \frac{2\hat{\rho} - \rho - \eta - \beta}{\hat{\rho} + \tau - \rho - \eta}\right)}} \right) \|x_{t+1} - x_t\|.$$

1168 *In particular,  $f_{1/\hat{\rho}}$  and  $\{x_t\}$  satisfy (B2).*

Author Proof

1169 **Proof** We have

$$1170 \quad \|\nabla f_{1/\hat{\rho}}(x_{t+1})\| \leq \|\nabla f_{1/\hat{\rho}}(x_t)\| + \max\left\{\hat{\rho}, \frac{\rho}{1-\rho/\hat{\rho}}\right\} \|x_{t+1} - x_t\|.$$

1171 Thus, we want to bound

$$1172 \quad \|\nabla f_{1/\hat{\rho}}(x_t)\| = \hat{\rho} \|\hat{x}_t - x_t\|$$

1173 by a multiple of  $\|x_{t+1} - x_t\|$ . This follows by Lemma 6.5:

$$1174 \quad \|\hat{x}_t - x_t\| \leq \|\hat{x}_t - y_t\| + \|y_t - x_t\| \leq \sqrt{\left(1 - \frac{2\hat{\rho} - \rho - \eta - \beta}{\hat{\rho} + \tau - \rho - \eta}\right)} \|x_t - \hat{x}_t\| + \|y_t - x_t\|$$

1175 Rearranging and using the definition  $x_t - y_t = (x_{t+1} - x_t)/\alpha$ , it holds

$$1176 \quad \begin{aligned} \|\hat{x}_t - x_t\| &\leq \frac{1}{1 - \sqrt{\left(1 - \frac{2\hat{\rho} - \rho - \eta - \beta}{\hat{\rho} + \tau - \rho - \eta}\right)}} \|y_t - x_t\| \\ 1177 \quad &= \frac{1}{\alpha} \frac{1}{1 - \sqrt{\left(1 - \frac{2\hat{\rho} - \rho - \eta - \beta}{\hat{\rho} + \tau - \rho - \eta}\right)}} \|x_{t+1} - x_t\|. \end{aligned}$$

1178 The proof is complete. as desired. □

1179 Finally, we can dispense with Assumption (B3), which is a simple consequence of  
1180 the continuity of  $f_{\hat{\rho}}$ .

1181 **Lemma 6.8** (Continuity Condition). *The function  $f_{\hat{\rho}}$  and the sequence  $\{x_t\}$  sat-*  
1182 *isfy (B3).*

1183 **Acknowledgements** We thank John Duchi for his insightful comments on an early version of the  
1184 manuscript. We also thank the anonymous referees for numerous suggestions that have improved the read-  
1185 ability of the paper.

## 1186 A Proofs of Theorems 2.9 and 5.2

1187 In this section, we prove Theorem 2.9. We should note that Theorem 2.9, appropriately  
1188 restated, holds much more broadly beyond the weakly convex function class. To sim-  
1189 plify the notational overhead, however, we impose the weak convexity assumption,  
1190 throughout.

1191 We will require some basic notation from variational analysis; for details, we refer  
1192 the reader to [57]. A set-valued map  $F: \mathbb{R}^d \rightrightarrows \mathbb{R}^m$  assigns to each point  $x \in \mathbb{R}^d$  a set  
1193  $F(x)$  in  $\mathbb{R}^m$ . The graph of  $F$  is defined by

$$1194 \quad \text{gph } F := \{(x, v) : v \in F(x)\}.$$

1195 A map  $F: \mathbb{R}^d \rightrightarrows \mathbb{R}^m$  is called *metrically regular at*  $(\bar{x}, \bar{v}) \in \text{gph } F$  if there exists a  
 1196 constant  $\kappa > 0$  such that the estimate holds:

1197 
$$\text{dist}(x, F^{-1}(v)) \leq \kappa \text{dist}(v, F(x))$$

1198 for all  $x$  near  $\bar{x}$  and all  $v$  near  $\bar{v}$ . If the graph  $\text{gph } F$  is a  $C^1$ -smooth manifold around  
 1199  $(\bar{x}, \bar{v})$ , then metric regularity at  $(\bar{x}, \bar{v})$  is equivalent to the condition [57, Theorem  
 1200 9.43(d)]:<sup>12</sup>

1201 
$$(0, u) \in N_{\text{gph } F}(\bar{x}, \bar{v}) \implies u = 0. \tag{A.1}$$

1202 We begin with the following lemma.

1203 **Lemma A.1** (Subdifferential metric regularity in smooth minimization). *Consider the*  
 1204 *optimization problem*

1205 
$$\min_{x \in \mathbb{R}^d} f(x) \text{ subject to } x \in \mathcal{M},$$

1206 where  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  is a  $C^2$ -smooth function and  $\mathcal{M}$  is a  $C^2$ -smooth manifold. Let  $\bar{x} \in$   
 1207  $\mathcal{M}$  satisfy the criticality condition  $0 \in \partial f_{\mathcal{M}}(\bar{x})$  and suppose that the subdifferential  
 1208 map  $\partial f_{\mathcal{M}}: \mathbb{R}^d \rightrightarrows \mathbb{R}^d$  is metrically regular at  $(\bar{x}, 0)$ . Then, the guarantee holds:

1209 
$$\inf_{u \in \mathbb{S}^{d-1} \cap T_{\mathcal{M}}(\bar{x})} d^2 f_{\mathcal{M}}(\bar{x})(u) \neq 0. \tag{A.2}$$

1210 **Proof** First, appealing to (A.1), we conclude that the implication holds:

1211 
$$(0, u) \in N_{\text{gph } \partial f_{\mathcal{M}}}(\bar{x}, 0) \implies u = 0. \tag{A.3}$$

1212 Let us now interpret the condition (A.3) in Lagrangian terms. To this end, let  $G = 0$   
 1213 be the local defining equations for  $\mathcal{M}$  around  $\bar{x}$ . Define the Lagrangian function

1214 
$$\mathcal{L}(x, \lambda) = f(x) + \langle G(x), \lambda \rangle,$$

1215 and let  $\bar{\lambda}$  be the unique Lagrange multiplier vector satisfying  $\nabla_x \mathcal{L}(\bar{x}, \bar{\lambda}) = 0$ . Accord-  
 1216 ing to [41, Corollary 2.9], we have the following expression:

1217 
$$(0, u) \in N_{\text{gph } \partial f_{\mathcal{M}}}(\bar{x}, 0) \iff u \in T_{\mathcal{M}}(\bar{x}) \text{ and } Lu \in N_{\mathcal{M}}(\bar{x}), \tag{A.4}$$

1218 where  $L := \nabla_{xx}^2 \mathcal{L}(\bar{x}, \bar{\lambda})$  denotes the Hessian of the Lagrangian. Combining (A.3) and  
 1219 (A.4), we deduce that the only vector  $u \in T_{\mathcal{M}}(\bar{x})$  satisfying  $Lu \in N_{\mathcal{M}}(\bar{x})$  is the zero  
 1220 vector  $u = 0$ .

1221 Now for the sake of contradiction, suppose that (A.2) fails. Then, the quadratic  
 1222 form  $Q(u) = \langle Lu, u \rangle$  is nonnegative on  $T_{\mathcal{M}}(\bar{x})$  and there exists  $0 \neq \bar{u} \in T_{\mathcal{M}}(\bar{x})$   
 1223 satisfying  $Q(\bar{u}) = 0$ . We deduce that  $\bar{u}$  minimizes  $Q(\cdot)$  on  $T_{\mathcal{M}}(\bar{x})$ , and therefore, the  
 1224 inclusion  $L\bar{u} \in N_{\mathcal{M}}(\bar{x})$  holds, a clear contradiction.  $\square$

<sup>12</sup> We should note that metric regularity of  $F$  at  $(\bar{x}, \bar{v})$  is equivalent to (A.1) for an arbitrary set-valued map  $F$  with closed graph, provided we interpret  $N_{\text{gph } F}(\bar{x}, \bar{v})$  as the limiting normal cone [57, Definition 6.3].



1225 The following corollary for active manifolds will now quickly follow.

1226 **Corollary A.2** (Subdifferential metric regularity and active manifolds). *Consider a*  
 1227 *closed and weakly convex function  $f: \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$ . Suppose that  $f$  admits a*  
 1228  *$C^2$ -smooth active manifold around a critical point  $\bar{x}$  and that the subdifferential map*  
 1229  *$\partial f: \mathbb{R}^d \rightrightarrows \mathbb{R}^d$  is metrically regular at  $(\bar{x}, 0)$ . Then,  $\bar{x}$  is either a strong local minimizer*  
 1230 *of  $f$  or satisfies the curvature condition  $d^2 f_{\mathcal{M}}(\bar{x})(u) < 0$  for some  $u \in T_{\mathcal{M}}(\bar{x})$ .*

1231 **Proof** The result [19, Proposition 10.2] implies that  $\text{gph } \partial f$  coincides with  $\text{gph } \partial f_{\mathcal{M}}$   
 1232 on a neighborhood of  $(\bar{x}, 0)$ . Therefore, the subdifferential map  $\partial f_{\mathcal{M}}: \mathbb{R}^d \rightrightarrows \mathbb{R}^d$  is  
 1233 metrically regular at  $(\bar{x}, 0)$ . Using Lemma A.1, we obtain the guarantee:

1234 
$$\inf_{u \in \mathbb{S}^{d-1} \cap T_{\mathcal{M}}(\bar{x})} d^2 f_{\mathcal{M}}(\bar{x})(u) \neq 0.$$

1235 If the infimum is strictly negative, the proof is complete. Otherwise, the infimum is  
 1236 strictly positive. In this case,  $\bar{x}$  is a strong local minimizer of  $f_{\mathcal{M}}$ , and therefore by  
 1237 [19, Proposition 7.2] a strong local minimizer of  $f$ .  $\square$

1238 We are now ready for the proofs of Theorems 2.9 and 5.2.

1239 **Proof of Theorem 2.9** The result [18, Corollary 4.8] shows that for almost all  $v \in \mathbb{R}^d$ ,  
 1240 the function  $g(x) := f(x) - \langle v, x \rangle$  has at most finitely many critical points. Moreover  
 1241 each such critical point  $\bar{x}$  lies on some  $C^2$  active manifold  $\mathcal{M}$  of  $g$  and the subdiffer-  
 1242 ential map  $\partial g: \mathbb{R}^d \rightrightarrows \mathbb{R}^d$  is metrically regular at  $(\bar{x}, 0)$ . Applying Corollary A.2 to  $g$   
 1243 for such generic vectors  $v$ , we deduce that every critical point  $\bar{x}$  of  $g$  is either a strong  
 1244 local minimizer or a strict saddle of  $g$ . The proof is complete.  $\square$

1245 **Proof of Theorem 5.2** The proof is identical to that of Theorem 2.9 with [18, Theorem  
 1246 5.2] playing the role of [18, Corollary 4.8].  $\square$

1247 **B Pathological Example**

1248 **Theorem B.1** *Consider the following function*

1249 
$$f(x, y) = \frac{1}{2}(|x| + |y|)^2 - \frac{\rho}{2}x^2$$

1250 *Assume that  $\lambda > \rho$ . Define a mapping  $T: \mathbb{R}^d \rightarrow \mathbb{R}$  by the following formula.*

1251 
$$S(x, y) = \begin{cases} 0 & \text{if } (x, y) = 0; \\ \left(0, \frac{\lambda}{1+\lambda}y\right) & \text{if } |x| \leq \frac{1}{1+\lambda}|y|; \\ \left(\frac{\lambda}{1+\lambda-\rho}x, 0\right) & \text{if } |y| \leq \frac{1}{1+\lambda-\rho}|x|, \end{cases}$$

1252 and if  $\frac{1}{(1+\lambda-\rho)}|x| < |y| < (1+\lambda)|x|$ , we have

Author Proof

$$1253 \quad S(x, y) = \begin{cases} \frac{\lambda}{(1+\lambda)(1+\lambda-\rho)-1} \begin{bmatrix} (1+\lambda) & -1 \\ -1 & (1+\lambda-\rho) \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} & \text{if } \text{sign}(x) = \text{sign}(y); \\ \frac{\lambda}{(1+\lambda)(1+\lambda-\rho)-1} \begin{bmatrix} (1+\lambda) & 1 \\ 1 & (1+\lambda-\rho) \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} & \text{if } \text{sign}(x) \neq \text{sign}(y). \end{cases}$$

1254 Then,  $\text{prox}_{(1/\lambda)f}(x, y) = S(x, y)$ .

1255 **Proof** Let us denote the components of  $S(x, y)$  by  $(x_+, y_+) = S(x, y)$ . By first-order  
1256 optimality conditions, we have  $\text{prox}_{(1/\lambda)f}(x, y) = (x_+, y_+)$  if and only if

$$1257 \quad \lambda(x - (1 - (1/\lambda)\rho)x_+, y - y_+) \in$$

$$1258 \quad \begin{cases} \{x_+ + \text{sign}(x_+)|y_+|\} \times \{\text{sign}(y_+)|x_+| + y_+\} & \text{if } x_+ \neq 0 \text{ and } y_+ \neq 0; \\ \{[-1, 1]y_+\} \times \{y_+\} & \text{if } x_+ = 0 \text{ and } y_+ \neq 0; \\ \{x_+\} \times \{[-1, 1]x_+\} & \text{if } x_+ \neq 0 \text{ and } y_+ = 0; \\ \{0\} \times \{0\} & \text{if } x_+ = 0 \text{ and } y_+ = 0. \end{cases}$$

1259 Let us show that  $(x_+, y_+)$  indeed satisfies this inclusion.

- 1260 1. If  $(x, y) = 0$ , then  $x_+ = y_+ = 0$ , and the pair satisfies the inclusion.  
1261 2. If  $|x| \leq \frac{1}{1+\lambda}|y|$  and  $y \neq 0$ , then  $x_+ = 0$ ,  $y_+ = \frac{\lambda}{1+\lambda}y$ , and

$$1262 \quad \lambda(x - (1 - (1/\lambda)\rho)x_+, y - y_+) = \lambda \left( x, \frac{1}{1+\lambda}y \right) \in \{[-1, 1]y_+\} \times \{y_+\}.$$

1263 Thus, the pair satisfies the inclusion.

- 1264 3. If  $|y| \leq \frac{1}{1+\lambda-\rho}|x|$  and  $x \neq 0$ , then  $x_+ = \frac{\lambda}{(1+\lambda-\rho)}x$ ,  $y_+ = 0$ , and

$$1265 \quad \lambda(x - (1 - (1/\lambda)\rho)x_+, y - y_+)$$

$$1266 \quad = \lambda \left( x - \frac{\lambda - \rho}{(1 + \lambda - \rho)}x, y \right) \in \{x_+\} \times \{[-1, 1]x_+\}.$$

1267 For the remaining two cases, let us assume that  $\frac{1}{(1+\lambda-\rho)}|x| < |y| < (1+\lambda)|x|$ .

- 1268 4. If  $\text{sign}(x) = \text{sign}(y)$ , let  $s = \text{sign}(x)$  and note that

$$1269 \quad \begin{bmatrix} x_+ \\ y_+ \end{bmatrix} = \frac{\lambda}{(1+\lambda)(1+\lambda-\rho)-1} \begin{bmatrix} (1+\lambda) & -1 \\ -1 & (1+\lambda-\rho) \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix}$$

$$1270 \quad = \frac{s\lambda}{(1+\lambda)(1+\lambda-\rho)-1} \begin{bmatrix} (1+\lambda)|x| - |y| \\ -|x| + (1+\lambda-\rho)|y| \end{bmatrix}$$

1271

1272 From this equation we learn  $\text{sign}(x_+) = \text{sign}(y_+) = s$ . Inverting the matrix, we  
 1273 also learn

$$\begin{aligned}
 1274 \quad \lambda \begin{bmatrix} x \\ y \end{bmatrix} &= \begin{bmatrix} (1 + \lambda - \rho) & 1 \\ 1 & (1 + \lambda) \end{bmatrix} \begin{bmatrix} x_+ \\ y_+ \end{bmatrix} = \begin{bmatrix} x_+ + \lambda(1 - \rho/\lambda)x_+ + y_+ \\ x_+ + y_+ + \lambda y_+ \end{bmatrix} \\
 1275 &= \begin{bmatrix} x_+ + \text{sign}(x_+)|y_+| + \lambda(1 - \rho/\lambda)x_+ \\ \text{sign}(y_+)|x_+| + y_+ + \lambda y_+ \end{bmatrix}.
 \end{aligned}$$

1277 Thus, the pair satisfies the inclusion.

1278 5. If  $\text{sign}(x) \neq \text{sign}(y)$ , let  $s = \text{sign}(x)$  and note that

$$\begin{aligned}
 1279 \quad \begin{bmatrix} x_+ \\ y_+ \end{bmatrix} &= \frac{\lambda}{(1 + \lambda)(1 + \lambda - \rho) - 1} \begin{bmatrix} (1 + \lambda) & 1 \\ 1 & (1 + \lambda - \rho) \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} \\
 1280 &= \frac{s\lambda}{(1 + \lambda)(1 + \lambda - \rho) - 1} \begin{bmatrix} (1 + \lambda)|x| - |y| \\ |x| - (1 + \lambda - \rho)|y| \end{bmatrix}
 \end{aligned}$$

1282 From this equation we learn  $\text{sign}(x_+) \neq \text{sign}(y_+)$ . Inverting the matrix we also  
 1283 learn

$$\begin{aligned}
 1284 \quad \lambda \begin{bmatrix} x \\ y \end{bmatrix} &= \begin{bmatrix} (1 + \lambda - \rho) & -1 \\ -1 & (1 + \lambda) \end{bmatrix} \begin{bmatrix} x_+ \\ y_+ \end{bmatrix} = \begin{bmatrix} x_+ + \lambda(1 - \rho/\lambda)x_+ - y_+ \\ -x_+ + y_+ + \lambda y_+ \end{bmatrix} \\
 1285 &= \begin{bmatrix} x_+ + \text{sign}(x_+)|y_+| + \lambda(1 - \rho/\lambda)x_+ \\ \text{sign}(y_+)|x_+| + y_+ + \lambda y_+ \end{bmatrix}.
 \end{aligned}$$

1287 Thus, the pair satisfies the inclusion.

1288 Therefore, the proof is complete. □

1289 **Corollary B.2** (Convergence to Saddles). *Assume the setting of Theorem B.1. Let  $\alpha \in$   
 1290  $(0, 1]$  and define the operator  $T := (1 - \alpha)I + \alpha S$  on  $\mathbb{R}^2$ . Then, the cone  $\mathcal{K} =$   
 1291  $\{(x, y) : |x| \leq (1 + \lambda)^{-1}|y|\}$  satisfies  $T\mathcal{K} \subseteq \mathcal{K}$ . Moreover, for any  $(x, y) \in \mathcal{K}$ , it holds  
 1292 that  $T^k(x, y) = ((1 - \alpha)^k x, (1 - \alpha(1 - \lambda(1 + \lambda)^{-1}))^k y)$  linearly converges to the  
 1293 origin as  $k$  tends to infinity.*

1294 **Proof** Since  $\mathcal{K}$  is convex, it suffices to show that  $S\mathcal{K} \subseteq \mathcal{K}$ . This follows from Theo-  
 1295 rem B.1. □

## 1296 References

- 1297 1. F. Al-Khayyal and J. Kyparisis. Finite convergence of algorithms for nonlinear programs and variational  
 1298 inequalities. *J. Optim. Theory Appl.*, 70(2):319–332, 1991.  
 1299 2. P. Albano and P. Cannarsa. Singularities of semiconcave functions in Banach spaces. In *Stochas-  
 1300 tic analysis, control, optimization and applications*, Systems Control Found. Appl., pages 171–190.  
 1301 Birkhäuser Boston, Boston, MA, 1999.  
 1302 3. H. Attouch, J. Bolte, and B.F. Svaiter. Convergence of descent methods for semi-algebraic and tame  
 1303 problems: proximal algorithms, forward–backward splitting, and regularized gauss–seidel methods.  
 1304 *Mathematical Programming*, 137(1-2):91–129, 2013.

- 1305 4. D. Avdiukhin, c. Jin, and G. Yaroslavtsev. Escaping saddle points with inequality constraints via noisy  
 1306 sticky projected gradient descent. *Optimization for Machine Learning Workshop*, 2019.
- 1307 5. A. Beck and M. Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse problems.  
 1308 *SIAM J. Imaging Sci.*, 2(1):183–202, 2009.
- 1309 6. S. Bhojanapalli, B. Neyshabur, and N. Srebro. Global optimality of local search for low rank matrix  
 1310 recovery. In *Advances in Neural Information Processing Systems*, pages 3873–3881, 2016.
- 1311 7. J. Bolte, A. Daniilidis, A. Lewis, and M. Shiota. Clarke subgradients of stratifiable functions. *SIAM*  
 1312 *Journal on Optimization*, 18(2):556–572, 2007.
- 1313 8. J.F. Bonnans and A. Shapiro. *Perturbation Analysis of Optimization Problems*. Springer, New York,  
 1314 2000.
- 1315 9. J.V. Burke. Descent methods for composite nondifferentiable optimization problems. *Math. Program-*  
 1316 *ming*, 33(3):260–279, 1985.
- 1317 10. J.V. Burke. On the identification of active constraints. II. The nonconvex case. *SIAM J. Numer. Anal.*,  
 1318 27(4):1081–1103, 1990.
- 1319 11. J.V. Burke and J.J. Moré. On the identification of active constraints. *SIAM J. Numer. Anal.*, 25(5):1197–  
 1320 1211, 1988.
- 1321 12. P.H. Calamai and J.J. Moré. Projected gradient methods for linearly constrained problems. *Math. Prog.*,  
 1322 39(1):93–116, 1987.
- 1323 13. V. Charisopoulos, Y. Chen, D. Davis, M. Díaz, L. Ding, and D. Drusvyatskiy. Low-rank matrix recovery  
 1324 with composite optimization: good conditioning and rapid convergence. *Foundations of Computational*  
 1325 *Mathematics*, pages 1–89, 2021.
- 1326 14. F.H. Clarke, Yu. Ledyaev, R.I. Stern, and P.R. Wolenski. *Nonsmooth Analysis and Control Theory*.  
 1327 Texts in Math. 178, Springer, New York, 1998.
- 1328 15. C. Criscitiello and N. Boumal. Efficiently escaping saddle points on manifolds. In *Advances in Neural*  
 1329 *Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- 1330 16. D. Davis and D. Drusvyatskiy. Stochastic model-based minimization of weakly convex functions.  
 1331 *SIAM Journal on Optimization*, 29(1):207–239, 2019.
- 1332 17. D. Drusvyatskiy. The proximal point method revisited. *SIAG/OPT Views and News*, 26(2), 2018.
- 1333 18. D. Drusvyatskiy, A.D. Ioffe, and A.S. Lewis. Generic minimizing behavior in semialgebraic optimiza-  
 1334 tion. *SIAM Journal on Optimization*, 26(1):513–534, 2016.
- 1335 19. D. Drusvyatskiy and A.S. Lewis. Optimality, identifiability, and sensitivity. *Math. Program.*, 147(1-2,  
 1336 Ser. A):467–498, 2014.
- 1337 20. D. Drusvyatskiy and A.S. Lewis. Error bounds, quadratic growth, and linear convergence of proximal  
 1338 methods. *Mathematics of Operations Research*, 43(3):919–948, 2018.
- 1339 21. D. Drusvyatskiy and C. Paquette. Efficiency of minimizing compositions of convex functions and  
 1340 smooth maps. *Mathematical Programming*, 178(1-2):503–558, 2019.
- 1341 22. S.S. Du, C. Jin, J.D. Lee, M.I. Jordan, A. Singh, and B. Póczos. Gradient descent can take exponential  
 1342 time to escape saddle points. In *Advances in neural information processing systems*, pages 1067–1077,  
 1343 2017.
- 1344 23. J.C. Duchi and F. Ruan. Stochastic methods for composite and weakly convex optimization problems.  
 1345 *SIAM Journal on Optimization*, 28(4):3229–3259, 2018.
- 1346 24. J.C. Dunn. On the convergence of projected gradient processes to singular critical points. *J. Optim.*  
 1347 *Theory Appl.*, 55(2):203–216, 1987.
- 1348 25. M.C. Ferris. Finite termination of the proximal point algorithm. *Math. Program.*, 50(3, (Ser. A)):359–  
 1349 366, 1991.
- 1350 26. S.D. Flåm. On finite convergence and constraint identification of subgradient projection methods.  
 1351 *Math. Program.*, 57:427–437, 1992.
- 1352 27. R. Ge, F. Huang, C. Jin, and Y. Yuan. Escaping from saddle points—online stochastic gradient for  
 1353 tensor decomposition. In *Conference on Learning Theory*, pages 797–842, 2015.
- 1354 28. R. Ge, C. Jin, and Y. Zheng. No spurious local minima in nonconvex low rank problems: A unified  
 1355 geometric analysis. In *Proceedings of the 34th International Conference on Machine Learning-Volume*  
 1356 *70*, pages 1233–1242. JMLR. org, 2017.
- 1357 29. R. Ge, J.D. Lee, and T. Ma. Matrix completion has no spurious local minimum. In *Advances in Neural*  
 1358 *Information Processing Systems*, pages 2973–2981, 2016.
- 1359 30. N. Hallak and M. Teboulle. Finding second-order stationary points in constrained minimization: A  
 1360 feasible direction approach. *Journal of Optimization Theory and Applications*, 186(2):480–503, 2020.
- 1361 31. W.L. Hare and A.S. Lewis. Identifying active manifolds. *Algorithmic Oper. Res.*, 2(2):75–82, 2007.

- 1362 32. C. Jin, P. Netrapalli, and M. Jordan. What is local optimality in nonconvex-nonconcave minimax  
1363 optimization? In *International Conference on Machine Learning*, pages 4880–4889. PMLR, 2020.
- 1364 33. R. Jin, C. Ge, P. Netrapalli, S.M. Kakade, and M.I. Jordan. How to escape saddle points efficiently. In  
1365 *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 1724–1732.  
1366 JMLR. org, 2017.
- 1367 34. J.D. Lee, I. Panageas, G. Piliouras, M. Simchowitz, M.I. Jordan, and B. Recht. First-order methods  
1368 almost always avoid strict saddle points. *Math. Program.*, 176(1-2):311–337, 2019.
- 1369 35. J.D. Lee, M. Simchowitz, M.I. Jordan, and B. Recht. Gradient descent only converges to minimizers.  
1370 In *Conference on learning theory*, pages 1246–1257, 2016a.
- 1371 36. J.M. Lee. Smooth manifolds. In *Introduction to Smooth Manifolds*, pages 1–31. Springer, 2013.
- 1372 37. Sangkyun Lee and Stephen J Wright. Manifold identification in dual averaging for regularized stochas-  
1373 tic online learning. *Journal of Machine Learning Research*, 13(Jun):1705–1744, 2012.
- 1374 38. C. Lemaréchal, F. Oustry, and C. Sagastizábal. The U-lagrangian of a convex function. *Trans. Amer.*  
1375 *Math. Soc.*, 352:711–729, 1996.
- 1376 39. A.S. Lewis. Active sets, nonsmoothness, and sensitivity. *SIAM J. Optim.*, 13(3):702–725 (electronic)  
1377 (2003), 2002.
- 1378 40. A.S. Lewis and S.J. Wright. A proximal method for composite minimization. *Math. Program.*, pages  
1379 1–46, 2015.
- 1380 41. A.S. Lewis and S. Zhang. Partial smoothness, tilt stability, and generalized Hessians. *SIAM Journal*  
1381 *on Optimization*, 23(1):74–94, 2013.
- 1382 42. B. Martinet. Régularisation d’inéquations variationnelles par approximations successives. *Rev.*  
1383 *Française Informat. Rech. Opérationnelle*, 4(Sér. R-3):154–158, 1970.
- 1384 43. B. Martinet. Détermination approchée d’un point fixe d’une application pseudo-contractante. Cas de  
1385 l’application prox. *C. R. Acad. Sci. Paris Sér. A-B*, 274:A163–A165, 1972.
- 1386 44. A. Mokhtari, A. Ozdaglar, and A. Jadbabaie. Escaping saddle points in constrained optimization. In  
1387 *Advances in Neural Information Processing Systems*, pages 3629–3639, 2018.
- 1388 45. B.S. Mordukhovich. *Variational analysis and generalized differentiation. I*, volume 330 of *Grundlehren*  
1389 *der Mathematischen Wissenschaften [Fundamental Principles of Mathematical Sciences]*. Springer-  
1390 Verlag, Berlin, 2006. Basic theory.
- 1391 46. J.-J. Moreau. Proximité et dualité dans un espace hilbertien. *Bull. Soc. Math. France*, 93:273–299,  
1392 1965.
- 1393 47. Yu. Nesterov. Modified Gauss–Newton scheme with worst case guarantees for global performance.  
1394 *Optimisation Methods and Software*, 22(3):469–483, 2007.
- 1395 48. Yu. Nesterov. Gradient methods for minimizing composite functions. *Mathematical Programming*,  
1396 140(1):125–161, 2013.
- 1397 49. M. Nouiehed, J.D. Lee, and M. Razaviyayn. Convergence to second-order stationarity for constrained  
1398 non-convex optimization. arXiv preprint 1810.02024, 2018.
- 1399 50. E.A. Nurminkii. The quasigradient method for the solving of the nonlinear programming problems.  
1400 *Cybernetics*, 9(1):145–150, 1973.
- 1401 51. I. Panageas and G. Piliouras. Gradient descent only converges to minimizers: Non-isolated critical  
1402 points and invariant regions. arXiv preprint 1605.00405, 2016.
- 1403 52. E. Pauwels. The value function approach to convergence analysis in composite optimization. *Opera-*  
1404 *tions Research Letters*, 44(6):790–795, 2016.
- 1405 53. R.A. Poliquin and R.T. Rockafellar. Prox-regular functions in variational analysis. *Trans. Amer. Math.*  
1406 *Soc.*, 348:1805–1838, 1996.
- 1407 54. R.T. Rockafellar. *Convex analysis*. Princeton Mathematical Series, No. 28. Princeton University Press,  
1408 Princeton, N.J., 1970.
- 1409 55. R.T. Rockafellar. Monotone operators and the proximal point algorithm. *SIAM J. Control Optimization*,  
1410 14(5):877–898, 1976.
- 1411 56. R.T. Rockafellar. Favorable classes of Lipschitz-continuous functions in subgradient optimization. In  
1412 *Progress in nondifferentiable optimization*, volume 8 of *IASA Collaborative Proc. Ser. CP-82*, pages  
1413 125–143. Int. Inst. Appl. Sys. Anal., Laxenburg, 1982.
- 1414 57. R.T. Rockafellar and R.J-B. Wets. *Variational Analysis*. Grundlehren der mathematischen Wis-  
1415 senschaften, Vol 317, Springer, Berlin, 1998.
- 1416 58. S. Rolewicz. On paraconvex multifunctions. In *Third Symposium on Operations Research (Univ.*  
1417 *Mannheim, Mannheim, 1978)*, Section I, volume 31 of *Operations Res. Verfahren*, pages 539–546.  
1418 Hain, Königstein/Ts., 1979.

- 1419 59. A. Shapiro. Second order sensitivity analysis and asymptotic theory of parametrized nonlinear pro-  
1420 grams. *Mathematical Programming*, 33(3):280–299, 1985.
- 1421 60. M. Shub. *Global stability of dynamical systems*. Springer Science & Business Media, 2013.
- 1422 61. J. Sun, Q. Qu, and J. Wright. When are nonconvex problems not scary? arXiv preprint 1510.06096,  
1423 2015.
- 1424 62. J. Sun, Q. Qu, and J. Wright. A geometric analysis of phase retrieval. *Foundations of Computational*  
1425 *Mathematics*, 18(5):1131–1198, 2018.
- 1426 63. Y. Sun, N. Flammarion, and M. Fazel. Escaping from saddle points on riemannian manifolds. In  
1427 *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- 1428 64. S.J. Wright. Identifiable surfaces in constrained optimization. *SIAM J. Control Optim.*, 31:1063–1079,  
1429 1993.

1430 **Publisher's Note** Springer Nature remains neutral with regard to jurisdictional claims in published maps  
1431 and institutional affiliations.

Journal: 10208 Article: 9516
---------------------------------

## Author Query Form

**Please ensure you fill out your response to the queries raised below  
and return this form along with your corrections**

Dear Author

During the process of typesetting your article, the following queries have arisen. Please check your typeset proof carefully against the queries listed below and mark the necessary changes either directly on the proof/online grid or in the 'Author's response' area provided below

Query	Details required	Author's response
1.	Please confirm if the corresponding author is correctly identified.	