

1                   **PROJECTED GRADIENT DESCENT ACCUMULATES AT**  
2                   **BOULIGAND STATIONARY POINTS\***

3                   GUILLAUME OLIKIER<sup>†</sup> AND IRÈNE WALDSPURGER<sup>‡</sup>

4           **Abstract.** This paper considers the projected gradient descent (PGD) algorithm for the problem  
5 of minimizing a continuously differentiable function on a nonempty closed subset of a Euclidean vector  
6 space. Without further assumptions, this problem is intractable and algorithms are only expected  
7 to find a stationary point. PGD is known to generate a sequence whose accumulation points are  
8 Mordukhovich stationary. In this paper, these accumulation points are proven to be Bouligand  
9 stationary, and even proximally stationary if the gradient is locally Lipschitz continuous. These are  
10 the strongest stationarity properties that can be expected for the considered problem.

11           **Key words.** projected gradient descent, stationarity, criticality, tangent and normal cones,  
12 Clarke regularity

13           **MSC codes.** 65K10, 49J53, 90C26, 90C30, 90C46

14           **1. Introduction.** Let  $\mathcal{E}$  be a Euclidean vector space,  $C$  a nonempty closed sub-  
15 set of  $\mathcal{E}$ , and  $f : \mathcal{E} \rightarrow \mathbb{R}$  a function satisfying at least the first of the following two  
16 properties:

17           (H1)  $f$  is differentiable on  $C$ , i.e., for every  $x \in C$ , there exists a (unique) vector  
18 in  $\mathcal{E}$ , denoted by  $\nabla f(x)$ , such that

$$19 \quad \lim_{\substack{y \rightarrow x \\ y \in \mathcal{E} \setminus \{x\}}} \frac{f(y) - f(x) - \langle \nabla f(x), y - x \rangle}{\|y - x\|} = 0,$$

20           and  $\nabla f : C \rightarrow \mathcal{E}$  is continuous;

21           (H2)  $f$  is differentiable on  $\mathcal{E}$  and  $\nabla f : \mathcal{E} \rightarrow \mathcal{E}$  is locally Lipschitz continuous.

22 This paper considers the problem

$$23 \quad (1.1) \quad \min_{x \in C} f(x)$$

24 of minimizing  $f$  on  $C$ . In general, without further assumptions on  $C$  or  $f$ , finding an  
25 exact or approximate global minimizer of problem (1.1) is intractable. Even finding  
26 an approximate local minimizer is not always feasible in polynomial time (unless  
27  $P = NP$ ) [2]. Therefore, algorithms are only expected to return a point satisfying a  
28 condition called *stationarity*, which is a tractable surrogate for local optimality.

29           A point  $x \in C$  is said to be stationary for (1.1) if  $-\nabla f(x)$  is normal to  $C$  at  $x$ .  
30 Several definitions of normality exist. Each one defines a different notion of station-  
31 arity, which is a surrogate for local optimality in the sense that, possibly under mild  
32 regularity assumptions on  $f$ , every local minimizer of  $f|_C$  is stationary for (1.1). In  
33 particular, each of the three notions of normality in [53, Definition 6.3 and Exam-  
34 ple 6.16], namely normality in the general sense, in the regular sense, and in the

---

\*This work was supported by the ERC grant #786854 G-Statistics from the European Research Council under the European Union’s Horizon 2020 research and innovation program and by the French government through the 3IA Côte d’Azur Investments ANR-19-P3IA-0002 and the PRAIRIE 3IA Institute ANR-19-P3IA-0001, managed by the National Research Agency.

<sup>†</sup>Université Côte d’Azur and Inria, Epione Project Team, 2004 route des Lucioles - BP 93, 06902 Sophia Antipolis Cedex, France (guillaume.olikier@inria.fr).

<sup>‡</sup>CNRS, Université Paris Dauphine, équipe-projet Mokaplan (Inria), place du Maréchal de Lattre de Tassigny, 75016 Paris, France (waldspurger@ceremade.dauphine.fr).

proximal sense, yields an important definition of stationarity. The sets of general, regular, and proximal normals to  $C$  at  $x \in C$  are respectively denoted by  $N_C(x)$ ,  $\widehat{N}_C(x)$ , and  $\widehat{\widehat{N}}_C(x)$ . These sets are reviewed in Section 2.2. Importantly, they are nested as follows: for every  $x \in C$ ,

$$(1.2) \quad \widehat{\widehat{N}}_C(x) \subseteq \widehat{N}_C(x) \subseteq N_C(x),$$

and  $C$  is said to be *Clarke regular* at  $x$  if the second inclusion is an equality. The definitions of stationarity based on these sets are given in Definition 1.1, and the terminology is discussed in Section 3.

DEFINITION 1.1. *For problem (1.1), a point  $x \in C$  is said to be:*

- Mordukhovich stationary (M-stationary) if  $-\nabla f(x) \in N_C(x)$ ;
- Bouligand stationary (B-stationary) if  $-\nabla f(x) \in \widehat{N}_C(x)$ ;
- proximally stationary (P-stationary) if  $-\nabla f(x) \in \widehat{\widehat{N}}_C(x)$ .

There are many practical examples of a set  $C$  for which at least one of the inclusions in (1.2) is strict, especially the second one. This is notably shown by the four examples studied in Section 7, where the second inclusion is strict at infinitely many points. The three notions of stationarity are therefore not equivalent. Actually, as explained next, B-stationarity and P-stationarity are the strongest necessary conditions for local optimality under different sets of assumptions on  $f$ , while M-stationarity is a weaker condition.

As pointed out in [13, §5], for problem (1.1) under the only assumption that  $f$  is differentiable on  $C$ , B-stationarity is the strongest necessary condition for local optimality. The same is true if  $f$  satisfies (H1). Indeed, by [53, Theorem 6.11], for all  $x \in C$ ,

$$(1.3) \quad \widehat{N}_C(x) = \left\{ -\nabla h(x) \mid \begin{array}{l} h : \mathcal{E} \rightarrow \mathbb{R} \text{ is differentiable at } x, \\ x \text{ is a local minimizer of } h|_C \end{array} \right\}$$

$$(1.4) \quad = \left\{ -\nabla h(x) \mid \begin{array}{l} h : \mathcal{E} \rightarrow \mathbb{R} \text{ satisfies (H1),} \\ x \text{ is a local minimizer of } h|_C \end{array} \right\}.$$

The inclusion  $\supseteq$  in (1.3) shows that every local minimizer of  $f|_C$  is B-stationary for (1.1). Thus,  $\widehat{N}_C(x)$  is sufficiently large to yield a necessary condition for local optimality. The inclusion  $\subseteq$  in (1.3) shows that replacing  $\widehat{N}_C(x)$  with one of its proper subsets would yield a condition that is not necessary for local optimality. The equality (1.4) shows that these observations also hold if  $f$  satisfies (H1).

P-stationarity is the strongest necessary condition for local optimality if  $f$  satisfies (H2). Indeed, by Theorem 2.5, for all  $x \in C$ ,

$$(1.5) \quad \widehat{\widehat{N}}_C(x) = \left\{ -\nabla h(x) \mid \begin{array}{l} h : \mathcal{E} \rightarrow \mathbb{R} \text{ satisfies (H2),} \\ x \text{ is a local minimizer of } h|_C \end{array} \right\}.$$

The inclusion  $\supseteq$  in (1.5) shows that, under (H2), every local minimizer of  $f|_C$  is P-stationary for (1.1). The inclusion  $\subseteq$  in (1.5) shows that replacing  $\widehat{\widehat{N}}_C(x)$  with one of its proper subsets would yield a condition that is not necessary for local optimality.

In comparison, M-stationarity is a weaker notion of stationarity which is considered unsatisfactory in [27, §4], [32, §1], and [50, §2.1]. Furthermore, as explained in [32], distinguishing convergence to a B-stationary point from convergence to an M-stationary point is difficult (a phenomenon formalized by the notion of *apocalypse* in [32]) in the sense that it cannot be done based on standard measures of

76 B-stationarity because of their possible lack of lower semicontinuity at points where  
 77 the feasible set is not Clarke regular, as also explained in [46, §§1.1 and 2.4].

78 Projected gradient descent, or PGD for short, is a basic algorithm aiming at  
 79 solving problem (1.1). To the best of our knowledge, the first article to have considered  
 80 PGD on a possibly nonconvex closed set was [6]. The nonmonotone backtracking  
 81 version considered in this paper is defined as Algorithm 4.2 and is based on [30,  
 82 Algorithm 3.1] and [11, Algorithm 3.1]. Given  $x \in C$  as input, the iteration map of  
 83 PGD, called the PGD map and defined as Algorithm 4.1, performs a backtracking  
 84 projected line search along the direction of  $-\nabla f(x)$ , i.e., computes a projection  $y$   
 85 of  $x - \alpha \nabla f(x)$  onto  $C$  for decreasing values of the step size  $\alpha \in (0, \infty)$  until  $y$   
 86 satisfies an Armijo condition. In the simplest version of PGD, called *monotone*,  
 87 the Armijo condition ensures that the value of  $f$  at the next iterate is smaller by a  
 88 specified amount than the value at the current iterate. Following the general settings  
 89 proposed in [30, 31] and [11], the value at the current iterate can be replaced with  
 90 the maximum value of  $f$  over a prefixed number of the previous iterates (“max” rule)  
 91 or with a weighted average of the values of  $f$  at the previous iterates (“average”  
 92 rule). This version of PGD is called *nonmonotone*. By [31, Theorem 3.1], monotone  
 93 PGD accumulates at M-stationary points of (1.1) if  $f$  is continuously differentiable  
 94 on  $\mathcal{E}$  and bounded from below on  $C$ . By [11, Theorem 4.6], the same result holds  
 95 for nonmonotone PGD with the “average” rule and, by [31, Theorem 4.1], also for  
 96 nonmonotone PGD with the “max” rule if  $f$  is further uniformly continuous on the  
 97 sublevel set

$$98 \quad (1.6) \quad \{x \in C \mid f(x) \leq f(x_0)\},$$

99 where  $x_0 \in C$  is the initial iterate given to the algorithm. However, as pointed out in  
 100 [32, §1], it is an open question whether the accumulation points of PGD are always  
 101 B-stationary for (1.1).

102 This paper answers positively the question by proving Theorem 1.2.

103 **THEOREM 1.2.** *Consider a sequence generated by PGD (Algorithm 4.2) when ap-*  
 104 *plied to problem (1.1).*

- 105 • *If this sequence is finite, then its last element is B-stationary for (1.1) un-*  
 106 *der (H1), and even P-stationary for (1.1) under (H2).*
- 107 • *If this sequence is infinite, then all of its accumulation points, if any, are B-*  
 108 *stationary for (1.1) under (H1), and even P-stationary for (1.1) under (H2).*

109 If  $\nabla f$  is globally Lipschitz continuous, then it is known in the literature that  
 110 every local minimizer of  $f|_C$  is P-stationary for (1.1) [61, Proposition 3.5(ii)] (the  
 111 result is given for a global minimizer but the proof shows that it also holds for a local  
 112 minimizer) and that PGD with a constant step size smaller than the inverse of the  
 113 Lipschitz constant accumulates at P-stationary points of (1.1) [61, Theorem 5.6(i)].  
 114 Indeed, the ZeroFPR algorithm proposed in [61] extends the proximal gradient (PG)  
 115 algorithm with a constant step size [61, Remark 5.5] which itself extends PGD with  
 116 a constant step size; problem (1.1) corresponds to [61, problem (1.1)] with  $g$  the  
 117 indicator function of our set  $C$ . These results were rediscovered in [50] where, in  
 118 addition, the distance from the negative gradient of the continuously differentiable  
 119 function to the regular subdifferential of the other function is proven to converge to  
 120 zero along the generated sequence, and a quadratic lower bound on  $f - f(\bar{x})$  at every  
 121 accumulation point  $\bar{x}$  of PG is obtained. The two results cited from [61] were already  
 122 stated in [5, Theorems 2.2 and 3.1] for  $C$  the set  $\mathbb{R}_{\leq s}^n$  of vectors of  $\mathbb{R}^n$  having at most  
 123  $s$  nonzero components for some positive integer  $s < n$ , and in [3, Proposition 1 and

124 Theorem 1] for  $C$  satisfying a regularity condition called *proximal smoothness* which  
 125 none of the four examples studied in Section 7 satisfies.

126 This paper is organized as follows. The necessary background in variational analy-  
 127 sis is introduced in Section 2. The literature on stationarity is partially surveyed in  
 128 Section 3. The PGD algorithm is reviewed in Section 4. It is analyzed under hy-  
 129 pothesis (H1) in Section 5 and under hypothesis (H2) in Section 6. Four practical  
 130 examples of a set  $C$  for which the first inclusion in (1.2) is an equality for all  $x \in C$   
 131 and the second is strict for infinitely many  $x \in C$  are given in Section 7. Theorem 1.2  
 132 is illustrated by a comparison between PGD and a first-order algorithm that is not  
 133 guaranteed to accumulate at B-stationary points of (1.1) in Section 8. Concluding  
 134 remarks are gathered in Section 9.

135 **2. Elements of variational analysis.** This section, mostly based on [53], re-  
 136 views background material in variational analysis that is used in the rest of the paper.  
 137 Section 2.1 concerns the projection map onto  $C$  and its main properties. Section 2.2  
 138 reviews the three notions of normality on which the three notions of stationarity  
 139 provided in Definition 1.1 are based.

140 Recall that, throughout the paper,  $\mathcal{E}$  is a Euclidean vector space and  $C \subseteq \mathcal{E}$  is  
 141 nonempty and closed. Moreover, for every  $x \in \mathcal{E}$  and  $\rho \in (0, \infty)$ ,  $B(x, \rho) := \{y \in \mathcal{E} \mid$   
 142  $\|x - y\| < \rho\}$  and  $B[x, \rho] := \{y \in \mathcal{E} \mid \|x - y\| \leq \rho\}$  are respectively the open and  
 143 closed balls of center  $x$  and radius  $\rho$  in  $\mathcal{E}$ . Following [53, §3B], a nonempty subset  $K$   
 144 of  $\mathcal{E}$  is called a *cone* if  $x \in K$  implies  $\alpha x \in K$  for all  $\alpha \in [0, \infty)$ .

145 **2.1. Projection map.** Given  $x \in \mathcal{E}$ , the distance from  $x$  to  $C$  is  $d(x, C) :=$   
 146  $\min_{y \in C} \|x - y\|$  and the projection of  $x$  onto  $C$  is  $P_C(x) := \operatorname{argmin}_{y \in C} \|x - y\|$ . By  
 147 [53, Example 1.20], the function  $\mathcal{E} \rightarrow \mathbb{R} : x \mapsto d(x, C)$  is continuous and, for every  
 148  $x \in \mathcal{E}$ , the set  $P_C(x)$  is nonempty and compact. Proposition 2.1 is invoked frequently  
 149 in the rest of the paper.

150 PROPOSITION 2.1. *For all  $x \in C$ ,  $v \in \mathcal{E}$ , and  $y \in P_C(x - v)$ ,*

$$151 \quad (2.1) \quad \|y - x\| \leq 2\|v\|,$$

$$152 \quad (2.2) \quad 2\langle v, y - x \rangle \leq -\|y - x\|^2,$$

153 *and the inequalities are strict if  $x \notin P_C(x - v)$ .*

154 *Proof.* By definition of the projection,  $\|y - (x - v)\| \leq \|x - (x - v)\| = \|v\|$  and  
 155 the inequality is strict if  $x \notin P_C(x - v)$ . Thus, on the one hand,

$$156 \quad \|y - x\| = \|y - (x - v) - v\| \leq \|y - (x - v)\| + \|-v\| \leq \|v\| + \|v\| = 2\|v\|,$$

157 and, on the other hand,  $\|y - (x - v)\|^2 \leq \|v\|^2$ , which is equivalent to (2.2).  $\square$

158 **2.2. Normality and stationarity.** Based on [53, Chapter 6], this section re-  
 159 views the three notions of normality on which the three notions of stationarity given  
 160 in Definition 1.1 are based.

161 Following [53, Definition 6.1], a vector  $v \in \mathcal{E}$  is said to be *tangent* to  $C$  at  $x \in C$   
 162 if there exist sequences  $(x_i)_{i \in \mathbb{N}}$  in  $C$  converging to  $x$  and  $(t_i)_{i \in \mathbb{N}}$  in  $(0, \infty)$  such that  
 163 the sequence  $(\frac{x_i - x}{t_i})_{i \in \mathbb{N}}$  converges to  $v$ . The set of all tangent vectors to  $C$  at  $x \in C$   
 164 is a closed cone [53, Proposition 6.2] called the *tangent cone* to  $C$  at  $x$  and denoted  
 165 by  $T_C(x)$ . Following [53, Definition 6.3 and Proposition 6.5], the *regular normal cone*  
 166 to  $C$  at  $x \in C$  is

$$167 \quad \widehat{N}_C(x) := \{v \in \mathcal{E} \mid \langle v, w \rangle \leq 0 \forall w \in T_C(x)\}$$

168 which is a closed convex cone. Following [53, Definition 6.3], a vector  $v \in \mathcal{E}$  is said to  
 169 be *normal (in the general sense)* to  $C$  at  $x \in C$  if there exist sequences  $(x_i)_{i \in \mathbb{N}}$  in  $C$   
 170 converging to  $x$  and  $(v_i)_{i \in \mathbb{N}}$  converging to  $v$  such that, for all  $i \in \mathbb{N}$ ,  $v_i \in \widehat{N}_C(x_i)$ . The  
 171 set of all normal vectors to  $C$  at  $x \in C$  is a closed cone [53, Proposition 6.5] called  
 172 the *normal cone* to  $C$  at  $x$  and denoted by  $N_C(x)$ . Following [53, Example 6.16], a  
 173 vector  $v \in \mathcal{E}$  is called a *proximal normal* to  $C$  at  $x \in C$  if there exists  $\bar{\alpha} \in (0, \infty)$  such  
 174 that  $x \in P_C(x + \bar{\alpha}v)$ , i.e.,  $\bar{\alpha}\|v\| = d(x + \bar{\alpha}v, C)$ , which implies that, for all  $\alpha \in [0, \bar{\alpha})$ ,  
 175  $P_C(x + \alpha v) = \{x\}$ . The set of all proximal normals to  $C$  at  $x \in C$  is a convex cone  
 176 called the *proximal normal cone* to  $C$  at  $x$  and denoted by  $\widehat{N}_C(x)$ .

177 As stated in (1.2), for all  $x \in C$ ,

$$178 \quad \widehat{N}_C(x) \subseteq \widehat{N}_C(x) \subseteq N_C(x).$$

179 Following [53, Definition 6.4],  $C$  is said to be *Clarke regular* at  $x \in C$  if  $\widehat{N}_C(x) =$   
 180  $N_C(x)$ . Thus, M-stationarity is equivalent to B-stationarity at a point  $x \in C$  if and  
 181 only if  $C$  is Clarke regular at  $x$ , which is not the case in many practical situations, as  
 182 shown by the four examples given in Section 7. For those examples, however, regular  
 183 normals are proximal normals (Proposition 7.1). An example of a set  $C$  and a point  
 184  $x \in C$  such that both inclusions in (1.2) are strict is given in Example 2.2.

185 **EXAMPLE 2.2.** Let  $\mathcal{E} := \mathbb{R}^2$  and  $C := \{(t, \max\{0, t^{3/5}\}) \mid t \in \mathbb{R}\}$  (inspired by [53,  
 186 Figure 6–12(a)]). Then,

$$187 \quad T_C(0, 0) = (\{0\} \times [0, \infty)) \cup ((-\infty, 0] \times \{0\}),$$

$$188 \quad \widehat{N}_C(0, 0) = [0, \infty) \times (-\infty, 0],$$

$$189 \quad \widehat{N}_C(0, 0) = \widehat{N}_C(0, 0) \setminus ((0, \infty) \times \{0\}),$$

$$190 \quad N_C(0, 0) = \widehat{N}_C(0, 0) \cup T_C(0, 0).$$

191 Thus,

$$192 \quad \widehat{N}_C(0, 0) \subsetneq \widehat{N}_C(0, 0) \subsetneq N_C(0, 0).$$

This is illustrated in Figure 1.

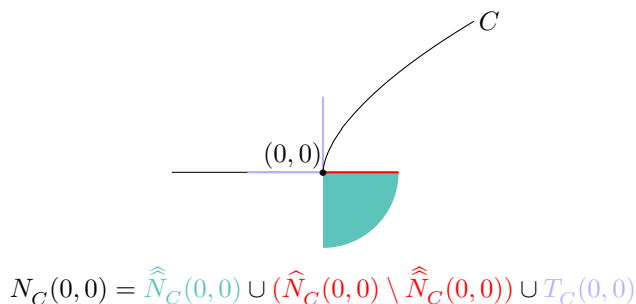


FIG. 1. Tangent and normal cones from Example 2.2.

193

194 As pointed out in Section 1, the regular and proximal normal cones enjoy gradient  
 195 characterizations which imply that B- and P-stationarity are the strongest necessary  
 196 conditions for local optimality under different sets of assumptions on  $f$ . Those given

197 in (1.3)–(1.4) come from [53, Theorem 6.11]. That given in (1.5) comes from Theo-  
198 rem 2.5, established at the end of this section.

199 As shown by (1.4), for problem (1.1), B-stationarity is the strongest necessary  
200 condition for local optimality if  $f$  is only assumed to satisfy (H1). In particular, under  
201 this assumption, P-stationarity is not necessary for local optimality, as illustrated by  
202 Example 2.3.

203 **EXAMPLE 2.3.** Let  $\mathcal{E} := \mathbb{R}^2$ ,  $C := \{(x_1, x_2) \in \mathbb{R}^2 \mid x_2 \geq \max\{0, x_1^{3/5}\}\}$  [53, Fig-  
204 ure 6–12(a)], and  $f : \mathbb{R}^2 \rightarrow \mathbb{R} : (x_1, x_2) \mapsto \frac{1}{2}(x_1 - 1)^2 + |x_2|^{3/2}$ . Then,  $f$  is continu-  
205 ously differentiable on  $\mathcal{E}$ , hence on  $C$ , and, for all  $(x_1, x_2) \in \mathbb{R}^2$ ,  $\nabla f(x_1, x_2) = (x_1 -$   
206  $1, \frac{3}{2}\text{sgn}(x_2)|x_2|^{1/2})$ . Thus,  $-\nabla f(0, 0) = (1, 0) \in \widehat{N}_C(0, 0) \setminus \widehat{\widehat{N}}_C(0, 0)$ , yet  $\text{argmin}_C f =$   
207  $\{(0, 0)\}$ .

208 Proposition 2.4 states that P-stationarity is necessary for local optimality if  $f$  is  
209 assumed to satisfy (H2), that is,  $f$  is differentiable on  $\mathcal{E}$  and  $\nabla f$  is locally Lipschitz  
210 continuous. The latter means that, for every open or closed ball  $\mathcal{B} \subseteq \mathcal{E}$ ,

$$211 \quad \text{Lip}(\nabla f) := \sup_{\substack{\mathcal{B} \\ x, y \in \mathcal{B} \\ x \neq y}} \frac{\|\nabla f(x) - \nabla f(y)\|}{\|x - y\|} < \infty,$$

212 which implies, by [44, Lemma 1.2.3], that, for all  $x, y \in \mathcal{B}$ ,

$$213 \quad (2.3) \quad |f(y) - f(x) - \langle \nabla f(x), y - x \rangle| \leq \frac{\text{Lip}_{\mathcal{B}}(\nabla f)}{2} \|y - x\|^2.$$

214 **PROPOSITION 2.4.** Assume that  $f$  satisfies (H2). If  $x \in C$  is a local minimizer of  
215  $f|_C$ , then  $-\nabla f(x) \in \widehat{\widehat{N}}_C(x)$ .

216 *Proof.* By contrapositive. Assume that  $-\nabla f(x) \notin \widehat{\widehat{N}}_C(x)$  for some  $x \in C$ . Let  
217  $\rho \in (0, \infty)$ . Then, for all  $\alpha \in (0, \frac{\rho}{2\|\nabla f(x)\|}]$ ,

$$218 \quad x \notin P_C(x - \alpha \nabla f(x)) \subseteq B(x, 2\alpha \|\nabla f(x)\|) \subseteq B(x, \rho),$$

219 where the first inclusion holds by (2.1). Thus, by (2.3) and (2.2), for all  $\alpha \in$   
220  $(0, \min\{\frac{\rho}{2\|\nabla f(x)\|}, \frac{1}{\text{Lip}_{B(x, \rho)}(\nabla f)}\})$  and  $y \in P_C(x - \alpha \nabla f(x))$ ,

$$221 \quad \begin{aligned} f(y) - f(x) &\leq \langle \nabla f(x), y - x \rangle + \frac{\text{Lip}_{B(x, \rho)}(\nabla f)}{2} \|y - x\|^2 \\ 222 \quad &< \left( -\frac{1}{2\alpha} + \frac{\text{Lip}_{B(x, \rho)}(\nabla f)}{2} \right) \|y - x\|^2 \\ 223 \quad &\leq 0. \end{aligned}$$

224 Hence,  $x$  is not a local minimizer of  $f|_C$ . □

225 Theorem 2.5 strengthens [53, Proposition 8.46(d)] by stating that (1.5) is valid,  
226 which shows that P-stationarity is the strongest necessary condition for local opti-  
227 mality under hypothesis (H2).

228 **THEOREM 2.5** (gradient characterization of proximal normals). For every  $x \in C$ ,  
229 (1.5) holds.

230 *Proof.* Let  $x \in C$ . The inclusion  $\supseteq$  holds by Proposition 2.4. For the inclusion  $\subseteq$ ,  
231 let  $v \in \widehat{\widehat{N}}_C(x)$ . By definition of the proximal normal cone, there exists  $\bar{\alpha} \in (0, \infty)$

232 such that  $x \in P_C(x + \bar{\alpha}v)$ . This is equivalent to the fact that  $x$  is a minimizer of  $h|_C$ ,  
 233 where  $h : \mathcal{E} \rightarrow \mathbb{R}$  is defined by

$$234 \quad h(y) := \frac{1}{2\bar{\alpha}} \|y - (x + \bar{\alpha}v)\|^2 \quad \forall y \in \mathcal{E}.$$

235 The function  $h$  is differentiable, its gradient is locally Lipschitz continuous (actually,  
 236 globally Lipschitz continuous, since it is an affine map), and

$$237 \quad -\nabla h(x) = v.$$

238 Since  $x$  is a global minimizer of  $h|_C$ , it is also a local minimizer of  $h|_C$ . This shows  
 239 that

$$240 \quad v \in \left\{ -\nabla h(x) \mid \begin{array}{l} h : \mathcal{E} \rightarrow \mathbb{R} \text{ satisfies (H2),} \\ x \text{ is a local minimizer of } h|_C \end{array} \right\},$$

241 which implies the inclusion  $\subseteq$  in (1.5). □

242 *Remark 2.6.* From our proof, we see that (1.5) is also true if we replace “local  
 243 minimizer” with “global minimizer”. The same holds for equations (1.3)–(1.4) [53,  
 244 Theorem 6.11]. However, in this section, we are interested in understanding the  
 245 closeness between the notions of stationarity and *local* optimality.

246 **3. Stationarity in the literature.** This section surveys the names given to the  
 247 stationarity notions provided in Definition 1.1 and attempts to offer a brief historical  
 248 perspective. The terms “B-stationarity” and “M-stationarity” first appeared in the  
 249 literature about mathematical programs with equilibrium constraints (MPECs), as  
 250 explained in Sections 3.1 and 3.2. In contrast, the term “P-stationarity” seems to be  
 251 new in the literature. P-stationarity is called “criticality” in [61, Definition 3.1(ii)];  
 252 problem (1.1) corresponds to [61, problem (1.1)] with  $g$  the indicator function of our  
 253 set  $C$ . We propose the name “P-stationarity” because this stationarity notion is based  
 254 on the proximal normal cone. It is closely related to the so-called  $\alpha$ -stationarity, as  
 255 explained in Section 3.3.

256 **3.1. A brief history of Bouligand stationarity.** Peano already knew that  
 257 B-stationarity is a necessary condition for optimality. The statement is implicit in  
 258 his 1887 book *Applicazioni geometriche del calcolo infinitesimale* and explicit in his  
 259 1908 book *Formulario Mathematico* where the formulation is based on the tangent  
 260 cone and the derivative defined in the same book; see the historical investigation in  
 261 [12, 13].

262 B-stationarity appears as a necessary condition for optimality in [62, Theorem 2.1]  
 263 and [23, Theorem 1], without any reference to Peano’s work. The latter theorem uses  
 264 the polar of the closure of the convex hull of the tangent cone which equals the polar of  
 265 the tangent cone by [53, Corollary 6.21]. Neither “stationary” nor “critical” appears  
 266 in [62] or [23].

267 The “Bouligand derivative”, or “B-derivative” for short, was introduced in [52].  
 268 It is a special case of the contingent derivative introduced by Aubin based on the  
 269 tangent cone. The name “Bouligand derivative” was chosen because the tangent  
 270 cone is generally attributed to Bouligand; see, e.g., [53, 41, 42] for recent references.  
 271 Differentiability implies B-differentiability.

272 In [58, §4], a point where a real-valued function is B-differentiable is called a  
 273 “Bouligand stationary (B-stationary) point” of the function if the B-derivative at that

274 point is nonnegative. This is a stationarity concept for unconstrained optimization,  
 275 which therefore does not apply to problem (1.1).

276 B-stationarity is called a “stationarity condition” and said to be “well known” in  
 277 [38, §4.1] where [23] is cited.

278 In [54, §2.1], the term “B-stationarity” is used to name the stationarity concept  
 279 for an MPEC that corresponds to the B-stationarity in the sense of [58, §4] for a  
 280 nonsmooth reformulation of the MPEC [54, Proposition 6]. As pointed out in [65,  
 281 §2.1] and [14, §3.3], the “B-stationarity” in the sense of [54, §2.1], which is specific  
 282 to MPECs, is not B-stationarity in the sense of Definition 1.1 and is called “MPEC-  
 283 linearized B-stationarity” in [14, §3.3] to avoid confusion. Nevertheless, this MPEC-  
 284 linearized B-stationarity appears under the name “B-stationarity” in [28, §1.1], [24,  
 285 Definition 2.2], and [64, Definition 3.2] which all cite [54].

286 The term “B-stationarity” was used to name the absence of descent directions in  
 287 the tangent cone (as in Definition 1.1) first in [48, §1]. It was used in this sense in  
 288 several subsequent works by various authors; see, e.g., [18, §2], [19, §2], [17, Defini-  
 289 tion 2.4], [65, Definition 2.2], [14, §§3.3 and 4], [15, §3], [47, §2], [59, Definition 2.4],  
 290 [21, Definition 3.4], [49, (18)], [7, Definition 3(1)], [8, Definition 4(i)], [27, §4], and [9,  
 291 Definition 6.1.1].

292 In [9, Definition 6.1.1], B-stationarity is defined for the problem of minimizing a  
 293 real-valued function that is B-differentiable on a nonempty closed subset of a Euclid-  
 294 ean vector space, thereby extending the concept introduced in [58, §4] to constrained  
 295 optimization. This more general definition reduces to that from Definition 1.1 if the  
 296 function is differentiable.

297 B-stationarity is also known under other names in the literature. First, in [16,  
 298 Definition 1(b)] and [40, §3], B-stationarity is called “strong stationarity”; [16, prob-  
 299 lem (4)] and [40, (P2)] reduce to problem (1.1) for  $F$  the identity map on  $\mathbb{R}^n$ . Second,  
 300 because the regular normal cone is also called the Fréchet normal cone, especially in  
 301 infinite-dimensional spaces [53, 41, 42], B-stationarity is called “Fréchet stationarity”,  
 302 or “F-stationarity” for short, in [35, Definition 4.1(ii)], [36, Definition 5.1(i)], and [37,  
 303 Definition 3.2(ii)]. Third, B-stationarity is simply called “stationarity” (or “critical-  
 304 ity”) in [55, §2.1], [25, §2.1.1], [32, Definition 2.3], [33, Definition 3.2(c)], and [20,  
 305 Definition 1].

306 **3.2. A brief history of Mordukhovich stationarity.** According to [16, §2],  
 307 the term “M-stationarity” was introduced in [56] for an MPEC. This name was chosen  
 308 because the corresponding stationarity condition was derived from the generalized  
 309 differential calculus of Mordukhovich. To the best of our knowledge, the term “M-  
 310 stationarity” was used to indicate that the negative gradient is in the normal cone (as  
 311 in Definition 1.1) first in [16, Definition 1(a)]; recall that [16, problem (4)] reduces to  
 312 problem (1.1) for  $F$  the identity map on  $\mathbb{R}^n$ . There, the name is motivated by the  
 313 presence of the normal cone which was introduced by Mordukhovich. M-stationarity  
 314 appears, under this name, in several subsequent works by various authors; see, e.g.,  
 315 [7, Definition 3(3)], [8, Definition 4(iii)], [27, §4], [40, §3], [31, §2], [30, §3], and [29,  
 316 §2.3].

317 **3.3. Proximal stationarity and  $\alpha$ -stationarity.** P-stationarity is related to  
 318  $\alpha$ -stationarity which was introduced in [5, Definition 2.3] for  $C = \mathbb{R}_{\leq s}^n$  and in [35,  
 319 Definition 4.1(i)], [25, §2.1.1], [36, Definition 5.1(ii)], [34, (4.2)], and [37, Defini-  
 320 tion 3.2(i)] for several low-rank sets. By definition of the proximal normal cone,  
 321 a point  $x \in C$  is P-stationary for (1.1) if and only if there exists  $\alpha \in (0, \infty)$  such  
 322 that  $x \in P_C(x - \alpha \nabla f(x))$ . In contrast, given  $\alpha \in (0, \infty)$ , a point  $x \in C$  is said to be



323  $\alpha$ -stationary for (1.1) if  $x \in P_C(x - \alpha \nabla f(x))$ . Thus, while  $\alpha$ -stationarity prescribes  
 324 the number  $\alpha \in (0, \infty)$ , P-stationarity merely requires the existence of such a number.  
 325 Furthermore,  $\alpha$ -stationarity should not be confused with the approximate stationarity  
 326 from [32, Definition 2.6].

327 **4. The PGD algorithm.** This section reviews the PGD algorithm, as defined  
 328 in [30, Algorithm 3.1] except that the “average” rule is allowed as an alternative to  
 329 the “max” rule. Its iteration map, called the PGD map, is defined as Algorithm 4.1.  
 330 PGD is defined as Algorithm 4.2 which uses Algorithm 4.1 as a subroutine. The  
 331 nonmonotonic behavior of PGD is described in Propositions 4.6 and 4.7.

---

**Algorithm 4.1** PGD map

---

**Require:**  $(\mathcal{E}, C, f, \underline{\alpha}, \bar{\alpha}, \beta, c)$  where  $\mathcal{E}$  is a Euclidean vector space,  $C$  is a nonempty  
 closed subset of  $\mathcal{E}$ ,  $f : \mathcal{E} \rightarrow \mathbb{R}$  is differentiable on  $C$ ,  $0 < \underline{\alpha} \leq \bar{\alpha} < \infty$ , and  
 $\beta, c \in (0, 1)$ .

**Input:**  $(x, \mu)$  with  $x \in C$  and  $\mu \in [f(x), \infty)$ .

**Output:**  $y \in \text{PGD}(x, \mu; \mathcal{E}, C, f, \underline{\alpha}, \bar{\alpha}, \beta, c)$ .

- 1: Choose  $\alpha \in [\underline{\alpha}, \bar{\alpha}]$  and  $y \in P_C(x - \alpha \nabla f(x))$ ;
  - 2: **while**  $f(y) > \mu + c \langle \nabla f(x), y - x \rangle$  **do**
  - 3:      $\alpha \leftarrow \alpha \beta$ ;
  - 4:     Choose  $y \in P_C(x - \alpha \nabla f(x))$ ;
  - 5: **end while**
  - 6: Return  $y$ .
- 

332 *Remark 4.1.* The Armijo condition

333 
$$f(y) \leq \mu + c \langle \nabla f(x), y - x \rangle$$

334 ensures that the decrease  $\mu - f(y)$  is at least a fraction  $c$  of the opposite of the  
 335 directional derivative of  $f$  at  $x$  along the update vector  $y - x$ . By (2.2), this condition  
 336 implies that

337 (4.1) 
$$f(y) \leq \mu - \frac{c}{2\alpha} \|y - x\|^2,$$

338 which is the condition used in [31, Algorithms 3.1 and 4.1] and [11, Algorithm 3.1].  
 339 Importantly, all results from [31] hold for both conditions, as is clear from the proofs.

340 *Remark 4.2.* By Proposition 5.3, if  $f$  satisfies (H1) and  $x$  is not B-stationary  
 341 for (1.1), then the while loop in Algorithm 4.1 is guaranteed to terminate, thereby  
 342 producing a point  $y$  such that  $f(y) < \mu$ ;  $y \neq x$  holds because  $x$  is not B-stationary  
 343 and hence not P-stationary. If  $f$  satisfies (H2), then the while loop is guaranteed to  
 344 terminate for every  $x \in C$ , by Corollary 6.2.

345 The PGD algorithm is defined as Algorithm 4.2. It is said to be monotone or  
 346 nonmonotone depending on whether  $\mu_i = f(x_i)$  for all  $i$  (that is,  $l = 0$  for the “max”  
 347 rule, or  $p = 1$  for the “average” rule) or not.

**Algorithm 4.2** PGD

**Require:**  $(\mathcal{E}, C, f, \underline{\alpha}, \bar{\alpha}, \beta, c, \text{“nonmonotonicity”})$  where  $\mathcal{E}$  is a Euclidean vector space,  $C$  is a nonempty closed subset of  $\mathcal{E}$ ,  $f : \mathcal{E} \rightarrow \mathbb{R}$  is differentiable on  $C$ ,  $0 < \underline{\alpha} \leq \bar{\alpha} < \infty$ ,  $\beta, c \in (0, 1)$ , and “nonmonotonicity”  $\in \{(\text{“max”}, l), (\text{“average”}, p)\}$  with  $l \in \mathbb{N}$  and  $p \in (0, 1]$ .

**Input:**  $x_0 \in C$ .

**Output:** a sequence in  $C$ .

```

1:  $i \leftarrow 0$ ;
2:  $\mu_{-1} \leftarrow f(x_0)$ ;
3: while  $-\nabla f(x_i) \notin \widehat{N}_C(x_i)$  do
4:   if “nonmonotonicity” = (‘‘max’’,  $l$ ) then
5:      $\mu_i \leftarrow \max_{j \in \{\max\{0, i-l\}, \dots, i\}} f(x_j)$ ;
6:   else if “nonmonotonicity” = (‘‘average’’,  $p$ ) then
7:      $\mu_i \leftarrow (1-p)\mu_{i-1} + pf(x_i)$ ;
8:   end if
9:   Choose  $x_{i+1} \in \text{PGD}(x_i, \mu_i; \mathcal{E}, C, f, \underline{\alpha}, \bar{\alpha}, \beta, c)$ ;
10:   $i \leftarrow i + 1$ ;
11: end while

```

348 *Remark 4.3.* For simplicity, we use a constant weight  $p$  in the ‘‘average’’ rule.  
349 However, we could allow the weight to change from one iteration to the other. It  
350 would then be denoted by  $p_i$ . The main results of the article would hold true in this  
351 more general setting, under the additional assumption that  $\inf_{i \in \mathbb{N}} p_i > 0$ .

352 *Remark 4.4.* If  $f$  satisfies (H2), then  $\widehat{N}_C(x_i)$  should be replaced with  $\widehat{\widehat{N}}_C(x_i)$  in  
353 line 3.

354 Examples of a set  $C$  for which the projection map and the regular and proximal  
355 normal cones can be described explicitly abound in the literature; see Section 7. For  
356 such examples, Algorithm 4.2 can be practically implemented.

357 *Remark 4.5.* From Remark 4.2, under (H1), the call to Algorithm 4.1 in line 9  
358 of PGD never results in an infinite loop. Consequently, by running PGD, one always  
359 encounters one of the following two situations:

- 360 • PGD generates a finite sequence, and the last element of this sequence is  
361 B-stationary for (1.1) if  $f$  satisfies (H1), and even P-stationary for (1.1) if  $f$   
362 satisfies (H2);
- 363 • PGD generates an infinite sequence.

364 The rest of this section and the next two concern the nontrivial case where PGD  
365 generates an infinite sequence. In that case, the stationarity of the accumulation  
366 points of the generated sequence, if any, is studied in Sections 5 and 6. Following  
367 [51, Remark 14], which states that it is usually better to determine whether an al-  
368 gorithm generates a sequence having at least one accumulation point by examining  
369 the algorithm in the light of the specific problem to which one wishes to apply it,  
370 no condition ensuring the existence of a convergent subsequence is imposed. As a  
371 reminder, a sequence  $(x_i)_{i \in \mathbb{N}}$  in  $\mathcal{E}$  has at least one accumulation point if and only if  
372  $\liminf_{i \rightarrow \infty} \|x_i\| < \infty$ .

373 A property of monotone PGD that is helpful for its analysis is the fact that  $f$  is  
374 strictly decreasing along the generated sequence. For nonmonotone PGD, this is not  
375 true. However, weaker properties, stated in the following two propositions, will be  
376 enough for our purposes.

377 PROPOSITION 4.6. Let  $(x_i)_{i \in \mathbb{N}}$  be a sequence generated by PGD (Algorithm 4.2)  
 378 using the “max” rule. For every  $i \in \mathbb{N}$ , let  $g(i) \in \operatorname{argmax}_{j \in \{\max\{0, i-l\}, \dots, i\}} f(x_j)$ .  
 379 Then:

- 380 1.  $(f(x_{g(i)}))_{i \in \mathbb{N}}$  is monotonically nonincreasing;
- 381 2.  $(x_i)_{i \in \mathbb{N}}$  is contained in the sublevel set (1.6);
- 382 3. if  $x \in C$  is an accumulation point of  $(x_i)_{i \in \mathbb{N}}$ , then  $(f(x_{g(i)}))_{i \in \mathbb{N}}$  converges to  
 383  $\varphi \in [f(x), f(x_0)]$ ;
- 384 4. if  $f$  is bounded from below and uniformly continuous on a set that contains  
 385  $(x_i)_{i \in \mathbb{N}}$ , then  $(f(x_i))_{i \in \mathbb{N}}$  converges to  $\varphi \in \mathbb{R}$ .

386 *Proof.* The first two statements are [31, Lemma 4.1 and Corollary 4.1]. For  
 387 the third one, let  $(x_{i_k})_{k \in \mathbb{N}}$  be a subsequence converging to  $x$ . Since the sequence  
 388  $(f(x_{g(i)}))_{i \in \mathbb{N}}$  is monotonically nonincreasing, it has a limit in  $\mathbb{R} \cup \{-\infty\}$ . Thus,

$$389 \quad \lim_{i \rightarrow \infty} f(x_{g(i)}) = \lim_{k \rightarrow \infty} f(x_{g(i_k)}) \geq \liminf_{k \rightarrow \infty} f(x_{i_k}) = f(x) > -\infty.$$

390 It remains to prove the fourth statement. From the first statement, and because  $f$   
 391 is bounded from below,  $(f(x_{g(i)}))_{i \in \mathbb{N}}$  converges to some limit  $\varphi \in \mathbb{R}$ . Assume, for  
 392 the sake of contradiction, that  $(f(x_i))_{i \in \mathbb{N}}$  does not converge to  $\varphi$ . Then, there exist  
 393  $\rho \in (0, \infty)$  and a subsequence  $(f(x_{i_j}))_{j \in \mathbb{N}}$  contained in  $\mathbb{R} \setminus [\varphi - \rho, \varphi + \rho]$ . For all  
 394  $j \in \mathbb{N}$ , define  $p_j := g(i_j + l) - i_j \in \{0, \dots, l\}$ . Then, there exist  $p \in \{0, \dots, l\}$   
 395 and a subsequence  $(p_{j_k})_{k \in \mathbb{N}}$  such that, for all  $k \in \mathbb{N}$ ,  $p_{j_k} = p$ . By [31, (27)] or [30,  
 396 (A.9)],  $(f(x_{g(i) - p}))_{i \in \mathbb{N}}$  converges to  $\varphi$ . Therefore,  $(f(x_{g(i+l) - p}))_{i \in \mathbb{N}}$  converges to  $\varphi$ .  
 397 Hence,  $(f(x_{g(i_{j_k} + l) - p}))_{k \in \mathbb{N}}$  converges to  $\varphi$ . This is a contradiction since, for all  $k \in \mathbb{N}$ ,  
 398  $f(x_{g(i_{j_k} + l) - p}) = f(x_{i_{j_k}})$ .  $\square$

399 PROPOSITION 4.7. Let  $(x_i)_{i \in \mathbb{N}}$  be a sequence generated by PGD (Algorithm 4.2)  
 400 using the “average” rule. Then:

- 401 1.  $(x_i)_{i \in \mathbb{N}}$  is contained in the sublevel set (1.6);
- 402 2. if  $(x_i)_{i \in \mathbb{N}}$  has an accumulation point, then  $(f(x_i))_{i \in \mathbb{N}}$  and  $(\mu_i)_{i \in \mathbb{N}}$  converge,  
 403 toward the same (finite) value.

404 *Proof.* The sequence  $(\mu_i)_{i \in \mathbb{N}}$  is monotonically nonincreasing since, for all  $i \in \mathbb{N}$ ,  
 405  $f(x_i) \leq \mu_{i-1}$ , hence  $\mu_i = (1 - p)\mu_{i-1} + pf(x_i) \leq \mu_{i-1}$ . Therefore, for all  $i \in \mathbb{N}$ ,

$$406 \quad f(x_i) \leq \mu_{i-1} \leq \mu_{-1} = f(x_0),$$

407 meaning that  $(x_i)_{i \in \mathbb{N}}$  is contained in the sublevel set (1.6).

408 Now, we prove the second item of the proposition. Let us assume that  $(x_i)_{i \in \mathbb{N}}$  has  
 409 an accumulation point  $x$ . Let  $(x_{i_k})_{k \in \mathbb{N}}$  be a subsequence converging to  $x$ . Observe  
 410 that

$$411 \quad \lim_{k \rightarrow \infty} f(x_{i_k}) = f(x),$$

412 since  $f$  is differentiable, and in particular continuous, at  $x$ . As  $(\mu_i)_{i \in \mathbb{N}}$  is monotonically  
 413 nonincreasing, it has a limit  $\varphi \in \mathbb{R} \cup \{-\infty\}$ . For all  $k \in \mathbb{N}$ ,

$$414 \quad f(x_{i_k}) \leq \mu_{i_k - 1}.$$

415 Letting  $k$  tend to infinity yields

$$416 \quad f(x) = \lim_{k \rightarrow \infty} f(x_{i_k}) \leq \lim_{k \rightarrow \infty} \mu_{i_k - 1} = \varphi.$$

417 In particular,  $\varphi$  is finite.

418 Now, we show that  $\varphi = \liminf_{i \rightarrow \infty} f(x_i)$ . Let  $(x_{j_k})_{k \in \mathbb{N}}$  be a subsequence such  
419 that

$$420 \quad \lim_{k \rightarrow \infty} f(x_{j_k}) = \liminf_{i \rightarrow \infty} f(x_i).$$

421 For all  $k \in \mathbb{N}$ , it holds that

$$422 \quad \mu_{j_k} = (1 - p)\mu_{j_k - 1} + pf(x_{j_k}).$$

423 The two sides of this equality must have the same limit:

$$424 \quad \varphi = (1 - p)\varphi + p \liminf_{i \rightarrow \infty} f(x_i).$$

425 As  $p > 0$ , this implies  $\varphi = \liminf_{i \rightarrow \infty} f(x_i)$  (and, in particular,  $\liminf_{i \rightarrow \infty} f(x_i) >$   
426  $-\infty$ ). To conclude, we observe that, for all  $k \in \mathbb{N}$ ,

$$427 \quad f(x_k) \leq \mu_{k-1}.$$

428 Hence,

$$429 \quad \limsup_{k \rightarrow \infty} f(x_k) \leq \lim_{k \rightarrow \infty} \mu_{k-1} = \varphi = \liminf_{k \rightarrow \infty} f(x_k).$$

430 Therefore,  $(f(x_k))_{k \in \mathbb{N}}$  converges to  $\varphi$ . □

431 **5. Convergence analysis for a continuous gradient.** In this section, PGD  
432 (Algorithm 4.2) is analyzed under hypothesis (H1). As mentioned after Remark 4.5,  
433 only the nontrivial case where an infinite sequence is generated is considered here.  
434 Specifically, the first part of the second item of Theorem 1.2, restated in Theorem 5.1  
435 for convenience, is proven.

436 **THEOREM 5.1.** *Let  $(x_i)_{i \in \mathbb{N}}$  be a sequence generated by PGD (Algorithm 4.2). If  
437  $f$  satisfies (H1), then all accumulation points of  $(x_i)_{i \in \mathbb{N}}$  are B-stationary for (1.1).  
438 If, moreover,  $(x_i)_{i \in \mathbb{N}}$  has an isolated accumulation point, then  $(x_i)_{i \in \mathbb{N}}$  converges.*

439 The proof is divided into three parts. First, in Section 5.1, we show that, in a  
440 neighborhood of every point that is not B-stationary for (1.1), the PGD map (Al-  
441 gorithm 4.1) terminates after a bounded number of iterations. Then, in Section 5.2,  
442 we prove that, if a subsequence  $(x_{i_k})_{k \in \mathbb{N}}$  converges, then  $(x_{i_k+1})_{k \in \mathbb{N}}$  also does, to the  
443 same limit. Finally, we combine the first two parts in Section 5.3: roughly, if  $(x_{i_k})_{k \in \mathbb{N}}$   
444 converges to  $x$ , then, from the second part,

$$445 \quad \|x_{i_k+1} - x_{i_k}\| \rightarrow 0 \quad \text{when } k \rightarrow \infty,$$

446 but, from the first part, if  $x$  is not B-stationary for (1.1), then the iterates of PGD  
447 move by at least a constant amount at each iteration. It is therefore impossible that  
448  $(x_{i_k})_{k \in \mathbb{N}}$  converges to a point that is not B-stationary for (1.1).

449 **5.1. First part: analysis of the PGD map.** In this section, we show that, if  
450  $\underline{x} \in C$  is not B-stationary for (1.1), then the while loop in Algorithm 4.1 terminates,  
451 in some neighborhood of  $\underline{x}$ , for nonvanishing values of  $\alpha$ . The intuition for this proof  
452 is that, for every  $x$  close to  $\underline{x}$  and for every  $y \in P_C(x - \alpha \nabla f(x))$ ,

$$453 \quad f(y) = f(x) + \langle \nabla f(x), y - x \rangle + \text{some remainder.}$$

454 The inner product  $\langle \nabla f(x), y - x \rangle$  is negative, and larger in absolute value than some  
 455 fraction of  $\|\nabla f(x)\| \|y - x\|$  (Proposition 5.2). On the other hand, if  $\alpha$  is small enough,  
 456 the remainder (upper bounded in Proposition 5.3) is smaller than some arbitrarily  
 457 small fraction of  $\|\nabla f(x)\| \|y - x\|$ . Therefore, for  $\alpha$  small enough,

$$458 \quad f(y) < f(x) + c \langle \nabla f(x), y - x \rangle.$$

459 PROPOSITION 5.2. *Assume that  $f$  satisfies (H1). Let  $\underline{x} \in C$  be non-B-stationary*  
 460 *for (1.1), and  $w \in T_C(\underline{x})$  be such that*

$$461 \quad (5.1) \quad \langle w, -\nabla f(\underline{x}) \rangle > 0.$$

462 Define  $\kappa := \sqrt{1 - \frac{\beta \langle w, -\nabla f(\underline{x}) \rangle^2}{8 \|w\|^2 \|\nabla f(\underline{x})\|^2}} \in (0, 1)$ . For every  $\varepsilon \in (0, \infty)$ , there exist  $\alpha_{\underline{x}} \in (0, \varepsilon]$   
 463 and  $\bar{\rho}(\alpha_{\underline{x}}) \in (0, \infty)$  such that, for all  $x \in B(\underline{x}, \bar{\rho}(\alpha_{\underline{x}})) \cap C$  and  $\alpha \in [\alpha_{\underline{x}}, \alpha_{\underline{x}}/\beta]$ ,

$$464 \quad d(x - \alpha \nabla f(x), C) \leq \kappa \alpha \|\nabla f(x)\|,$$

465 which implies, for all  $y \in P_C(x - \alpha \nabla f(x))$ ,

$$466 \quad \langle \nabla f(x), y - x \rangle \leq -\sqrt{1 - \kappa^2} \|\nabla f(x)\| \|y - x\|.$$

467 *Proof.* Let  $\varepsilon \in (0, \infty)$  be fixed. We show that there exist  $\alpha_{\underline{x}} \in (0, \varepsilon]$  and  $\bar{\rho}(\alpha_{\underline{x}}) \in$   
 468  $(0, \infty)$  satisfying the required property.

469 Let  $(w_i)_{i \in \mathbb{N}}$  be a sequence in  $C$  converging to  $\underline{x}$ , and  $(t_i)_{i \in \mathbb{N}}$  be a sequence in  
 470  $(0, \infty)$  such that

$$471 \quad \frac{w_i - \underline{x}}{t_i} \xrightarrow{i \rightarrow \infty} w.$$

472 From the definition of  $w$  in (5.1), it holds for all  $i \in \mathbb{N}$  large enough that

$$473 \quad (5.2) \quad \langle w_i - \underline{x}, -\nabla f(\underline{x}) \rangle > 0.$$

474 As  $\frac{1}{t_i} \frac{\|w_i - \underline{x}\|^2}{\langle w_i - \underline{x}, -\nabla f(\underline{x}) \rangle} \xrightarrow{i \rightarrow \infty} \frac{\|w\|^2}{\langle w, -\nabla f(\underline{x}) \rangle}$  and  $t_i \xrightarrow{i \rightarrow \infty} 0$ , it also holds for all  $i \in \mathbb{N}$  large  
 475 enough that

$$476 \quad (5.3) \quad \frac{\|w_i - \underline{x}\|^2}{\langle w_i - \underline{x}, -\nabla f(\underline{x}) \rangle} < \varepsilon.$$

477 Similarly, it holds for all  $i \in \mathbb{N}$  large enough that

$$478 \quad (5.4) \quad \frac{\langle w_i - \underline{x}, -\nabla f(\underline{x}) \rangle^2}{\|w_i - \underline{x}\|^2} > \frac{\langle w, -\nabla f(\underline{x}) \rangle^2}{2 \|w\|^2}.$$

479 Fix  $i \in \mathbb{N}$  satisfying (5.2), (5.3), and (5.4). Pick  $\alpha_{\underline{x}}$  such that

$$480 \quad \frac{\alpha_{\underline{x}}}{2} < \frac{\|w_i - \underline{x}\|^2}{\langle w_i - \underline{x}, -\nabla f(\underline{x}) \rangle} < \alpha_{\underline{x}} < \varepsilon.$$

481 Since  $\nabla f$  is continuous at  $\underline{x}$ , there exists  $\rho_0 \in (0, \infty)$  such that, for all  $x \in B[\underline{x}, \rho_0] \cap C$ ,

482

$$483 \quad (5.5a) \quad \langle w_i - \underline{x}, -\nabla f(x) \rangle > 0,$$

$$484 \quad (5.5b) \quad \frac{\alpha_{\underline{x}}}{2} < \frac{\|w_i - \underline{x}\|^2}{\langle w_i - \underline{x}, -\nabla f(x) \rangle} < \alpha_{\underline{x}},$$

$$485 \quad (5.5c) \quad \frac{\langle w_i - \underline{x}, -\nabla f(x) \rangle^2}{\|w_i - \underline{x}\|^2 \|\nabla f(x)\|^2} > \frac{\langle w, -\nabla f(\underline{x}) \rangle^2}{2\|w\|^2 \|\nabla f(\underline{x})\|^2}.$$

486 We now establish the first inequality we have to prove: for an adequate value of  $\bar{\rho}(\alpha_{\underline{x}})$ ,  
487 it holds for all  $x \in B(\underline{x}, \bar{\rho}(\alpha_{\underline{x}})) \cap C$  and  $\alpha \in [\alpha_{\underline{x}}, \alpha_{\underline{x}}/\beta]$  that

$$488 \quad \|x - \alpha \nabla f(x) - y\| \leq \kappa \alpha \|\nabla f(x)\|, \quad \forall y \in P_C(x - \alpha \nabla f(x)),$$

489 which is equivalent to  $d(x - \alpha \nabla f(x), C) \leq \kappa \alpha \|\nabla f(x)\|$ .

490 Let us for the moment consider any  $\bar{\rho}(\alpha_{\underline{x}}) \in (0, \rho_0]$ . For all  $x \in B(\underline{x}, \bar{\rho}(\alpha_{\underline{x}})) \cap C$ ,  
491  $\alpha \in [\alpha_{\underline{x}}, \alpha_{\underline{x}}/\beta]$ , and  $y \in P_C(x - \alpha \nabla f(x))$ ,

$$\begin{aligned} 492 \quad \|x - \alpha \nabla f(x) - y\|^2 &\leq \|x - \alpha \nabla f(x) - w_i\|^2 \\ 493 \quad &= \|\underline{x} - \alpha \nabla f(x) - w_i\|^2 + 2 \langle \underline{x} - x, \alpha \nabla f(x) + w_i - \underline{x} \rangle + \|\underline{x} - x\|^2 \\ 494 \quad &\leq \|\underline{x} - \alpha \nabla f(x) - w_i\|^2 \\ 495 \quad &\quad + 2\bar{\rho}(\alpha_{\underline{x}}) (\alpha \|\nabla f(x)\| + \|w_i - \underline{x}\|) + \bar{\rho}(\alpha_{\underline{x}})^2 \\ 496 \quad &\leq \|\underline{x} - \alpha \nabla f(x) - w_i\|^2 \\ 497 \quad &\quad + 2\bar{\rho}(\alpha_{\underline{x}}) \left( \alpha \max_{z \in B[\underline{x}, \rho_0] \cap C} \|\nabla f(z)\| + \|w_i - \underline{x}\| \right) + \bar{\rho}(\alpha_{\underline{x}})^2 \\ 498 \quad &= \alpha^2 \|\nabla f(x)\|^2 - 2\alpha \langle w_i - \underline{x}, -\nabla f(x) \rangle + \|w_i - \underline{x}\|^2 \\ 499 \quad &\quad + 2\bar{\rho}(\alpha_{\underline{x}}) \left( \alpha \max_{z \in B[\underline{x}, \rho_0] \cap C} \|\nabla f(z)\| + \|w_i - \underline{x}\| \right) + \bar{\rho}(\alpha_{\underline{x}})^2 \\ 500 \quad &\leq \alpha^2 \|\nabla f(x)\|^2 - \alpha \langle w_i - \underline{x}, -\nabla f(x) \rangle \\ 501 \quad &\quad + 2\bar{\rho}(\alpha_{\underline{x}}) \left( \frac{\alpha_{\underline{x}}}{\beta} \max_{z \in B[\underline{x}, \rho_0] \cap C} \|\nabla f(z)\| + \|w_i - \underline{x}\| \right) + \bar{\rho}(\alpha_{\underline{x}})^2 \end{aligned}$$

502 where the last inequality follows from (5.5b) and the fact that  $\alpha_{\underline{x}} \leq \alpha \leq \frac{\alpha_{\underline{x}}}{\beta}$ . Choose  
503  $\bar{\rho}(\alpha_{\underline{x}}) \in (0, \rho_0]$  small enough to ensure

$$\begin{aligned} 504 \quad &2\bar{\rho}(\alpha_{\underline{x}}) \left( \frac{\alpha_{\underline{x}}}{\beta} \max_{z \in B[\underline{x}, \rho_0] \cap C} \|\nabla f(z)\| + \|w_i - \underline{x}\| \right) + \bar{\rho}(\alpha_{\underline{x}})^2 \\ 505 \quad &\leq \frac{\alpha_{\underline{x}}}{2} \min_{z \in B[\underline{x}, \rho_0] \cap C} \langle w_i - \underline{x}, -\nabla f(z) \rangle. \end{aligned}$$

506 Note that the right-hand side of this inequality is positive, from (5.5a). Combining

507 this definition with the previous inequality, we arrive at

$$\begin{aligned}
508 \quad \|x - \alpha \nabla f(x) - y\|^2 &\leq \alpha^2 \|\nabla f(x)\|^2 - \frac{\alpha}{2} \langle w_i - \underline{x}, -\nabla f(x) \rangle \\
509 \quad &= \alpha^2 \|\nabla f(x)\|^2 \left( 1 - \frac{\langle w_i - \underline{x}, -\nabla f(x) \rangle}{2\alpha \|\nabla f(x)\|^2} \right) \\
510 \quad &\leq \alpha^2 \|\nabla f(x)\|^2 \left( 1 - \frac{\beta \langle w_i - \underline{x}, -\nabla f(x) \rangle}{2\alpha_{\underline{x}} \|\nabla f(x)\|^2} \right) \text{ as } \alpha \leq \frac{\alpha_{\underline{x}}}{\beta} \\
511 \quad &\leq \alpha^2 \|\nabla f(x)\|^2 \left( 1 - \frac{\beta \langle w_i - \underline{x}, -\nabla f(x) \rangle^2}{4\|w_i - \underline{x}\|^2 \|\nabla f(x)\|^2} \right) \text{ from (5.5b)} \\
512 \quad &\leq \alpha^2 \|\nabla f(x)\|^2 \left( 1 - \frac{\beta \langle w, -\nabla f(\underline{x}) \rangle^2}{8\|w\|^2 \|\nabla f(\underline{x})\|^2} \right) \text{ from (5.5c)} \\
513 \quad &= \kappa^2 \alpha^2 \|\nabla f(x)\|^2.
\end{aligned}$$

514 In other words, for all  $x \in B(\underline{x}, \bar{\rho}(\alpha_{\underline{x}})) \cap C$ ,  $\alpha \in [\alpha_{\underline{x}}, \alpha_{\underline{x}}/\beta]$ , and  $y \in P_C(x - \alpha \nabla f(x))$ ,  
515 it holds that

$$516 \quad \|x - \alpha \nabla f(x) - y\| \leq \kappa \alpha \|\nabla f(x)\|.$$

517 To conclude, we show that this inequality implies

$$518 \quad (5.6) \quad \left\langle \frac{y - x}{\|y - x\|}, \frac{\nabla f(x)}{\|\nabla f(x)\|} \right\rangle \leq -\sqrt{1 - \kappa^2}.$$

519 Indeed, if we define  $\theta \in \mathbb{R}$  such that  $\left\langle \frac{y - x}{\|y - x\|}, \frac{\nabla f(x)}{\|\nabla f(x)\|} \right\rangle = \cos(\theta)$ , we have

$$520 \quad \|y - x\|^2 + 2\alpha \|\nabla f(x)\| \|y - x\| \cos(\theta) + \alpha^2 \|\nabla f(x)\|^2 \leq \alpha^2 \kappa^2 \|\nabla f(x)\|^2.$$

521 This already shows that  $\cos(\theta) < 0$ . In addition, if we minimize the left-hand side  
522 over all possible values of  $\|y - x\|$ , we get

$$523 \quad -\alpha^2 \|\nabla f(x)\|^2 \cos^2(\theta) + \alpha^2 \|\nabla f(x)\|^2 \leq \alpha^2 \kappa^2 \|\nabla f(x)\|^2,$$

524 hence  $\cos^2(\theta) \geq 1 - \kappa^2$ , which establishes (5.6).  $\square$

525 **PROPOSITION 5.3.** *Let  $\underline{\alpha} \in (0, \infty)$  and  $c \in (0, 1)$ . Assume that  $f$  satisfies (H1).  
526 Let  $\underline{x} \in C$  be non-B-stationary for (1.1). There exists  $\alpha_{\underline{x}} \in (0, \underline{\alpha}]$  and  $\rho \in (0, \infty)$  such  
527 that, for all  $x \in B(\underline{x}, \rho) \cap C$ ,  $\alpha \in [\alpha_{\underline{x}}, \alpha_{\underline{x}}/\beta]$ , and  $y \in P_C(x - \alpha \nabla f(x))$ ,*

$$528 \quad f(y) < f(x) + c \langle \nabla f(x), y - x \rangle.$$

529 *Proof.* Define  $\kappa$  as in Proposition 5.2. Let  $\delta \in (0, \infty)$  be small enough to ensure

$$530 \quad (5.7a) \quad \sup_{y \in B[\underline{x}, \frac{7\delta}{2\beta} \|\nabla f(\underline{x})\|] \cap C \setminus \{x\}} \frac{|f(y) - f(\underline{x}) - \langle \nabla f(\underline{x}), y - \underline{x} \rangle|}{\|y - \underline{x}\|} < \frac{(1 - c)\sqrt{1 - \kappa^2} \|\nabla f(\underline{x})\|}{4 \left(1 + \frac{8}{3(1 - \kappa)}\right)},$$

$$531 \quad (5.7b) \quad \sup_{y \in B[\underline{x}, \frac{7\delta}{2\beta} \|\nabla f(\underline{x})\|] \cap C} \|\nabla f(y) - \nabla f(\underline{x})\| < \frac{(1 - c)\sqrt{1 - \kappa^2}}{4} \|\nabla f(\underline{x})\|.$$

532 These inequalities are satisfied by all  $\delta$  small enough, from the definition of the gra-  
 533 dient for the first one, and because the gradient is continuous at  $\underline{x}$  for the second  
 534 one.

535 Then, define  $\varepsilon := \min\{\underline{\alpha}, \delta\}$  and let  $\alpha_{\underline{x}} \in (0, \varepsilon]$  and  $\bar{\rho}(\alpha_{\underline{x}}) \in (0, \infty)$  be as in  
 536 Proposition 5.2. Define

$$537 \quad \rho := \min\{\bar{\rho}(\alpha_{\underline{x}}), \alpha_{\underline{x}}\|\nabla f(\underline{x})\|\}.$$

538 Note that, for all  $x \in B(\underline{x}, \rho) \cap C$ ,

$$539 \quad \|x - \underline{x}\| < \rho \leq \alpha_{\underline{x}}\|\nabla f(\underline{x})\| < \frac{7\alpha_{\underline{x}}}{2\beta}\|\nabla f(\underline{x})\| \leq \frac{7\delta}{2\beta}\|\nabla f(\underline{x})\|,$$

540 so that from (5.7b),  $\|\nabla f(x) - \nabla f(\underline{x})\| < \frac{\|\nabla f(\underline{x})\|}{4}$ , which implies

$$\begin{aligned} 541 \quad & \frac{3}{4}\|\nabla f(\underline{x})\| < \|\nabla f(\underline{x})\| - \|\nabla f(x) - \nabla f(\underline{x})\| \\ 542 \quad & \leq \|\nabla f(x)\| \\ 543 \quad & \leq \|\nabla f(\underline{x})\| + \|\nabla f(x) - \nabla f(\underline{x})\| \\ 544 \quad (5.8) \quad & < \frac{5}{4}\|\nabla f(\underline{x})\|. \end{aligned}$$

545 For all  $x \in B(\underline{x}, \rho) \cap C$ ,  $\alpha \in [\alpha_{\underline{x}}, \alpha_{\underline{x}}/\beta]$ , and  $y \in P_C(x - \alpha\nabla f(x))$ ,

$$\begin{aligned} 546 \quad & f(y) = f(x) + \langle \nabla f(\underline{x}), y - x \rangle \\ 547 \quad & \quad + (f(\underline{x}) - f(x) - \langle \nabla f(\underline{x}), \underline{x} - x \rangle) \\ 548 \quad & \quad + (f(y) - f(\underline{x}) - \langle \nabla f(\underline{x}), y - \underline{x} \rangle) \\ 549 \quad (5.9) \quad & \leq f(x) + \langle \nabla f(\underline{x}), y - x \rangle + \frac{(1-c)\sqrt{1-\kappa^2}\|\nabla f(\underline{x})\|}{4\left(1 + \frac{8}{3(1-\kappa)}\right)} (\|\underline{x} - x\| + \|y - \underline{x}\|). \end{aligned}$$

550 The last inequality follows from (5.7a); observe that

$$\begin{aligned} 551 \quad & \|y - \underline{x}\| \leq \|y - x\| + \|x - \underline{x}\| \\ 552 \quad & \leq 2\alpha\|\nabla f(x)\| + \rho \text{ from (2.1)} \\ 553 \quad & \leq \frac{2\alpha_{\underline{x}}}{\beta}\|\nabla f(x)\| + \alpha_{\underline{x}}\|\nabla f(\underline{x})\| \\ 554 \quad & < \frac{5\alpha_{\underline{x}}}{2\beta}\|\nabla f(\underline{x})\| + \alpha_{\underline{x}}\|\nabla f(\underline{x})\| \text{ from (5.8)} \\ 555 \quad & = \frac{7\alpha_{\underline{x}}}{2\beta}\|\nabla f(\underline{x})\|. \end{aligned}$$



556 We continue from (5.9):

$$\begin{aligned}
557 \quad f(y) &\leq f(x) + \langle \nabla f(\underline{x}), y - x \rangle + \frac{(1-c)\sqrt{1-\kappa^2}\|\nabla f(\underline{x})\|}{4\left(1 + \frac{8}{3(1-\kappa)}\right)} (2\|\underline{x} - x\| + \|y - x\|) \\
558 \quad &\stackrel{(a)}{<} f(x) + \langle \nabla f(\underline{x}), y - x \rangle + \frac{(1-c)\sqrt{1-\kappa^2}\|\nabla f(\underline{x})\|}{4} \|y - x\| \\
559 \quad &\leq f(x) + \langle \nabla f(x), y - x \rangle + \|\nabla f(\underline{x}) - \nabla f(x)\| \|y - x\| \\
560 \quad &\quad + \frac{(1-c)\sqrt{1-\kappa^2}\|\nabla f(\underline{x})\|}{4} \|y - x\| \\
561 \quad &\leq f(x) + \langle \nabla f(x), y - x \rangle + \frac{(1-c)\sqrt{1-\kappa^2}\|\nabla f(\underline{x})\|}{2} \|y - x\| \text{ from (5.7b)} \\
562 \quad &< f(x) + \langle \nabla f(x), y - x \rangle + (1-c)\sqrt{1-\kappa^2}\|\nabla f(x)\| \|y - x\| \text{ from (5.8)} \\
563 \quad &\leq f(x) + \langle \nabla f(x), y - x \rangle - (1-c)\langle \nabla f(x), y - x \rangle \text{ from Proposition 5.2} \\
564 \quad &= f(x) + c\langle \nabla f(x), y - x \rangle.
\end{aligned}$$

565 Inequality (a) is true because

$$\begin{aligned}
566 \quad \|y - x\| &\geq \alpha\|\nabla f(x)\| - \|x - \alpha\nabla f(x) - y\| \text{ by the triangle inequality} \\
567 \quad &= \alpha\|\nabla f(x)\| - d(x - \alpha\nabla f(x), C) \\
568 \quad &\geq (1-\kappa)\alpha\|\nabla f(x)\| \text{ from Proposition 5.2} \\
569 \quad &\geq \frac{3}{4}(1-\kappa)\rho \text{ from (5.8), the definition of } \rho, \text{ and } \alpha_{\underline{x}} \leq \alpha \\
570 \quad &> \frac{3}{4}(1-\kappa)\|x - \underline{x}\|. \quad \square
\end{aligned}$$

## 571 5.2. Second part: convergence of successive iterates.

572 PROPOSITION 5.4. *Assume that  $f$  satisfies (H1). Let  $(x_i)_{i \in \mathbb{N}}$  be a sequence gen-*  
573 *erated by PGD (Algorithm 4.2), and  $x$  be an accumulation point. Then, for every*  
574 *subsequence  $(x_{i_k})_{k \in \mathbb{N}}$  converging to  $x$ , the sequence  $(x_{i_k+1})_{k \in \mathbb{N}}$  also converges to  $x$ .*

575 *Proof.* Let  $(x_{i_k})_{k \in \mathbb{N}}$  be a subsequence converging to  $x$ . We show that  $(x_{i_k+1})_{k \in \mathbb{N}}$   
576 also converges to  $x$ .

577 If the nonmonotonicity rule is set to “average”, this is a direct consequence of  
578 Proposition 4.7. Indeed, for all  $i \in \mathbb{N}$ , from (4.1),

$$579 \quad f(x_{i+1}) \leq \mu_i - \frac{c}{2\bar{\alpha}} \|x_{i+1} - x_i\|^2 \leq \mu_i.$$

580 From Proposition 4.7,  $(f(x_{i+1}))_{i \in \mathbb{N}}$  and  $(\mu_i)_{i \in \mathbb{N}}$  converge to the same limit. Therefore,

$$581 \quad \left( \mu_i - \frac{c}{2\bar{\alpha}} \|x_{i+1} - x_i\|^2 \right)_{i \in \mathbb{N}}$$

582 also converges to this limit. This implies that  $(\|x_{i+1} - x_i\|)_{i \in \mathbb{N}}$  converges to 0, hence  
583  $(\|x_{i_k+1} - x_{i_k}\|)_{k \in \mathbb{N}}$  converges to 0, and  $(x_{i_k+1})_{k \in \mathbb{N}}$  converges to the same limit as  
584  $(x_{i_k})_{k \in \mathbb{N}}$ , that is,  $x$ .

585 Now, let us consider the “max” rule case. It suffices to show that  $x$  is an accumula-  
586 tion point of every subsequence of  $(x_{i_k+1})_{k \in \mathbb{N}}$ . In other words, we show the following:  
587 for every subsequence  $(i_{j_k})_{k \in \mathbb{N}}$  of  $(i_k)_{k \in \mathbb{N}}$ , there exists a subsequence of  $(x_{i_{j_k}+1})_{k \in \mathbb{N}}$

588 that converges to  $x$ . Let  $(i_{j_k})_{k \in \mathbb{N}}$  be a subsequence of  $(i_k)_{k \in \mathbb{N}}$ . For all  $i \in \mathbb{N}$ , define  
 589  $g(i) \in \operatorname{argmax}_{j \in \{\max\{0, i-l\}, \dots, i\}} f(x_j)$ , as in Proposition 4.6. By the third statement of  
 590 Proposition 4.6, the sequence  $(f(x_{g(i)}))_{i \in \mathbb{N}}$  converges to  $\varphi \in [f(x), f(x_0)]$ . For every  
 591  $k \in \mathbb{N}$ , letting  $\alpha_{i_{j_k}} \in (0, \bar{\alpha}]$  be the number such that  $x_{i_{j_k}+1} \in P_C(x_{i_{j_k}} - \alpha_{i_{j_k}} \nabla f(x_{i_{j_k}}))$ ,  
 592 by (2.1),

$$593 \quad \|x_{i_{j_k}+1} - x_{i_{j_k}}\| \leq 2\alpha_{i_{j_k}} \|\nabla f(x_{i_{j_k}})\| \leq 2\bar{\alpha} \|\nabla f(x_{i_{j_k}})\|.$$

594 Thus, since  $(x_{i_{j_k}})_{k \in \mathbb{N}}$  is bounded and  $\nabla f$  is locally bounded (as it is continuous), the  
 595 sequence  $(x_{i_{j_k}+1})_{k \in \mathbb{N}}$  is bounded. If we replace  $(i_{j_k})_{k \in \mathbb{N}}$  by a subsequence, we can  
 596 assume that  $(x_{i_{j_k}+1})_{k \in \mathbb{N}}$  converges.

597 Iterating the reasoning, we can assume that  $(x_{i_{j_k}+s})_{k \in \mathbb{N}}$  converges to some  $x^s \in C$   
 598 for every  $s \in \{0, \dots, l+1\}$ . By definition of  $x$ ,  $x^0 = x$ .

599 Observe that, from the continuity of  $f$ ,

$$600 \quad f(x_{g(i_{j_k}+l+1)}) = \max\{f(x_{i_{j_k}+1}), \dots, f(x_{i_{j_k}+l+1})\}$$

$$601 \quad \rightarrow \max\{f(x^1), \dots, f(x^{l+1})\} \text{ when } k \rightarrow \infty.$$

602 In particular, there exists  $s_1 \in \{1, \dots, l+1\}$  such that

$$603 \quad (5.10) \quad f(x^{s_1}) = \varphi.$$

604 Let  $s_1$  be the smallest such integer. For all  $k \in \mathbb{N}$ , from the condition in line 2 of  
 605 Algorithm 4.1 and (4.1),

$$606 \quad f(x_{i_{j_k}+s_1}) \leq f(x_{g(i_{j_k}+s_1-1)}) - \frac{c}{2\bar{\alpha}} \|x_{i_{j_k}+s_1} - x_{i_{j_k}+s_1-1}\|^2.$$

607 Letting  $k$  tend to infinity yields

$$608 \quad \varphi = f(x^{s_1}) \leq \varphi - \frac{c}{2\bar{\alpha}} \|x^{s_1} - x^{s_1-1}\|^2.$$

609 Consequently,  $x^{s_1} = x^{s_1-1}$ . In particular,  $f(x^{s_1-1}) = f(x^{s_1}) = \varphi$ . Therefore,  $s_1 = 1$ ,  
 610 otherwise it would not be the smallest integer satisfying (5.10). The equality  $x^{s_1} =$   
 611  $x^{s_1-1}$  then rewrites as  $x^1 = x^0 = x$  and, when  $k \rightarrow \infty$ ,

$$612 \quad x_{i_{j_k}+1} \rightarrow x^1 = x. \quad \square$$

613 **5.3. Third part: proof of Theorem 5.1.** Let  $\underline{x}$  be an accumulation point of  
 614  $(x_i)_{i \in \mathbb{N}}$ . Assume, for the sake of contradiction, that  $\underline{x}$  is not B-stationary for (1.1).  
 615 Let  $(x_{i_k})_{k \in \mathbb{N}}$  be a subsequence converging to  $\underline{x}$ .

616 Let  $\alpha_{\underline{x}}$  and  $\rho$  be as in Proposition 5.3. For all  $k \in \mathbb{N}$  large enough,  $x_{i_k} \in B(\underline{x}, \rho) \cap$   
 617  $C$ . Thus, when Algorithm 4.1 is called at point  $x_{i_k}$ , the condition in line 2 stops being  
 618 fulfilled for some  $\alpha_{i_k} \geq \alpha_{\underline{x}}$ , meaning that

$$619 \quad x_{i_k+1} \in P_C(x_{i_k} - \alpha_{i_k} \nabla f(x_{i_k})) \text{ for some } \alpha_{i_k} \in [\alpha_{\underline{x}}, \bar{\alpha}].$$

620 If we replace  $(i_k)_{k \in \mathbb{N}}$  with a subsequence, we can assume that  $(\alpha_{i_k})_{k \in \mathbb{N}}$  converges to  
 621 some  $\alpha_{\lim} \in [\alpha_{\underline{x}}, \bar{\alpha}]$ .

622 For all  $k \in \mathbb{N}$ , we have

$$623 \quad \|x_{i_k} - \alpha_{i_k} \nabla f(x_{i_k}) - x_{i_k+1}\| = d(x_{i_k} - \alpha_{i_k} \nabla f(x_{i_k}), C)$$

624 and since the distance to a nonempty closed set is a continuous function, we can  
 625 take this equality to the limit. We use the fact that  $x_{i_k+1} \rightarrow \underline{x}$  when  $k \rightarrow \infty$ , from  
 626 Proposition 5.4. This yields

$$627 \quad \|\alpha_{\text{lim}} \nabla f(\underline{x})\| = d(\underline{x} - \alpha_{\text{lim}} \nabla f(\underline{x}), C),$$

628 which means that  $\underline{x} \in P_C(\underline{x} - \alpha_{\text{lim}} \nabla f(\underline{x}))$ . In particular,  $-\nabla f(\underline{x}) \in \widehat{N}_C(\underline{x}) \subseteq \widehat{N}_C(\underline{x})$ ,  
 629 which contradicts our assumption that  $\underline{x}$  is not B-stationary for (1.1). We have  
 630 therefore proven that every accumulation point is B-stationary.

631 Finally, if  $(x_i)_{i \in \mathbb{N}}$  has an isolated accumulation point, then the sequence  $(x_i)_{i \in \mathbb{N}}$   
 632 converges, from Proposition 5.4 and [43, Lemma 4.10].

633 **6. Convergence analysis for a locally Lipschitz continuous gradient.** In  
 634 this section, PGD (Algorithm 4.2) is analyzed under hypothesis (H2). As mentioned  
 635 after Remark 4.5, only the nontrivial case where an infinite sequence is generated  
 636 is considered here. Specifically, the second part of the second item of Theorem 1.2,  
 637 restated in Theorem 6.3 for convenience, is proven based on Proposition 6.1 and  
 638 Corollary 6.2 which state that, for every  $\underline{x} \in C$  and every input  $x$  sufficiently close to  
 639  $\underline{x}$ , the PGD map (Algorithm 4.1) terminates after at most a given number of iterations  
 640 which depends only on  $\underline{x}$ .

641 **PROPOSITION 6.1.** *Assume that  $f$  satisfies (H2). Let  $\underline{x} \in C$ ,  $\bar{\alpha} \in (0, \infty)$ ,  $c \in$   
 642  $(0, 1)$ , and  $\rho \in (0, \infty)$ . Let  $\bar{\rho} \in [\rho + 2\bar{\alpha} \max_{x \in B[\underline{x}, \rho] \cap C} \|\nabla f(x)\|, \infty)$  and define  $\alpha_* :=$   
 643  $(1 - c) / \text{Lip}_{B[\underline{x}, \bar{\rho}]}(\nabla f)$ . Then, for all  $x \in B[\underline{x}, \rho] \cap C$ ,  $\alpha \in [0, \min\{\alpha_*, \bar{\alpha}\}]$ , and  $y \in$   
 644  $P_C(x - \alpha \nabla f(x))$ ,*

$$645 \quad f(y) \leq f(x) + c \langle \nabla f(x), y - x \rangle.$$

646 *Proof.* For all  $x \in B[\underline{x}, \rho] \cap C$  and  $\alpha \in [0, \bar{\alpha}]$ ,  $P_C(x - \alpha \nabla f(x)) \subseteq B[\underline{x}, \bar{\rho}]$ ; indeed,  
 647 for all  $y \in P_C(x - \alpha \nabla f(x))$ ,

$$648 \quad \|y - \underline{x}\| \leq \|y - x\| + \|x - \underline{x}\| \leq 2\alpha \|\nabla f(x)\| + \rho \leq \bar{\rho},$$

649 where the second inequality follows from (2.1). Thus, by (2.3) and (2.2), for all  
 650  $x \in B[\underline{x}, \rho] \cap C$ ,  $\alpha \in [0, \min\{\alpha_*, \bar{\alpha}\}]$ , and  $y \in P_C(x - \alpha \nabla f(x))$ ,

$$651 \quad f(y) \leq f(x) + \langle \nabla f(x), y - x \rangle + \frac{1}{2} \text{Lip}_{B[\underline{x}, \bar{\rho}]}(\nabla f) \|y - x\|^2$$

$$652 \quad \leq f(x) + \left(1 - \alpha \text{Lip}_{B[\underline{x}, \bar{\rho}]}(\nabla f)\right) \langle \nabla f(x), y - x \rangle$$

$$653 \quad \leq f(x) + c \langle \nabla f(x), y - x \rangle. \quad \square$$

654 **COROLLARY 6.2.** *Consider Algorithm 4.1 under hypothesis (H2). Given  $\underline{x} \in C$   
 655 and  $\rho \in (0, \infty)$ , let  $\bar{\rho}$  be as in Proposition 6.1. Then, for every  $x \in B[\underline{x}, \rho] \cap C$ , the  
 656 while loop terminates with a step size  $\alpha \in \left[\min\left\{\underline{\alpha}, \frac{\beta(1-c)}{\text{Lip}_{B[\underline{x}, \bar{\rho}]}(\nabla f)}\right\}, \bar{\alpha}\right]$  and hence after  
 657 at most*

$$658 \quad \max\left\{0, \left\lceil \ln\left(\frac{1-c}{\alpha_0 \text{Lip}_{B[\underline{x}, \bar{\rho}]}(\nabla f)}\right) / \ln(\beta) \right\rceil\right\}$$

659 iterations, where  $\alpha_0$  is the step size chosen in line 1.

660 *Proof.* At the latest, the while loop ends after iteration  $i \in \mathbb{N} \setminus \{0\}$  with  $\alpha = \alpha_0 \beta^i$   
 661 such that  $\frac{\alpha}{\beta} > \frac{1-c}{\text{Lip}_{B[\underline{x}, \bar{\rho}]}(\nabla f)}$ . In that case,  $i < 1 + \ln(\frac{1-c}{\alpha_0 \text{Lip}_{B[\underline{x}, \bar{\rho}]}(\nabla f)}) / \ln(\beta)$  and thus  
 662  $i \leq \lceil \ln(\frac{1-c}{\alpha_0 \text{Lip}_{B[\underline{x}, \bar{\rho}]}(\nabla f)}) / \ln(\beta) \rceil$ .  $\square$

663 **THEOREM 6.3.** *Assume that  $f$  satisfies (H2). Let  $(x_i)_{i \in \mathbb{N}}$  be a sequence generated*  
 664 *by PGD (Algorithm 4.2). Then, all accumulation points of  $(x_i)_{i \in \mathbb{N}}$  are  $P$ -stationary*  
 665 *for (1.1). Moreover, for every convergent subsequence  $(x_{i_j})_{j \in \mathbb{N}}$ ,*

$$666 \quad (6.1) \quad \lim_{j \rightarrow \infty} d(-\nabla f(x_{i_j+1}), \widehat{N}_C(x_{i_j+1})).$$

667 *Proof.* Assume that a subsequence  $(x_{i_j})_{j \in \mathbb{N}}$  converges to  $\underline{x} \in C$ . Given  $\rho \in (0, \infty)$ ,  
 668 let  $\bar{\rho}$  be as in Proposition 6.1. Define

$$669 \quad I := \left[ \min \left\{ \underline{\alpha}, \frac{\beta(1-c)}{\text{Lip}_{B[\underline{x}, \bar{\rho}]}(\nabla f)} \right\}, \bar{\alpha} \right].$$

670 There exists  $j_* \in \mathbb{N}$  such that, for all integers  $j \geq j_*$ ,  $x_{i_j} \in B[\underline{x}, \rho]$ , thus, by Corol-  
 671 lary 6.2,  $x_{i_j+1} \in P_C(x_{i_j} - \alpha_{i_j} \nabla f(x_{i_j}))$  with  $\alpha_{i_j} \in I$ , and hence

$$672 \quad \|x_{i_j+1} - (x_{i_j} - \alpha_{i_j} \nabla f(x_{i_j}))\| = d(x_{i_j} - \alpha_{i_j} \nabla f(x_{i_j}), C).$$

673 Since  $I$  is compact, a subsequence  $(\alpha_{i_{j_k}})_{k \in \mathbb{N}}$  converges to  $\alpha \in I$ . Moreover, there exists  
 674  $k_* \in \mathbb{N}$  such that  $j_{k_*} \geq j_*$ . Furthermore, by Proposition 5.4,  $(x_{i_j+1})_{j \in \mathbb{N}}$  converges to  
 675  $\underline{x}$ . Therefore, for all integers  $k \geq k_*$ ,

$$676 \quad \|x_{i_{j_k}+1} - (x_{i_{j_k}} - \alpha_{i_{j_k}} \nabla f(x_{i_{j_k}}))\| = d(x_{i_{j_k}} - \alpha_{i_{j_k}} \nabla f(x_{i_{j_k}}), C),$$

677 and letting  $k$  tend to infinity yields

$$678 \quad \|\underline{x} - (\underline{x} - \alpha \nabla f(\underline{x}))\| = d(\underline{x} - \alpha \nabla f(\underline{x}), C).$$

679 It follows that  $\underline{x} \in P_C(\underline{x} - \alpha \nabla f(\underline{x}))$ , which implies that  $-\nabla f(\underline{x}) \in \widehat{N}_C(\underline{x})$ .

680 We now establish (6.1). Recall that, for all integers  $j \geq j_*$ , since  $x_{i_j+1} \in P_C(x_{i_j} -$   
 681  $\alpha_{i_j} \nabla f(x_{i_j}))$  with  $\alpha_{i_j} \in I$ , it holds that  $\frac{1}{\alpha_{i_j}}(x_{i_j} - x_{i_j+1}) - \nabla f(x_{i_j}) \in \widehat{N}_C(x_{i_j+1})$ , and  
 682 thus

$$683 \quad d(-\nabla f(x_{i_j+1}), \widehat{N}_C(x_{i_j+1})) \leq \left\| -\nabla f(x_{i_j+1}) - \left( \frac{1}{\alpha_{i_j}}(x_{i_j} - x_{i_j+1}) - \nabla f(x_{i_j}) \right) \right\|$$

$$684 \quad \leq \frac{1}{\alpha_{i_j}} \|x_{i_j+1} - x_{i_j}\| + \|\nabla f(x_{i_j+1}) - \nabla f(x_{i_j})\|$$

$$685 \quad \rightarrow 0 \text{ when } j \rightarrow \infty,$$

686 by Proposition 5.4 and the fact that  $(\alpha_{i_j})_{j \in \mathbb{N}}$  is bounded away from zero.  $\square$

687 Proposition 6.4 considers the case where PGD generates a bounded sequence.

688 **PROPOSITION 6.4.** *Assume that  $f$  satisfies (H2). Let  $(x_i)_{i \in \mathbb{N}}$  be a sequence gen-*  
 689 *erated by PGD (Algorithm 4.2). If  $(x_i)_{i \in \mathbb{N}}$  is bounded, which is the case if the sublevel*  
 690 *set (1.6) is bounded, then all of its accumulation points, of which there exists at least*  
 691 *one, are  $P$ -stationary for (1.1) and have the same image by  $f$ , and*

$$692 \quad (6.2) \quad \lim_{i \rightarrow \infty} d(-\nabla f(x_i), \widehat{N}_C(x_i)) = 0.$$

693 *Proof.* Assume that  $(x_i)_{i \in \mathbb{N}}$  is bounded. It suffices to establish (6.2) and to prove  
 694 that all accumulation points of  $(x_i)_{i \in \mathbb{N}}$  have the same image by  $f$ ; the other statements  
 695 follow from Theorem 6.3.

696 The proof that all accumulation points of  $(x_i)_{i \in \mathbb{N}}$  have the same image by  $f$  is  
 697 based on the argument given in the proof of [51, Theorem 65]. Assume that  $(x_{i_k})_{k \in \mathbb{N}}$   
 698 and  $(x_{j_k})_{k \in \mathbb{N}}$  converge respectively to  $\underline{x}$  and  $\bar{x}$ . Being bounded, the sequence  $(x_i)_{i \in \mathbb{N}}$  is  
 699 contained in a compact set. By Propositions 4.6 and 4.7, the sequence  $(f(x_i))_{i \in \mathbb{N}}$  con-  
 700 verges; Proposition 4.6 applies because a continuous real-valued function is bounded  
 701 from below and uniformly continuous on every compact set [63, Propositions 1.3.3 and  
 702 1.3.5]. Therefore,  $f(\underline{x}) = \lim_{k \rightarrow \infty} f(x_{i_k}) = \lim_{i \rightarrow \infty} f(x_i) = \lim_{k \rightarrow \infty} f(x_{j_k}) = f(\bar{x})$ .

703 Let us establish (6.2). Assume, for the sake of contradiction, that (6.2) does not  
 704 hold. Then, there exist  $\varepsilon \in (0, \infty)$  and a subsequence  $(x_{i_j})_{j \in \mathbb{N}}$  such that  $i_0 \geq 1$  and  
 705  $d(-\nabla f(x_{i_j}), \widehat{N}_C(x_{i_j})) > \varepsilon$  for all  $j \in \mathbb{N}$ . Since  $(x_{i_{j-1}})_{j \in \mathbb{N}}$  is bounded, it contains a  
 706 subsequence  $(x_{i_{j_k-1}})_{k \in \mathbb{N}}$  that converges to a point  $\underline{x} \in C$ . Therefore, by (6.1),

$$707 \quad \lim_{k \rightarrow \infty} d(-\nabla f(x_{i_{j_k}}, \widehat{N}_C(x_{i_{j_k}})) = 0,$$

708 a contradiction.  $\square$

709 **7. Examples of feasible sets on which PGD can be practically imple-**  
 710 **mented.** Examples of a set  $C$  on which PGD can be practically implemented include:

- 711 1. the closed cone  $\mathbb{R}_{\leq s}^n$  of  $s$ -sparse vectors of  $\mathbb{R}^n$ , i.e., those having at most  $s$   
 712 nonzero components,  $n$  and  $s$  being positive integers such that  $s < n$ ;
- 713 2. the closed cone  $\mathbb{R}_{\leq s}^n \cap \mathbb{R}_+^n$  of nonnegative  $s$ -sparse vectors of  $\mathbb{R}^n$ ;
- 714 3. the determinantal variety [26, Lecture 9]

$$715 \quad \mathbb{R}_{\leq r}^{m \times n} := \{X \in \mathbb{R}^{m \times n} \mid \text{rank } X \leq r\},$$

716  $m, n$ , and  $r$  being positive integers such that  $r < \min\{m, n\}$ ;

- 717 4. the closed cone

$$718 \quad \mathbb{S}_{\leq r}^+(n) := \{X \in \mathbb{R}_{\leq r}^{n \times n} \mid X^\top = X, X \succeq 0\}$$

719 of order- $n$  real symmetric positive-semidefinite matrices of rank at most  $r$ ,  $n$   
 720 and  $r$  being positive integers such that  $r < n$ .

721 Indeed, for every set in this list, the projection map, the tangent cone, the regular  
 722 normal cone, and the normal cone are explicitly known; see [45, §§6 and 7.4] and the  
 723 references therein. In particular, it is known that these sets are not Clarke regular at  
 724 infinitely many points. In this section, we prove that, for these sets, regular normals  
 725 are proximal normals.

726 As detailed in [45], if  $C$  is a set in this list, then there exist a positive integer  $p$   
 727 and disjoint nonempty smooth submanifolds  $S_0, \dots, S_p$  of  $\mathcal{E}$  such that  $\overline{S_p} = C$  and,  
 728 for all  $i \in \{0, \dots, p\}$ ,  $\overline{S_i} = \bigcup_{j=0}^i S_j$ . This implies that  $\{S_0, \dots, S_p\}$  is a *stratification*  
 729 of  $C$  satisfying the *condition of the frontier* [39, §5]. Thus,  $C$  is called a *stratified set*  
 730 and  $S_0, \dots, S_p$  are called the *strata* of  $\{S_0, \dots, S_p\}$ .

731 **PROPOSITION 7.1.** *Let  $C$  be a set in the list. For all  $x \in C$ ,*

$$732 \quad \widehat{N}_C(x) = \widehat{N}_C(x)$$

733 *and, if  $x \notin S_p$ , then*

$$734 \quad \widehat{N}_C(x) \subsetneq N_C(x).$$

735 Since the proof of Proposition 7.1 relies on significantly different concepts than  
736 those previously used, we present it in Appendix A.

737 **8. Comparison of PGD and P<sup>2</sup>GD on a simple example.** P<sup>2</sup>GD, which  
738 is short for projected-projected gradient descent, was introduced in [55, Algorithm 3]  
739 for  $C := \mathbb{R}_{\leq r}^{m \times n}$  and extended to an arbitrary set  $C$  in [45, Algorithm 5.1]. It works  
740 like PGD except that it involves an additional projection: given  $x \in C$  as input, the  
741 P<sup>2</sup>GD map [45, Algorithm 5.1] performs a backtracking projected line search along a  
742 projection  $g$  of  $-\nabla f(x)$  onto  $T_C(x)$ , i.e., computes a projection  $y$  of  $x + \alpha g$  onto  $C$   
743 for decreasing values of the step size  $\alpha \in (0, \infty)$  until  $y$  satisfies an Armijo condition.

744 As pointed out in [57, §3.2], the convergence of optimization algorithms that  
745 use descent directions in the tangent cone, such as P<sup>2</sup>GD, often suffers from the  
746 discontinuity of the tangent cone. In [32, §2.2], on an instance of (1.1) where  $\mathcal{E} := \mathbb{R}^{3 \times 3}$   
747 and  $C := \mathbb{R}_{\leq 2}^{3 \times 3}$ , P<sup>2</sup>GD is proven to generate a sequence converging to a point of  
748 rank one that is M-stationary but not B-stationary. Several methods are compared  
749 numerically on this instance in [46, §8.2].

750 In this section, monotone PGD and P<sup>2</sup>GD are compared analytically on the in-  
751 stance of (1.1) where  $\mathcal{E} := \mathbb{R}^2$ ,  $C := \mathbb{R}_{\leq 1}^2$ ,  $f(x) := \frac{1}{2}\|x - x_*\|^2$  for all  $x \in \mathbb{R}^2$ ,  
752  $x_* := (a, 0)$ , and  $a \in \mathbb{R} \setminus \{0\}$ . For all  $x \in \mathbb{R}^2$ ,  $\nabla f(x) = x - x_*$ . Thus, the global  
753 Lipschitz constant of  $\nabla f$  is 1; in particular,  $f$  satisfies (H2). Both algorithms are used  
754 with  $\underline{\alpha} := \bar{\alpha} := \alpha \in (0, 2)$  and an arbitrary  $\beta \in (0, 1)$ . The initial iterate is  $(0, b)$  for  
755 some  $b \in \mathbb{R} \setminus \{0\}$ .

756 We recall from [45, Proposition 7.13] that  $T_{\mathbb{R}_{\leq 1}^2}(0, 0) = \mathbb{R}_{\leq 1}^2$  and, for all  $t \in \mathbb{R} \setminus \{0\}$ ,

$$757 \quad T_{\mathbb{R}_{\leq 1}^2}(0, t) = \{0\} \times \mathbb{R}, \quad T_{\mathbb{R}_{\leq 1}^2}(t, 0) = \mathbb{R} \times \{0\},$$

758 from [45, Propositions 7.16 and 7.17] that

$$759 \quad \widehat{N}_{\mathbb{R}_{\leq 1}^2}(0, 0) = \{(0, 0)\} \subsetneq \mathbb{R}_{\leq 1}^2 = N_{\mathbb{R}_{\leq 1}^2}(0, 0)$$

760 and, for all  $t \in \mathbb{R} \setminus \{0\}$ ,

$$761 \quad \widehat{N}_{\mathbb{R}_{\leq 1}^2}(0, t) = \mathbb{R} \times \{0\}, \quad \widehat{N}_{\mathbb{R}_{\leq 1}^2}(t, 0) = \{0\} \times \mathbb{R},$$

762 and from Proposition 7.1 that  $\widehat{N}_{\mathbb{R}_{\leq 1}^2}(x) = \widehat{N}_{\mathbb{R}_{\leq 1}^2}(x)$  for all  $x \in \mathbb{R}_{\leq 1}^2$ .

763 Proposition 8.1 explicitly describes the sequences generated by PGD and P<sup>2</sup>GD  
764 for small values of  $c$ . We omit its proof, which consists in elementary computations.

765 **PROPOSITION 8.1.** *If  $\alpha = 1$  and  $c \in (0, \frac{1}{2}]$ , then PGD and P<sup>2</sup>GD generate the  
766 finite sequences  $((0, b), (a, 0))$  and  $((0, b), (0, 0), (a, 0))$ , respectively. If  $\alpha \neq 1$ , then  
767 both algorithms generate infinite sequences.*

768 • For every  $c \in (0, \frac{2-\alpha}{2}]$ , P<sup>2</sup>GD generates the sequence  $((0, (1-\alpha)^i b))_{i \in \mathbb{N}}$  which  
769 converges to  $(0, 0)$ .

770 • For every  $c \in (0, \frac{2-\alpha}{4})$ :

771 – if  $\alpha|a|/|b| > |1-\alpha|$ , then PGD generates the sequence

$$772 \quad ((0, b), (a(1 - (1-\alpha)^{i+1}), 0))_{i \in \mathbb{N}};$$

773 – if  $\alpha|a|/|b| \leq |1-\alpha|$ , then  $i_* := \left\lceil \frac{\ln(\alpha|a|/|b|)}{\ln(1-\alpha)} \right\rceil \in \mathbb{N} \setminus \{0\}$  and PGD generates  
774 the sequence

$$775 \quad (((0, (1-\alpha)^i b))_{i=0}^{i=i_*}, (a(1 - (1-\alpha)^{i+1}), 0))_{i \in \mathbb{N}})$$

776

if  $\frac{\ln(\alpha|a|/|b|)}{\ln(1-\alpha)} \notin \mathbb{N}$  and

777

$$(((0, (1-\alpha)^i b))_{i=0}^{i=i_*}, (a(1-(1-\alpha)^{i_*+1}), 0))_{i \in \mathbb{N}}$$

778

$$\text{or } (((0, (1-\alpha)^i b))_{i=0}^{i=i_*+1}, (a(1-(1-\alpha)^{i_*+1}), 0))_{i \in \mathbb{N}}$$

779

if  $\frac{\ln(\alpha|a|/|b|)}{\ln(1-\alpha)} \in \mathbb{N}$ .

780

Thus, every sequence generated by PGD converges to  $(a, 0)$ .

781

In conclusion, if  $\alpha \neq 1$ , then P<sup>2</sup>GD converges to  $(0, 0)$ , which is M-stationary but not B-stationary, while PGD converges to  $(a, 0)$ , which is P-stationary and even a global minimizer of  $f|_{\mathbb{R}_{\leq 1}^2}$  (and  $f$ ). This is illustrated in Figure 2 for some choice of

784

$a, b$ , and  $\alpha$ .

785

By Proposition 8.1, for every sequence  $(x_i)_{i \in \mathbb{N}}$  generated by P<sup>2</sup>GD, it holds that

786

$$\lim_{i \rightarrow \infty} d(-\nabla f(x_i), \widehat{N}_{\mathbb{R}_{\leq 1}^2}(x_i)) = 0.$$

787

Thus, the measure of B-stationarity  $\mathbb{R}_{\leq 1}^2 \rightarrow \mathbb{R} : x \mapsto d(-\nabla f(x), \widehat{N}_{\mathbb{R}_{\leq 1}^2}(x))$  is not lower semicontinuous at  $(0, 0)$ , and the convergence to an M-stationary point that is not B-stationary cannot be suspected based on the mere observation of this limit. In the terminology of [32],  $((0, 0), (x_i)_{i \in \mathbb{N}}, f)$  is an apocalypse.

790

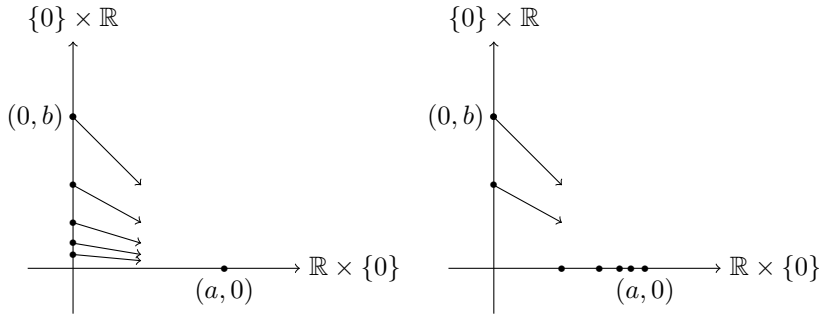


FIG. 2. First few iterates generated by PGD (right) and P<sup>2</sup>GD (left) on the instance of (1.1) studied in Section 8 with  $a := b := 1$  and  $\alpha := 0.45$ . The arrows represent  $x_i - \alpha \nabla f(x_i)$ . The point  $(a, 0)$ , which is the unique global minimizer, is also represented. It is already visible from the first few iterates that P<sup>2</sup>GD converges to the M-stationary point  $(0, 0)$  while PGD converges to the global minimizer.

791

**9. Conclusion.** The main contribution of this paper is the proof of Theorem 1.2.

792

This theorem ensures that PGD (Algorithm 4.2) enjoys the strongest stationarity properties that can be expected for problem (1.1) under the considered assumptions.

793

A sufficient condition for the convergence of a sequence generated by PGD is provided in Theorem 5.1. However, if satisfied, this condition does not offer a characterization of the rate of convergence. This important matter is addressed in [29] for monotone PGD under the assumption that  $f$  satisfies (H2) and a Kurdyka–Lojasiewicz property.

798

Two possible extensions of this work are left for future research. First, can Theorem 1.2 be extended to an algorithm that uses more general search directions than PGD? For example, a search direction at a point  $x \in C$  that is not B-stationary for (1.1) could be a vector  $v \notin \widehat{N}_C(x)$  that satisfies [22, conditions (2) and (3)], i.e.,  $\langle \nabla f(x), v \rangle \leq -c_1 \|\nabla f(x)\|^2$  and  $\|v\| \leq c_2 \|\nabla f(x)\|$  with  $c_1, c_2 \in (0, \infty)$ .

803

804 Second, can Theorem 1.2 be extended to the proximal gradient algorithm as defined in [31, Algorithm 4.1] or [11, Algorithm 3.1]? The first step toward such an extension would be defining suitable stationarity notions for the corresponding problem whose objective function is not differentiable. Furthermore, significant adaptations would be needed, e.g., because inequality (2.1), which plays an instrumental role in our analysis, does not seem to admit a straightforward extension.

810 **Appendix A. Proof of Proposition 7.1.** The strict inclusion follows from [45, Proposition 7.16] and [4, Theorem 3.9] if  $C = \mathbb{R}_{\leq s}^n$ , from [45, Proposition 6.7] and [60, Theorem 3.4] if  $C = \mathbb{R}_{\leq s}^n \cap \mathbb{R}_+^n$ , from [27, Corollary 2.3 and Theorem 3.1] if  $C = \mathbb{R}_{\leq r}^{m \times n}$ , and from [45, Proposition 6.28] and [60, Theorem 3.12] if  $C = \mathbb{S}_{\leq r}^+(n)$ .

814 By (1.2), it remains to prove that, for all  $x \in C$ ,  $\widehat{N}_C(x) \supseteq \widehat{N}_C(x)$ . This follows from [1, Lemma 4] if  $x \in S_p$ . Let  $x \in C \setminus S_p$ . If  $C$  is  $\mathbb{R}_{\leq s}^n$  or  $\mathbb{R}_{\leq r}^{m \times n}$ , then, by [45, Proposition 7.16] and [27, Corollary 2.3],  $\widehat{N}_C(x) = \{0\}$  and the result follows. If  $C$  is  $\mathbb{R}_{\leq s}^n \cap \mathbb{R}_+^n$  or  $\mathbb{S}_{\leq r}^+(n)$ , then the result follows from [45, Proposition 6.7] and [60, Proposition 3.2] or [45, Proposition 6.28] and [10, Corollary 17]; the detail is given below for completeness.

820 Assume that  $C$  is  $\mathbb{R}_{\leq s}^n \cap \mathbb{R}_+^n$ . Let  $\text{supp}(x) := \{i \in \{1, \dots, n\} \mid x_i \neq 0\}$ . By [45, Proposition 6.7],

$$822 \quad \widehat{N}_{\mathbb{R}_{\leq s}^n \cap \mathbb{R}_+^n}(x) = \{v \in \mathbb{R}_+^n \mid \text{supp}(v) \subseteq \{1, \dots, n\} \setminus \text{supp}(x)\}.$$

823 Thus, by [60, Proposition 3.2], for every  $v \in \widehat{N}_{\mathbb{R}_{\leq s}^n \cap \mathbb{R}_+^n}(x)$ ,  $P_{\mathbb{R}_{\leq s}^n \cap \mathbb{R}_+^n}(x + v) = \{x\}$ .

824 Assume now that  $C$  is  $\mathbb{S}_{\leq r}^+(n)$ . By [45, Proposition 6.28],

$$825 \quad \widehat{N}_{\mathbb{S}_{\leq r}^+(n)}(X) = \mathbb{S}(n)^\perp + \{Z \in \mathbb{S}^-(n) \mid XZ = 0_{n \times n}\},$$

826 with  $\mathbb{S}(n) := \{X \in \mathbb{R}^{n \times n} \mid X^\top = X\}$ ,  $\mathbb{S}(n)^\perp = \{X \in \mathbb{R}^{n \times n} \mid X^\top = -X\}$ , and  
827  $\mathbb{S}^-(n) := \{X \in \mathbb{S}(n) \mid X \preceq 0\}$ . Let  $Z \in \widehat{N}_{\mathbb{S}_{\leq r}^+(n)}(X)$  and  $Z_{\text{sym}} := \frac{1}{2}(Z + Z^\top)$ . Then,  
828 by [10, Corollary 17],  $P_{\mathbb{S}_{\leq r}^+(n)}(X + Z) = P_{\mathbb{S}_{\leq r}^+(n)}(X + Z_{\text{sym}})$ . Let  $\underline{r} := \text{rank } X$  and  
829  $\tilde{r} := \text{rank } Z_{\text{sym}}$ . Since  $\text{im } Z_{\text{sym}} \subseteq \ker X$ ,  $\tilde{r} \leq n - \underline{r}$  and there exists  $U \in \text{O}(n)$  such  
830 that

$$831 \quad X = U \text{diag}(\lambda_1(X), \dots, \lambda_{\underline{r}}(X), 0_{n-\underline{r}})U^\top$$

832 and

$$833 \quad Z_{\text{sym}} = U \text{diag}(0_{n-\tilde{r}}, \lambda_{n-\tilde{r}+1}(Z_{\text{sym}}), \dots, \lambda_n(Z_{\text{sym}}))U^\top$$

834 are eigendecompositions. Thus,

$$835 \quad X + Z_{\text{sym}} = U \text{diag}(\lambda_1(X), \dots, \lambda_{\underline{r}}(X), 0_{n-\underline{r}-\tilde{r}}, \lambda_{n-\tilde{r}+1}(Z_{\text{sym}}), \dots, \lambda_n(Z_{\text{sym}}))U^\top$$

836 is an eigendecomposition. Hence, by [10, Corollary 17],  $P_{\mathbb{S}_{\leq r}^+(n)}(X + Z_{\text{sym}}) = \{X\}$ .

837 **Acknowledgments.** The authors thank two anonymous referees for several  
838 helpful comments that improved the quality of the paper.



- 840 [1] P.-A. ABSIL AND J. MALICK, *Projection-like retractions on matrix manifolds*, SIAM Journal  
841 on Optimization, 22 (2012), pp. 135–158, <https://doi.org/10.1137/100802529>.
- 842 [2] A. A. AHMADI AND J. ZHANG, *On the complexity of finding a local minimizer of a quadratic*  
843 *function over a polytope*, Mathematical Programming, 195 (2022), pp. 783–792, <https://doi.org/10.1007/s10107-021-01714-2>.
- 844 [3] M. V. BALASHOV, B. T. POLYAK, AND A. A. TREMBA, *Gradient projection and conditional gra-*  
845 *dient methods for constrained nonconvex minimization*, Numerical Functional Analysis and  
846 Optimization, 41 (2020), pp. 822–849, <https://doi.org/10.1080/01630563.2019.1704780>.
- 847 [4] H. H. BAUSCHKE, D. R. LUKE, H. M. PHAN, AND X. WANG, *Restricted normal cones and*  
848 *sparsity optimization with affine constraints*, Foundations of Computational Mathematics,  
849 14 (2014), pp. 63–83, <https://doi.org/10.1007/s10208-013-9161-0>.
- 850 [5] A. BECK AND Y. C. ELДАР, *Sparsity constrained nonlinear optimization: Optimality conditions*  
851 *and algorithms*, SIAM Journal on Optimization, 23 (2013), pp. 1480–1509, <https://doi.org/10.1137/120869778>.
- 852 [6] A. BECK AND M. TEBOLLE, *Fixed-Point Algorithms for Inverse Problems in Science and En-*  
853 *gineering*, vol. 49 of Springer Optimization and Its Applications, Springer New York, 2011,  
854 ch. A Linearly Convergent Algorithm for Solving a Class of Nonconvex/Affine Feasibility  
855 Problems, pp. 33–48, [https://doi.org/10.1007/978-1-4419-9569-8\\_3](https://doi.org/10.1007/978-1-4419-9569-8_3).
- 856 [7] M. BENKO AND H. GFRERER, *On estimating the regular normal cone to constraint systems*  
857 *and stationarity conditions*, Optimization, 66 (2017), pp. 61–92, <https://doi.org/10.1080/02331934.2016.1252915>.
- 858 [8] M. BENKO AND H. GFRERER, *New verifiable stationarity concepts for a class of mathematical*  
859 *programs with disjunctive constraints*, Optimization, 67 (2018), pp. 1–23, <https://doi.org/10.1080/02331934.2017.1387547>.
- 860 [9] Y. CUI AND J.-S. PANG, *Modern Nonconvex Nondifferentiable Optimization*, MOS-SIAM Series  
861 on Optimization, Society for Industrial and Applied Mathematics, Philadelphia, PA, 2021,  
862 <https://doi.org/10.1137/1.9781611976748>.
- 863 [10] A. DAX, *Low-rank positive approximants of symmetric matrices*, Advances in Linear Algebra  
864 & Matrix Theory, 4 (2014), pp. 172–185, <https://doi.org/10.4236/alamt.2014.43015>.
- 865 [11] A. DE MARCHI, *Proximal gradient methods beyond monotony*, Journal of Nonsmooth Analysis  
866 and Optimization, 4 (2023), <https://doi.org/10.46298/jnsao-2023-10290>.
- 867 [12] S. DOLECKI AND G. H. GRECO, *Towards historical roots of necessary conditions of optimality:*  
868 *Regula of Peano*, Control and Cybernetics, 36 (2007), pp. 491–518.
- 869 [13] S. DOLECKI AND G. H. GRECO, *Tangency vis-à-vis differentiability by Peano, Severi and*  
870 *Guareschi*, Journal of Convex Analysis, 18 (2011), p. 301–339.
- 871 [14] M. L. FLEGEL AND C. KANZOW, *Abadie-type constraint qualification for mathematical programs*  
872 *with equilibrium constraints*, Journal of Optimization Theory and Applications, 124 (2005),  
873 pp. 595–614, <https://doi.org/10.1007/s10957-004-1176-x>.
- 874 [15] M. L. FLEGEL AND C. KANZOW, *On the Guignard constraint qualification for mathematical*  
875 *programs with equilibrium constraints*, Optimization, 54 (2005), pp. 517–534, <https://doi.org/10.1080/02331930500342591>.
- 876 [16] M. L. FLEGEL, C. KANZOW, AND J. V. OUTRATA, *Optimality conditions for disjunctive*  
877 *programs with application to mathematical programs with equilibrium constraints*,  
878 Set-Valued and Variational Analysis, 15 (2007), pp. 139–162, <https://doi.org/10.1007/s11228-006-0033-5>.
- 879 [17] M. FUKUSHIMA AND G.-H. LIN, *Smoothing methods for mathematical programs with equilib-*  
880 *rium constraints*, in International Conference on Informatics Research for Development of  
881 Knowledge Society Infrastructure, 2004 (ICKS 2004), Kyoto, Japan, 2004, IEEE, pp. 206–  
882 213, <https://doi.org/10.1109/ICKS.2004.1313426>.
- 883 [18] M. FUKUSHIMA AND J.-S. PANG, *Convergence of a smoothing continuation method for mathe-*  
884 *matical programs with complementarity constraints*, in Ill-posed Variational Problems and  
885 Regularization Techniques, M. Théra and R. Tichatschke, eds., vol. 477 of Lecture Notes  
886 in Economics and Mathematical Systems, Springer, Berlin, Heidelberg, 1999, pp. 99–110,  
887 [https://doi.org/10.1007/978-3-642-45780-7\\_7](https://doi.org/10.1007/978-3-642-45780-7_7).
- 888 [19] M. FUKUSHIMA AND P. TSENG, *An implementable active-set algorithm for computing a B-*  
889 *stationary point of a mathematical program with linear complementarity constraints*,  
890 SIAM Journal on Optimization, 12 (2002), pp. 724–739, <https://doi.org/10.1137/S1052623499363232>.
- 891 [20] B. GAO, R. PENG, AND Y.-X. YUAN, *Low-rank optimization on Tucker tensor varieties*, (2023),  
892 <https://arxiv.org/abs/2311.18324>.
- 893 [21] H. GFRERER, *Optimality conditions for disjunctive programs based on generalized differentia-*  
894 *tion with application to mathematical programs with equilibrium constraints*, SIAM Journal  
895

- 902 on Optimization, 24 (2014), pp. 898–931, <https://doi.org/10.1137/130914449>.
- 903 [22] L. GRIPPO, F. LAMPARIELLO, AND S. LUCIDI, *A nonmonotone line search technique for Newton's*  
904 *method*, SIAM Journal on Numerical Analysis, 23 (1986), pp. 707–716, <https://doi.org/10.1137/0723046>.
- 905 [23] M. GUIGNARD, *Generalized Kuhn–Tucker conditions for mathematical programming problems*  
906 *in a Banach space*, SIAM Journal on Control, 7 (1969), pp. 232–241, <https://doi.org/10.1137/0307016>.
- 907 [24] L. GUO AND G.-H. LIN, *Notes on some constraint qualifications for mathematical programs*  
908 *with equilibrium constraints*, Journal of Optimization Theory and Applications, 156 (2013),  
909 pp. 600–616, <https://doi.org/10.1007/s10957-012-0084-8>.
- 910 [25] W. HA, H. LIU, AND R. F. BARBER, *An equivalence between critical points for rank constraints*  
911 *versus low-rank factorizations*, SIAM Journal on Optimization, 30 (2020), pp. 2927–2955,  
912 <https://doi.org/10.1137/18M1231675>.
- 913 [26] J. HARRIS, *Algebraic Geometry*, vol. 133 of Graduate Texts in Mathematics, Springer-Verlag  
914 New York, 1992, <https://doi.org/10.1007/978-1-4757-2189-8>.
- 915 [27] S. HOSSEINI, D. R. LUKE, AND A. USCHMAJEW, *Nonsmooth Optimization and Its Applica-*  
916 *tions*, vol. 170 of International Series of Numerical Mathematics, Birkhäuser Cham, 2019,  
917 ch. Tangent and Normal Cones for Low-Rank Matrices, pp. 45–53, [https://doi.org/10.1007/978-3-030-11370-4\\_3](https://doi.org/10.1007/978-3-030-11370-4_3).
- 918 [28] X. M. HU AND D. RALPH, *Convergence of a penalty method for mathematical programming*  
919 *with complementarity constraints*, Journal of Optimization Theory and Applications, 123  
920 (2004), pp. 365–390, <https://doi.org/10.1007/s10957-004-5154-0>.
- 921 [29] X. JIA, C. KANZOW, AND P. MEHLITZ, *Convergence analysis of the proximal gradient method*  
922 *in the presence of the Kurdyka–Lojasiewicz property without global Lipschitz assump-*  
923 *tions*, SIAM Journal on Optimization, 33 (2023), pp. 3038–3056, <https://doi.org/10.1137/23M1548293>.
- 924 [30] X. JIA, C. KANZOW, P. MEHLITZ, AND G. WACHSMUTH, *An augmented Lagrangian method for*  
925 *optimization problems with structured geometric constraints*, Mathematical Programming,  
926 199 (2023), pp. 1365–1415, <https://doi.org/10.1007/s10107-022-01870-z>.
- 927 [31] C. KANZOW AND P. MEHLITZ, *Convergence properties of monotone and nonmonotone proximal*  
928 *gradient methods revisited*, Journal of Optimization Theory and Applications, 195 (2022),  
929 pp. 624–646, <https://doi.org/10.1007/s10957-022-02101-3>.
- 930 [32] E. LEVIN, J. KILEEL, AND N. BOUMAL, *Finding stationary points on bounded-rank matrices: a*  
931 *geometric hurdle and a smooth remedy*, Mathematical Programming, 199 (2023), pp. 831–  
932 864, <https://doi.org/10.1007/s10107-022-01851-2>.
- 933 [33] E. LEVIN, J. KILEEL, AND N. BOUMAL, *The effect of smooth parametrizations on nonconvex*  
934 *optimization landscapes*, Mathematical Programming, (2024), <https://doi.org/10.1007/s10107-024-02058-3>.
- 935 [34] X. LI AND Z. LUO, *Normal cones intersection rule and optimality analysis for low-rank matrix*  
936 *optimization with affine manifolds*, SIAM Journal on Optimization, 33 (2023), pp. 1333–  
937 1360, <https://doi.org/10.1137/22M147863X>.
- 938 [35] X. LI, W. SONG, AND N. XIU, *Optimality conditions for rank-constrained matrix optimization*,  
939 Journal of the Operations Research Society of China, 7 (2019), pp. 285–301, <https://doi.org/10.1007/s40305-019-00245-0>.
- 940 [36] X. LI, N. XIU, AND S. ZHOU, *Matrix optimization over low-rank spectral sets: Stationary points*  
941 *and local and global minimizers*, Journal of Optimization Theory and Applications, 184  
942 (2020), pp. 895–930, <https://doi.org/10.1007/s10957-019-01606-8>.
- 943 [37] Z. LUO AND L. QI, *Optimality conditions for Tucker low-rank tensor optimization*, Computa-  
944 tional Optimization and Applications, 86 (2023), pp. 1275–1298, <https://doi.org/10.1007/s10589-023-00465-4>.
- 945 [38] Z.-Q. LUO, J.-S. PANG, D. RALPH, AND S.-Q. WU, *Exact penalization and stationarity condi-*  
946 *tions of mathematical programs with equilibrium constraints*, Mathematical Programming,  
947 75 (1996), pp. 19–76, <https://doi.org/10.1007/BF02592205>.
- 948 [39] J. MATHER, *Notes on topological stability*, Bulletin of the American Mathematical Society, 49  
949 (2012), pp. 475–506, <https://doi.org/10.1090/S0273-0979-2012-01383-6>.
- 950 [40] P. MEHLITZ, *Stationarity conditions and constraint qualifications for mathematical programs*  
951 *with switching constraints*, Mathematical Programming, 181 (2020), pp. 149–186, <https://doi.org/10.1007/s10107-019-01380-5>.
- 952 [41] B. S. MORDUKHOVICH, *Variational Analysis and Generalized Differentiation I*, vol. 330 of  
953 Grundlehren der mathematischen Wissenschaften, Springer-Verlag Berlin Heidelberg, 2006,  
954 <https://doi.org/10.1007/3-540-31247-1>.
- 955 [42] B. S. MORDUKHOVICH, *Variational Analysis and Applications*, Springer Monographs in Math-

- 964 ematics, Springer Cham, 2018, <https://doi.org/10.1007/978-3-319-92775-6>.
- 965 [43] J. J. MORÉ AND D. C. SORENSSEN, *Computing a trust region step*, SIAM Journal on Scientific  
966 and Statistical Computing, 4 (1983), pp. 553–572, <https://doi.org/10.1137/0904038>.
- 967 [44] Y. NESTEROV, *Lectures on Convex Optimization*, vol. 137 of Springer Optimization and Its  
968 Applications, Springer, Cham, 2nd ed., 2018, <https://doi.org/10.1007/978-3-319-91578-4>.
- 969 [45] G. OLIKIER, K. A. GALLIVAN, AND P.-A. ABSIL, *First-order optimization on stratified sets*,  
970 (2023), <https://arxiv.org/abs/2303.16040>.
- 971 [46] G. OLIKIER, K. A. GALLIVAN, AND P.-A. ABSIL, *Low-rank optimization methods based on*  
972 *projected-projected gradient descent that accumulate at bouligand stationary points*, (2024),  
973 <https://arxiv.org/abs/2201.03962v2>.
- 974 [47] J.-S. PANG, *Partially B-regular optimization and equilibrium problems*, Mathematics of Oper-  
975 ations Research, 32 (2007), pp. 687–699, <https://doi.org/10.1287/moor.1070.0262>.
- 976 [48] J.-S. PANG AND M. FUKUSHIMA, *Complementarity constraint qualifications and simplified B-*  
977 *stationarity conditions for mathematical programs with equilibrium constraints*, Computa-  
978 tional Optimization and Applications, 13 (1999), pp. 111–136, [https://doi.org/10.1023/A:](https://doi.org/10.1023/A:1008656806889)  
979 1008656806889.
- 980 [49] J.-S. PANG, M. RAZAVIYAYN, AND A. ALVARADO, *Computing B-stationary points of nonsmooth*  
981 *DC programs*, Mathematics of Operations Research, 42 (2017), pp. 95–118, [https://doi.](https://doi.org/10.1287/moor.2016.0795)  
982 [org/10.1287/moor.2016.0795](https://doi.org/10.1287/moor.2016.0795).
- 983 [50] E. PAUWELS, *Generic Fréchet stationarity in constrained optimization*, (2024), [https://arxiv.](https://arxiv.org/abs/2402.09831v1)  
984 [org/abs/2402.09831v1](https://arxiv.org/abs/2402.09831v1).
- 985 [51] E. POLAK, *Computational Methods in Optimization*, vol. 77 of Mathematics in Science and  
986 Engineering, Academic Press, 1971.
- 987 [52] S. M. ROBINSON, *Nonlinear Analysis and Optimization*, vol. 30 of Mathematical Programming  
988 Studies, Springer Berlin Heidelberg, 1987, ch. Local structure of feasible sets in nonlinear  
989 programming, Part III: Stability and sensitivity, pp. 45–66, [https://doi.org/10.1007/](https://doi.org/10.1007/BFb0121154)  
990 BFb0121154.
- 991 [53] R. T. ROCKAFELLAR AND R. J.-B. WETS, *Variational Analysis*, vol. 317 of Grundlehren der  
992 mathematischen Wissenschaften, Springer-Verlag Berlin Heidelberg, 1998, [https://doi.org/](https://doi.org/10.1007/978-3-642-02431-3)  
993 10.1007/978-3-642-02431-3. Corrected 3rd printing 2009.
- 994 [54] H. SCHEEL AND S. SCHOLTES, *Mathematical programs with complementarity constraints: Sta-*  
995 *tionarity, optimality, and sensitivity*, Mathematics of Operations Research, 25 (2000),  
996 pp. 1–22, <https://doi.org/10.1287/moor.25.1.1.15213>.
- 997 [55] R. SCHNEIDER AND A. USCHMAJEV, *Convergence results for projected line-search methods on*  
998 *varieties of low-rank matrices via Lojasiewicz inequality*, SIAM Journal on Optimization,  
999 25 (2015), pp. 622–646, <https://doi.org/10.1137/140957822>.
- 1000 [56] S. SCHOLTES, *Convergence properties of a regularization scheme for mathematical programs*  
1001 *with complementarity constraints*, SIAM Journal on Optimization, 11 (2001), pp. 918–936,  
1002 <https://doi.org/10.1137/S1052623499361233>.
- 1003 [57] S. SCHOLTES, *Nonconvex structures in nonlinear programming*, Operations Research, 52 (2004),  
1004 pp. 368–383, <https://doi.org/10.1287/opre.1030.0102>.
- 1005 [58] S. SCHOLTES AND M. STÖHR, *Exact penalization of mathematical programs with equilibrium*  
1006 *constraints*, SIAM Journal on Optimization, 37 (1999), pp. 617–652, [https://doi.org/10.](https://doi.org/10.1137/S0363012996306121)  
1007 1137/S0363012996306121.
- 1008 [59] S. STEFFENSEN AND M. ULBRICH, *A new relaxation scheme for mathematical programs with*  
1009 *equilibrium constraints*, SIAM Journal on Optimization, 20 (2010), pp. 2504–2539, [https:](https://doi.org/10.1137/090748883)  
1010 [//doi.org/10.1137/090748883](https://doi.org/10.1137/090748883).
- 1011 [60] M. K. TAM, *Regularity properties of non-negative sparsity sets*, Journal of Mathematical Analy-  
1012 sis and Applications, 447 (2017), pp. 758–777, <https://doi.org/10.1016/j.jmaa.2016.10.040>.
- 1013 [61] A. THEMELIS, L. STELLA, AND P. PATRINOS, *Forward-backward envelope for the sum of two*  
1014 *nonconvex functions: Further properties and nonmonotone linesearch algorithms*, SIAM  
1015 Journal on Optimization, 28 (2018), pp. 2274–2303, <https://doi.org/10.1137/16M1080240>.
- 1016 [62] P. P. VARAIYA, *Nonlinear programming in Banach space*, SIAM Journal on Applied Mathe-  
1017 matics, 15 (1967), pp. 284–293, <https://doi.org/10.1137/0115028>.
- 1018 [63] M. WILLEM, *Functional Analysis: Fundamentals and Applications*, Cornerstones, Birkhäuser  
1019 New York, 2013, <https://doi.org/10.1007/978-1-4614-7004-5>.
- 1020 [64] J. WU, L. ZHANG, AND Y. ZHANG, *Mathematical programs with semidefinite cone complemen-*  
1021 *tarity constraints: Constraint qualifications and optimality conditions*, Set-Valued and  
1022 Variational Analysis, 22 (2014), pp. 155–187, <https://doi.org/10.1007/s11228-013-0242-7>.
- 1023 [65] J. J. YE, *Necessary and sufficient optimality conditions for mathematical programs with equilib-*  
1024 *rium constraints*, Journal of Mathematical Analysis and Applications, 307 (2005), pp. 350–  
1025 369, <https://doi.org/10.1016/j.jmaa.2004.10.032>.