

# A PROXIMAL ORACLE SEQUENTIAL NLP ALGORITHM FOR CONSTRAINED NONSMOOTH OPTIMIZATION\*

COSMIN G. PETRA<sup>†</sup>

**Abstract.** We develop a proximal sequential nonlinear programming (NLP) method for constrained finite-sum optimization with nonsmooth, nonconvex objective terms. Each nonsmooth summand is assumed only to be proper, lower semicontinuous, and prox-bounded, and is accessed through a proximal oracle requiring no derivatives or sensitivity information. The method applies Moreau envelope regularization separately to each summand, preserving separability while retaining the original smooth nonlinear constraints. We establish consistency of this distributed Moreau continuation: as the regularization parameter vanishes, optimal values and global minimizers of the regularized problems converge, along subsequences, to those of the original problem. We then propose a continuation algorithm whose inner loop solves sequential NLP majorization models. These models enforce the smooth constraints exactly and represent the regularized nonsmooth terms by proximal oracle subgradient linearizations plus quadratic stabilization, using the upper- $C^2$  property of Moreau envelopes; an acceptance-ratio test adaptively updates the stabilization parameter. For fixed regularization, the cluster points of the inner iteration satisfy the stationarity conditions for the regularized problem, and the inner loop terminates finitely for any positive stationarity tolerance. For prescribed positive regularization and stationarity tolerances, the full continuation scheme also terminates finitely and returns an approximately limiting sum-stationary point with an explicit error bound. As the regularization and tolerances vanish, cluster points satisfy first-order limiting stationarity for the original nonsmooth problem under standard constraint qualifications and an attentive convergence condition, which simplify for locally Lipschitz terms. Finally, we show how the algorithm applies when the nonsmooth summands are optimal value functions of parametric optimization problems, for which the proximal oracles can be implemented by quadratically stabilized subproblems.

**1. Introduction.** We consider optimization problems of the form

$$(1.1) \quad \min_{x \in \mathbb{R}^n} F(x) := f(x) + \sum_{i \in \mathcal{K}} r_i(x),$$

$$(1.2) \quad \text{s.t. } h(x) = 0, \quad x_\ell \leq x \leq x_u,$$

where  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  and  $h : \mathbb{R}^n \rightarrow \mathbb{R}^m$  are  $C^2$ , and each  $r_i : \mathbb{R}^n \rightarrow \mathbb{R}$  is proper, lower semicontinuous (lsc), finite-valued, and possibly nonconvex and nonsmooth, for each  $i$  in a finite index set  $\mathcal{K}$ . Problems with this structure arise in several applications; in particular, our main motivation comes from stochastic programming with recourse and security-constrained power system optimization, where the functions  $r_i$  are optimal value functions of parametric optimization problems, as elaborated in Section 5.

We assume that the functions  $r_i$  are accessed through an oracle. In many applications, evaluating  $r_i(x)$  already requires solving an optimization subproblem, while derivatives of  $r_i$  with respect to  $x$  are unavailable, unreliable, or expensive to compute. Rather than requiring such sensitivities, we assume access to a proximal oracle, namely the value of the Moreau envelope [30] of  $r_i$

$$(1.3) \quad e_\lambda r_i(x) := \inf_{z \in \mathbb{R}^n} \left\{ r_i(z) + \frac{1}{2\lambda} \|z - x\|^2 \right\},$$

together with a proximal point

$$w_i \in \text{prox}_\lambda r_i(x) := \arg \min_{z \in \mathbb{R}^n} \left\{ r_i(z) + \frac{1}{2\lambda} \|z - x\|^2 \right\}, \quad \forall i \in \mathcal{K},$$

---

\*This work was supported by the U.S. Department of Energy through the Scientific Discovery through Advanced Computing (SciDAC) program under the ASCR-OE SLOPE-Grid partnership.

<sup>†</sup>Center for Applied Scientific Computing, Lawrence Livermore National Laboratory, Livermore, CA, USA. Email: petra1@llnl.gov. Release number: LLNL-JRNL-2019378

42 for  $\lambda > 0$ . When  $r_i$  is defined through an optimization problem, existing numerical  
 43 machinery can often be adapted to evaluate this oracle with little structural change.

44 Our starting point is the observation that the Moreau envelope regularizes each  
 45 nonsmooth term. Under the mild assumption of prox-boundedness,  $e_\lambda r_i$  converges  
 46 pointwise to  $r_i$  as  $\lambda \downarrow 0$  and is locally Lipschitz and upper- $C^2$  for sufficiently small  $\lambda$ .  
 47 This suggests replacing (1.1)–(1.2) by the family of regularized problems

$$48 \quad (1.4) \quad \min_{x \in \mathbb{R}^n} F_\lambda(x) := f(x) + \sum_{i \in \mathcal{K}} e_\lambda r_i(x)$$

$$49 \quad (1.5) \quad \text{s.t. } h(x) = 0 \text{ and } x_\ell \leq x \leq x_u,$$

50 and driving  $\lambda \downarrow 0$  in a continuation scheme.

51 For fixed  $\lambda$ , the objective in (1.4) is a sum of a smooth term and upper- $C^2$  terms,  
 52 and is therefore upper- $C^2$ . At first sight, this places the fixed- $\lambda$  problem within  
 53 the scope of the sequential quadratic programming (SQP) framework of Wang and  
 54 Petra [33] for nonsmooth problems with upper- $C^2$  objectives. The present work,  
 55 however, is not a direct application of that framework. In [33], the nonsmooth ob-  
 56 jective (1.4) is treated as an upper- $C^2$  function for which Clarke subgradients are  
 57 assumed to be available; here, we do not assume access to subgradients of either  $r_i$  or  
 58  $e_\lambda r_i$ . The only information used by the algorithm is a proximal oracle for each  $r_i$ .

59 The use of Moreau envelopes leads to a different stationarity mechanism, which  
 60 ultimately provides limiting/Mordukhovich subgradients stationarity for the original  
 61 problem and is stronger than the Clarke stationarity of [33] even for the regularized  
 62 problems, despite apparently not using subgradients. The key idea is to use proximal  
 63 points  $w_i \in \text{prox}_\lambda r_i(x)$ , produced by the oracle, to form the vectors

$$64 \quad v_i = \frac{1}{\lambda}(x - w_i).$$

65 For general lsc  $r_i$ , the vectors  $v_i$  may not be limiting subgradients of  $e_\lambda r_i$  at  $x$ , but:

66 (i) satisfy a majorization bound specific to upper- $C^2$  functions; and

67 (ii) are limiting subgradients of  $r_i$  at the nearby proximal point  $w_i$ .

68 The property (ii) is noteworthy and plays an important role in the analysis of the  
 69 algorithm proposed here. The property stems from the use of Moreau envelopes and  
 70 is not present for general upper- $C^2$  functions. We refer to the vectors  $v_i$  as *proximal*  
 71 *oracle subgradients* to emphasize their Moreau envelope origin and to distinguish them  
 72 from proximal subgradients previously used in the literature [30].

73 Consequently, property (i) allows us to develop a sequential proximal oracle NLP  
 74 algorithm for the regularized, fixed- $\lambda$  constrained problems (1.4)–(1.5). The method  
 75 builds convex quadratic stabilized models for the nonsmooth objective term(s) using  
 76 only proximal points and uses an NLP model for step computations. The choice of  
 77 an NLP model instead of an SQP model is made to enforce the smooth constraints  
 78 and objective term exactly, as required in important classes of problems, such as  
 79 power grid optimization, which motivates this work. An acceptance ratio test is  
 80 also devised to adjust the stabilization parameter and obtain a simple descent-type  
 81 of scheme. We prove that the inner, fixed- $\lambda$  iterations result in cluster points that  
 82 satisfy stationarity conditions for the regularized problem and terminate finitely for  
 83 any positive stationarity tolerance.

84 Furthermore, the property (ii) above yields an outer-loop convergence analysis  
 85 for the continuation scheme that carries, as  $\lambda \downarrow 0$ , the stationarity certificates for the  
 86 fixed- $\lambda$  problems to limiting subgradients-based stationarity of the original nonsmooth

87 problem. More specifically, under standard constraint qualifications (CQs) and a  
 88 mild attentive convergence condition—which holds automatically in several important  
 89 cases, including locally Lipschitz  $r_i$ —we show that every cluster point generated as  
 90  $\lambda \downarrow 0$  satisfies first-order stationarity conditions for (1.1)–(1.2).

91 Applying Moreau regularization separately to individual summands  $r_i$  preserves  
 92 the finite-sum structure in the regularized problems, allows parallel proximal com-  
 93 putations, and results in single-level decomposition with significant computational  
 94 parallelism. For example, we have previously shown that this single-level decomposi-  
 95 tion is particularly well suited when each  $r_i$  is the optimal value of a computationally  
 96 expensive optimization subproblem and can scale up to exascale supercomputers [20].  
 97 However, since infimal convolution does not in general commute with summation,  
 98 replacing each  $r_i$  by  $e_\lambda r_i$  is not merely a rewriting of the Moreau envelope of the full  
 99 nonsmooth objective. We show that, under mild assumptions, the regularized prob-  
 100 lems (1.4)–(1.5) are nevertheless consistent with the original problem, in the sense  
 101 that global minimizers of the regularized problems converge, along subsequences, to  
 102 global minimizers of (1.1)–(1.2) as  $\lambda \downarrow 0$ .

103 **1.1. Related work.** Our approach is related to proximal point algorithms (PPAs) [4,  
 104 5, 26], but differs in an essential way. Classical PPAs regularize the entire objective  
 105 and would require evaluating the proximal mapping of  $F$ , or equivalently solving sub-  
 106 problems involving  $e_\lambda F$ . In our setting, this would couple all objective terms and  
 107 destroy the separable structure and the parallelization potential of the single-level  
 108 decomposition associated with the nonsmooth part.

109 Compared with Douglas-Rachford splitting methods [21, 31, 32, 17] and related al-  
 110 ternating direction method of multipliers schemes [34, 15, 1, 2, 16], the present frame-  
 111 work is designed to work directly with the native structure of (1.1)–(1.2). Splitting  
 112 methods are typically most natural after reformulating the problem into a two-block  
 113 or linearly constrained consensus form, which may require aggregating the nonsmooth  
 114 terms or introducing auxiliary variables and coupling constraints. In contrast, we keep  
 115 the smooth nonlinear constraints explicit and preserve the finite-sum structure.

116 The closest broad framework is the constrained composite augmented Lagrangian  
 117 approach of [12], which treats problems with a single nonsmooth composite term and  
 118 general set-membership constraints. The framework of [12] is more general, but it does  
 119 not take advantage of the particular oracle and separability structures that motivate  
 120 this work. Our setting is deliberately more structured: each nonsmooth summand is  
 121 accessed through its proximal oracle, and the original smooth nonlinear constraints  
 122 are retained explicitly rather than absorbed into a reformulated composite constraint,  
 123 which would obscure the term-wise subgradients needed in the continuation analysis.

124 The inner-loop sequential NLP approach here is related to nonsmooth SQP and  
 125 gradient-sampling methods; see, for example, [35, 10, 6, 8, 9]. These methods build  
 126 SQP models using sampled gradients or quasi-Newton approximations for locally Lip-  
 127 schitz nonsmooth constrained problems. However, our stationarity analysis relies on  
 128 proximal points and deterministic nearby limiting subgradients rather than on Clarke  
 129 subgradient sampling or penalty-SQP stationarity. Moreover, our setting assumes no  
 130 derivative or sampled-gradient access to the nonsmooth summands.

131 **1.2. Contributions.** The paper’s main contribution is a proximal-oracle se-  
 132 quential NLP algorithm for constrained finite-sum nonsmooth, nonconvex optimiza-  
 133 tion under very general assumptions: each nonsmooth summand need only be proper,  
 134 lsc, and prox-bounded, with no convexity, weak convexity, prox-regularity, or dif-  
 135 ferentiability required. The algorithm uses proximal oracle evaluations rather than

136 sensitivity information and preserves the single-level decomposition of the original  
 137 problem. It provides finite inner-loop termination and stationarity for the regularized  
 138 problems by using the upper- $C^2$  property and proximal oracle subgradients features  
 139 specific to Moreau envelopes. The stationarity certificates from the regularized prob-  
 140 lems are general enough to be carried via an outer continuation analysis and to obtain  
 141 limiting first-order sum-stationarity for the original nonsmooth problem under stan-  
 142 dard CQs and a mild attentive convergence condition. The algorithm applies naturally  
 143 to and is well suited for nonsmooth terms defined by optimal value functions, where  
 144 proximal oracles can be implemented via stabilized subproblems.

145 **1.3. Organization.** The remainder of the paper is organized as follows. Section  
 146 2 reviews the variational analysis tools used throughout, with emphasis on upper- $C^2$   
 147 functions and subdifferential properties of Moreau envelopes. Section 3 establishes  
 148 the well-posedness and consistency of the Moreau envelope continuation approach  
 149 as the regularization parameter tends to zero. Section 4 presents the algorithm and  
 150 analyzes both the fixed- $\lambda$  inner iterations and the outer continuation scheme. Section  
 151 5 discusses how the framework applies naturally to parametric recourse optimization  
 152 models, with a case study for electric power grids. The focus here is theoretical  
 153 rather than computational, as the present paper establishes the prerequisites needed  
 154 to justify and support future computational developments. Section 6 summarizes  
 155 some of these developments and provides concluding remarks.

156 **2. Preliminary technical results.** Let  $r : \mathbb{R}^n \rightarrow \overline{\mathbb{R}}$ , where  $\overline{\mathbb{R}} := \mathbb{R} \cup \{+\infty\} \cup$   
 157  $\{-\infty\}$ , and let  $\bar{x}$  be such  $r(\bar{x}) \in \mathbb{R}$ . A vector  $v \in \mathbb{R}^n$  is in the regular (Fréchet)  
 158 subdifferential [30] (denoted by  $\hat{\partial}r(\bar{x})$ ) if  $r(x) \geq r(\bar{x}) + \langle v, x - \bar{x} \rangle + o(\|x - \bar{x}\|)$ .

159 Following again [30, Definition 8.2], we define the limiting (also known as Mor-  
 160 dukhovich) subdifferential, denoted by  $\partial r(\bar{x})$ . A vector  $v \in \mathbb{R}^n$  is a limiting subgra-  
 161 dient of  $r$  at  $\bar{x}$  (written  $v \in \partial r(\bar{x})$ ) if there exist sequences

$$162 \quad x^\nu \rightarrow_r \bar{x} \text{ and } v^\nu \in \hat{\partial}r(x^\nu) \text{ such that } v^\nu \rightarrow v.$$

163 Here,  $x^\nu \rightarrow_r \bar{x}$  denotes the  $r$ -attentive convergence and holds whenever  $x^\nu \rightarrow \bar{x}$  and  
 164  $r(x^\nu) \rightarrow r(\bar{x})$  [30, Def. 8.2]. The *horizon subdifferential* is

$$165 \quad \partial^\infty r(\bar{x}) = \left\{ v^\infty : \exists x^\nu \rightarrow_r \bar{x}, t_\nu \downarrow 0, \text{ and } v^\nu \in \hat{\partial}r(x^\nu) \text{ such that } t^\nu v_\nu \rightarrow v^\infty, \right\}.$$

166 For a closed set  $C \subset \mathbb{R}^n$ , the *limiting normal cone* is defined as  $N_C(x) := \partial\delta_C(x)$ ,  
 167 where  $\delta_C$  is the indicator of  $C$ , defined as  $\delta_C(x) = 0$  if  $x \in C$  and  $+\infty$  otherwise.

168 For a Lipschitz continuous function  $r$  we also use Clarke subdifferentials [30, 29, 7]  
 169 defined by

$$170 \quad \partial^C r(x) = \{v \in \mathbb{R}^n \mid r^\circ(x; d) \geq \langle v, d \rangle \text{ for all } d \in \mathbb{R}^n\}.$$

171 where  $r^\circ(x; d)$  is the directional Clarke derivative along direction  $d \in \mathbb{R}^n$ , defined as

$$172 \quad r^\circ(x; d) = \limsup_{\substack{y \rightarrow x \\ t \downarrow 0}} \frac{r(y + td) - r(y)}{t}.$$

173 The function  $r$  is called *prox-bounded* if there exists some  $\lambda > 0$  such that  $e_\lambda r(x) >$   
 174  $-\infty$  for some  $x \in \mathbb{R}^n$ . The prox-threshold of  $r$  is the supremum over all such  $\lambda > 0$ .  
 175 For lsc, proper functions, prox-boundedness is equivalent to the function  $x \mapsto r +$   
 176  $\frac{1}{2}\alpha\|\cdot\|^2$  being bounded from below for some  $\alpha > 0$  [30, Exercise 1.24]. We remark

177 that any function bounded from below is prox-bounded (with prox-threshold  $+\infty$ ).  
 178 Also, for every  $x \in \mathbb{R}^n$ ,  $\text{prox}_\lambda r(x)$  is nonempty and compact [30, Theorem 1.25] for  
 179 any  $\lambda > 0$  smaller than the prox-threshold.

180 We build on the concept of upper- $C^k$  ( $k \in \{1, 2, \dots\}$ ) functions in this work.  
 181 Following Rockafellar and Wets [30, Definition 10.29], a real-valued function  $f$  on an  
 182 open set  $O \in \mathbb{R}^n$  is upper- $C^k$  if locally  $r$  is the minimum of a family of  $C^k$  functions  
 183 indexed over a compact set, namely,  $r(x) = \min_{t \in T} f_t(x)$  for all  $x$  in a neighborhood  
 184 of  $\bar{x} \in O$ , with  $T$  compact,  $f_t \in C^k$ , and the values and all partial derivatives up  
 185 to order  $k$  depending jointly continuously on  $(t, x)$ . The minimum of finitely many  
 186  $C^k$  functions is therefore upper- $C^k$ , but may fail to be differentiable. There is no  
 187 distinction for orders  $k \geq 2$  since upper- $C^2$  implies upper- $C^\infty$  [30, Corollary 10.34].  
 188 Some prominent members of the upper- $C^2$  class are squared distance functions to  
 189 closed sets and Moreau envelopes of lsc, proper, prox-bounded functions for sufficiently  
 190 small regularization parameters.

191 A useful equivalent characterization of upper- $C^2$  functions is via local difference  
 192 of smooth and finite convex function, namely a function  $r$  is upper- $C^2$  if and only if  
 193 locally  $r = r_{sm} - r_{co}$  with  $r_{co}$  a finite convex function and  $r_{sm} \in C^2$  [30, Theorem  
 194 10.33]. In fact, one may take  $r_{sm} = \frac{\rho}{2} \|x\|_2^2$  locally [30, Theorem 10.33], which  
 195 is equivalent to  $r - \frac{\rho}{2} \|\cdot\|_2^2$  being concave locally. Recently, this property has been  
 196 exploited to devise simplified nonsmooth algorithms for upper- $C^2$  objectives that need  
 197 not maintain bundles, *e.g.*, [33]; we also leverage it here (see quadratic majorization  
 198 Lemma 4.3) to obtain convergence to stationary points of the regularized problems  
 199 (see Proposition 4.5) using proximal oracle subgradients (see Lemma 2.1 below).

200 Conversely, a function  $r$  is lower- $C^k$  if  $-r$  is upper- $C^k$ . Lower- $C^2$  functions enjoy  
 201 better regularity than upper- $C^2$ , in particular they are subdifferentially regular (also  
 202 known as lower regular) and, for example, the concepts of regular, limiting, and Clarke  
 203 subdifferentials coincide [30]. On the other hand, the strict inclusion  $\hat{\partial} \subset \partial \subset \partial^C$  can  
 204 hold for upper- $C^2$  functions in general.

205 A salient feature of Moreau envelopes of lsc, proper, and prox-bounded functions,  
 206 not available for general upper- $C^2$  functions, is that their limiting subdifferentials can  
 207 be characterized by means of proximal points as shown by the following lemma.

208 LEMMA 2.1. *Let  $r : \mathbb{R}^n \rightarrow (-\infty, +\infty]$  be proper lsc and prox-bounded, and let*  
 209  *$\lambda > 0$  be below the prox-threshold. We have the following characterizations of limiting*  
 210 *subdifferentials of the Moreau envelopes at a given  $x \in \mathbb{R}^n$ :*

- 211 (i)  $\text{prox}_\lambda r(x) \subset (I + \lambda \partial r)^{-1}(x)$ ;
- 212 (ii)  $\partial e_\lambda r(x) \subset \frac{1}{\lambda}(x - \text{prox}_\lambda r(x))$ ;
- 213 (iii) *For every  $v \in \partial e_\lambda r(x)$  there exists  $y \in \text{prox}_\lambda r(x)$  such that  $v = \frac{x-y}{\lambda}$  and*  
 214  *$v \in \partial r(y)$ .*

215 *Proof.* (i) and (ii) are Examples 10.2 and 10.32 from [30], respectively.  
 216 (iii) Consider  $v \in \partial e_\lambda r(x)$ . (ii) implies there exists  $y \in \text{prox}_\lambda r(x)$  such that  
 217  $v = \frac{x-y}{\lambda}$ . Also, (i) indicates that  $x \in y + \lambda \partial r(y)$ , which implies  $v \in \partial r(y)$  and  
 218 completes the proof.  $\square$

219 Lemma 2.1 (ii) above shows that any (limiting) subgradient  $v \in \partial e_\lambda r(x)$  cor-  
 220 responds to a minimizer  $y \in \text{prox}_\lambda r(x)$  of the Moreau envelope optimization prob-  
 221 lem (1.3). However, the converse is not necessarily true as the inclusion in (ii) can be  
 222 strict. Intuitively, this can occur, for example, when  $\text{prox}_\lambda r(x)$  is not a singleton, as  
 223 illustrated by the following example.

224 *Example 2.2.* Consider  $e_{\frac{1}{2}} r(x) = \frac{1}{2}x^2 - |x| + \frac{1}{2}$ , the Moreau envelope of  $r(z) =$

225  $\frac{1}{2} - z^2$  if  $|z| \leq \frac{1}{2}$  and  $r(z) = (|z| - 1)^2$  if  $|z| \geq \frac{1}{2}$ . For this, we observe that

226 
$$\partial e_{1/2}r(0) = \{-1, 1\} \subsetneq [-1, 1] = 2(0 - \text{prox}_{1/2}r(0)).$$

227 A complication caused by the possibly strict inclusion in (ii) of Lemma 2.1 is that  
 228 minimizers/proximal points  $y \in \text{prox}_\lambda r(x)$  of (1.3) do not necessarily provide a sub-  
 229 gradient in  $\partial e_\lambda r(x)$ . In the algorithm proposed in Section 4, we circumvent this  
 230 limitation by using the proximal oracle subgradients  $v = \frac{1}{\lambda}(x - y)$  instead of the true  
 231 limiting subgradients of  $\partial e_\lambda r$  at  $x$  to solve the Moreau regularized problem for fixed  
 232  $\lambda > 0$ . Notably, Lemma 2.1 (i) shows that these proximal oracle subgradients are  
 233 limiting subgradients of  $r$  at proximal points  $y \in \text{prox}_\lambda r(x)$ . This fact and the con-  
 234 vergence consistency between proximal points and algorithm iterates (see Lemma 4.8)  
 235 are key ingredients in obtaining (asymptotic) convergence of the algorithm to a lim-  
 236 iting stationary point to the target problem (1.1)–(1.2) without using the limiting  
 237 subgradients of the Moreau envelope approximants from (1.4)–(1.5).

238 *Remark 2.3.* If  $r$  is convex, or more generally weakly convex [11], then for  $\lambda$   
 239 small enough,  $\text{prox}_\lambda r$  is single valued and  $e_\lambda r$  is  $C^1$ , so Lemma 2.1 (ii) holds with  
 240 equality. For prox-regular [28] or lower- $C^2$  functions, same conclusion holds locally  
 241 on neighborhoods where the proximal mapping is single-valued [30].

242 In the remainder of the paper, we work under the following two assumptions.

243 *Assumption 2.4.* Functions  $f$  and  $h$  are  $C^2$  and functions  $r_i$  are lsc, finite valued,  
 244 and prox-bounded.

245 We remark that prox-boundedness is a general assumption (for example, holds for any  
 246 function bounded from below) and that we do not assume convexity or weak-convexity  
 247 (e.g., lower- $C^1$ , lower- $C^2$ ), prox-regularity, or differentiability of  $r_i$ .

248 Let us define the shorthand

249 
$$r_\lambda(x) = \sum_{i \in \mathcal{K}} e_\lambda r_i(x)$$

250 and denote the feasible set of the master problem (1.1)–(1.2) (also of (1.4)–(1.5)) by

251 
$$C := \{x \in \mathbb{R}^n : h(x) = 0, x_\ell \leq x \leq x_u\}.$$

252 For simplicity, we assume that the box bounds  $x_\ell$  and  $x_u$  are finite so that  $C$  is  
 253 compact. A more general assumption would be that the feasible sublevel sets of  $F$   
 254 are compact, but this would introduce considerable additional technical complexity  
 255 without adding much insight.

256 *Assumption 2.5.* The constraint set  $C$  is nonempty and compact.

257 **3. Moreau continuation: well-posedness and consistency.** This section is  
 258 concerned with the well-posedness of the continuation and proves that the minimizers  
 259 of the Moreau envelope-regularized problem (1.4)–(1.5) converge, and every cluster  
 260 point is a minimizer of the original problem (1.1)–(1.2) as  $\lambda \downarrow 0$ .

261 Let us define the extended-real objectives

262 
$$G(x) = F(x) + \delta_C(x) \text{ and } G_\lambda(x) = F_\lambda(x) + \delta_C(x).$$

263 **THEOREM 3.1.** *Choose  $\bar{\lambda} > 0$  below the prox-thresholds of  $r_i$ ,  $i \in \mathcal{K}$ , and consider*  
 264 *a sequence  $\{\lambda_k\}_{k \in \mathbb{N}}$  such that  $\lambda_k \downarrow 0$  and  $\lambda_k < \bar{\lambda}$ . Then the following properties hold*  
 265 *under Assumptions 2.4 and 2.5.*

- 266 (i)  $\{G_{\lambda_k}\}$  converges pointwise and monotonically to  $G$ ,  $G_{\lambda_k}(x) \uparrow G(x)$  for all  $x$ .  
 267 (ii)  $\{G_{\lambda_k}\}$  epi-converges to  $G$ , denoted as  $G_{\lambda_k} \xrightarrow{e} G$ .  
 268 (iii) Optimal values converge,

269 
$$\inf_{x \in C} F_{\lambda_k}(x) = \inf_x G_{\lambda_k}(x) \rightarrow \inf_x G(x) = \inf_{x \in C} F(x).$$

- 270 (iv) Moreover, for each  $k$ , the regularized problem (1.4)–(1.5) admits a global  
 271 minimizer, and the original problem (1.1)–(1.2) admits a global minimizer.  
 272 For any choice of minimizers  $x(\lambda_k)$  of (1.4)–(1.5), the sequence  $\{x(\lambda_k)\}$  is  
 273 bounded and every cluster point  $x^*$  is a minimizer of (1.1)–(1.2). If  $x^*$  is the  
 274 unique solution of (1.1)–(1.2), then  $x(\lambda_k) \rightarrow x^*$ .

275 *Proof.* (i) For  $x \notin C$ , one has  $G_{\lambda_k}(x) = G(x) = \infty$ . For  $x \in C$ , since  $r_i$  are  
 276 prox-bounded proper lsc functions, the Moreau envelope satisfies [30, Theorem 1.25]

277 
$$e_{\lambda_k} r_i(x) \uparrow r_i(x) \quad \text{as } \lambda_k \downarrow 0,$$

278 which immediately proves  $G_{\lambda_k}(x) \uparrow G(x)$ .

279 (ii) The sequence  $\{G_{\lambda_k}\}$  is pointwise nondecreasing and each  $G_{\lambda_k}$  is lsc, being the  
 280 sum of a smooth function, upper- $C^2$  functions, and the indicator function of the closed  
 281 set  $C$ . The monotone epi-limit equals the pointwise supremum of the lsc functions [30,  
 282 Proposition 7.4(d)], namely,  $G_{\lambda_k} \xrightarrow{e} G$ .

283 (iii) We first remark that the domains of  $G_{\lambda_k}$  are the compact set  $C$  for all  $k$   
 284 and hence are uniformly bounded. This ensures that all  $\{G_{\lambda_k}\}$  are eventually level-  
 285 bounded [30, Exercise 7.32(a)]. Using epi-convergence and eventual level-boundedness,  
 286 [30, Theorem 7.33] implies that  $\inf G_{\lambda_k} \rightarrow \inf G$ , which proves (iii). ■

287 (iv) Since  $C$  is compact and  $G_{\lambda_k}$  and  $G$  are lsc with effective domain  $C$ , the exist-  
 288 ence of minimizers for all regularized problems and for the original problem follows  
 289 from the Weierstrass theorem. Under Assumption 2.5,  $\{x(\lambda_k)\}$  is bounded and we  
 290 consider any cluster point  $x^*$  of it. Theorem 7.33 of [30] also implies that  $x^*$  is a  
 291 minimizer of  $G$ . If  $\operatorname{argmin} G$  is a singleton, then every convergent subsequence has  $x^*$   
 292 as the cluster point, hence  $\{x(\lambda_k)\} \rightarrow x^*$ . □

293 Theorem 3.1 can be viewed as a continuation consistency result, namely, *global*  
 294 minimizers of the regularized problems converge (along subsequences) to *global* min-  
 295 imizers of the original problem as  $\lambda \downarrow 0$ . A stronger result holds when the Moreau  
 296 envelope is over the whole objective, *i.e.*, one uses  $e_\lambda(f + \sum_{i \in \mathcal{K}} r_i)$  instead of regu-  
 297 larizing each term  $r_i$ . Minimizers of the whole-objective regularized problem coincide  
 298 with minimizers of  $f + \sum_{i \in \mathcal{K}} r_i$  for  $\lambda > 0$  [30, Example 1.46] under the assumption of  
 299 prox-boundedness. Furthermore, the local minimizers are also preserved under similar  
 300 assumptions as recently proved in [19].

301 *Remark 3.2.* In our setup, with the Moreau envelope distributed inside the min-  
 302 imization, such results involving local minimizers do not necessarily hold in the non-  
 303 convex case as illustrated by the following counterexample. Take  $f(x) = x^3 + x^4 - \frac{1}{2}x^2$ ,  
 304  $r_1(x) = \frac{1}{2}x^2$ , and  $\mathcal{K} = \{1\}$ . Using that  $e_\lambda r_1(x) = \frac{x^2}{2(1+\lambda)}$ , we see that  $F_\lambda$  has three  
 305 stationary points: a global minimizer  $x_1(\lambda) = \frac{1}{8} \left( -3 - \sqrt{\frac{9+25\lambda}{1+\lambda}} \right)$ , a local minimizer  
 306  $x_2(\lambda) = \frac{1}{8} \left( -3 + \sqrt{\frac{9+25\lambda}{1+\lambda}} \right)$ , and a local maximizer 0, for any  $\lambda > 0$ . As  $\lambda \downarrow 0$ ,  
 307  $x_1(\lambda)$  converges to  $x^* = -\frac{3}{4}$ , the global minimizer of  $F$ , consistent with Theorem 3.1;  
 308 however, the local minimizer  $x_2(\lambda)$  collapses to zero, which is a stationary point of  $F$ .

309 **4. Proposed algorithm.** In this section, we start by describing the continuation  
 310 algorithm, which is summarized in listing Algorithm 4.1. Then, in Sections 4.1 and  
 311 4.2, we analyze its convergence properties for the fixed- $\lambda$  problems (1.4)–(1.5) and of  
 312 the overall continuation scheme, respectively.

313 The continuation strategy consists of using a sequence of regularization param-  
 314 eters  $\lambda_s \downarrow 0$  and iteratively solving fixed- $\lambda$  problems (1.4)–(1.5) to approximate sta-  
 315 tionarity. In what follows, for each  $\lambda_s$  the iterations of the algorithm are denoted by  
 316  $\tilde{x}^k$ ,  $k = \{0, 1, 2, \dots\}$  and referred to as inner iterations. For better readability, the  
 317 iterates corresponding to each  $\lambda_s$  are denoted as  $x^s$  and referred to as outer iterates.

318 Consider at some fixed  $\lambda_s > 0$  the inner iterate  $\tilde{x}^k$ . A new inner iteration  $\tilde{x}^{k+1}$  is  
 319 obtained by means of an inexact stationary point  $\bar{x}^{k+1}$  of a NLP model approximation  
 320 of (1.4)–(1.5) in the form

$$321 \quad (4.1) \quad \min_{x \in \mathbb{R}^n} M(x; \tilde{x}^k, g^k; \lambda_s, \alpha_k) := f(x) + r_{\lambda_s}(\tilde{x}^k) + (g^k)^\top (x - \tilde{x}^k) + \alpha_k \|x - \tilde{x}^k\|^2$$

$$322 \quad (4.2) \quad \text{s.t.} \quad h(x) = 0, \quad x_\ell \leq x \leq x_u.$$

323 Similar to the SQP approach of Wang and Petra [33] we use a convex quadratic model  
 324 for the nonsmooth objective that uses first-order derivative information  $g^k$  defined  
 325 below in (4.4) and a quadratic penalty term for step/update stabilization controlled  
 326 by the scalar  $\alpha_k > 0$ . At variance with the SQP approach, we enforce the smooth  
 327 constraints (4.2) exactly. This is done for practical considerations since it allows  
 328 using off-the-shelf NLP solvers to compute the new iterate, but it also simplifies the  
 329 analysis. We believe that the SQP approach [33] can be revisited and used to replace  
 330 the NLP approach present here at the cost of increased complexity in the analysis and  
 331 implementation since it requires using feasibility restoration to deal with potentially  
 332 inconsistent SQP linearizations of the constraints. In contrast, the sequential NLP  
 333 approach here leaves this complication to the underlying NLP solver.

334 We do not require that the candidate  $\bar{x}^{k+1}$  be a local or global minimizer of the  
 335 NLPs (4.1) – (4.2) since this may be unrealistic. Instead, we only require a stationary  
 336 point  $\tilde{x}^{k+1}$  of the NLP model that satisfies a weak descent condition

$$337 \quad (4.3) \quad M(\tilde{x}^{k+1}; \tilde{x}^k, g^k; \lambda_s, \alpha_k) \leq M(\tilde{x}^k; \tilde{x}^k, g^k; \lambda_s, \alpha_k).$$

338 We remark that there is always a path in  $C$  along which the descent condition holds,  
 339 unless  $\tilde{x}^k$  is a stationary point. Indeed, if  $\tilde{x}^k$  is not a critical point of the NLP  
 340 model, then there must exist a direction  $d$  in the tangent cone of  $C$  to  $\tilde{x}^k$  such that  
 341  $d^\top \nabla M(\tilde{x}^k) < 0$ ; otherwise, if every tangent direction  $d$  satisfies  $d^\top \nabla M(\tilde{x}^k) \geq 0$ ,  
 342 we would have  $-\nabla M(\tilde{x}^k) \in N_C(\tilde{x}^k)$ , meaning  $\tilde{x}^k$  is stationary. Moreover, using  
 343 definition of tangents, there is a path  $x(\beta) = \tilde{x}^k + \beta d + o(\beta) \in C$ ,  $\beta \in [0, \bar{\beta}]$  such  
 344 that  $M(x(\beta); \tilde{x}^k, g^k; \lambda, \alpha) = M(\tilde{x}^k; \tilde{x}^k, g^k; \lambda, \alpha) + \beta d^\top \nabla M(\tilde{x}^k) + o(\beta) < M(\tilde{x}^k)$  for all  
 345  $\beta \in [0, \bar{\beta}]$  and small enough  $\bar{\beta} > 0$ . Therefore, the NLP model has a feasible path of  
 346 descent and the descent condition (4.3) is a mild requirement for the underlying NLP  
 347 solver and less intrusive than requiring a minimizer.

348 Above,  $g^k$  encapsulates “derivative“ information, namely,

$$349 \quad (4.4) \quad g^k = \sum_{i \in \mathcal{K}} g_i^k \quad \text{with} \quad g_i^k = \frac{1}{\lambda_s} (\tilde{x}^k - w_i^k) \quad \text{and} \quad w_i^k \in \text{prox}_{\lambda_s} r_i(\tilde{x}^k).$$

350 We use slightly less regular and more general subgradients than Wang and Petra [33],  
 351 namely, we rely on proximal oracle subgradients  $g_i^k$ , which are limiting subgradients

352 at “nearby” points (see Lemma 2.1) and are defined for discontinuous functions, while  
 353 the mentioned previous work employs Clarke subgradients. The proximal oracle sub-  
 354 gradients used in  $g^k$  above are a subset of the Clarke subgradients of the Moreau  
 355 envelopes in general, see Theorem 4.6 and its proof of Clarke stationarity. The choice  
 356 of proximal oracle subgradients  $g_i^k$ s (as opposed to choosing  $g_i^k$  more restrictively,  
 357  $g_i^k \in \partial e_{\lambda_s} r_i(\tilde{x}^k)$ ) stems from the representations of subdifferentials of Moreau en-  
 358 velopes of general functions satisfying Assumption 2.4, as discussed after Lemma 2.1.

359 We use an accept-reject rule based on a majorization-ratio test that updates the  
 360 penalty parameter  $\alpha^k > 0$  in the NLP model solved at each iteration. Upon solving  
 361 the NLP model (4.1)–(4.2), we define the predicted and actual reductions at the  
 362 candidate  $\bar{x}^{k+1}$  as

$$363 \quad (4.5) \quad \begin{aligned} \text{Pred}_k(\bar{x}^{k+1}) &= M(\tilde{x}^k; \tilde{x}^k, g^k; \lambda, \alpha_k) - M(\bar{x}^{k+1}; \tilde{x}^k, g^k; \lambda, \alpha_k) \text{ and} \\ \text{Ared}_k(\bar{x}^{k+1}) &= F_\lambda(\tilde{x}^k) - F_\lambda(\bar{x}^{k+1}). \end{aligned}$$

364 We define the ratio

$$365 \quad \rho_k = \frac{\text{Ared}_k(\bar{x}^{k+1})}{\text{Pred}_k(\bar{x}^{k+1}) + \sigma \|\bar{x}^{k+1} - \tilde{x}^k\|^2}.$$

366 This is analogous to trust-region acceptance, but specialized to our exact con-  
 367 straints model. The iterate  $\tilde{x}^{k+1} = \bar{x}^{k+1}$  is accepted whenever  $\rho_k \geq \bar{\rho}$  with  $\bar{\rho} \in (0, 1)$ ;  
 368 otherwise, the step is rejected and  $\alpha_k$  increases by a factor of two.

369 At each outer iteration  $s$ , the inner iterates  $\tilde{x}^k$  of Algorithm 4.1 (steps 5 – 20) aim  
 370 to reach an approximate stationary point for (1.4)–(1.5). Inexactness is controlled by  
 371 the tolerance  $\varepsilon_s$  and has the goal to balance inner solve accuracy with the continuation  
 372 error. The stopping criterion of the inner iterates use the lagged stationarity measure

$$373 \quad \hat{R}(\tilde{x}^k) := \nabla f(\tilde{x}^k) + g^{k-1} + \nabla h(\tilde{x}^k)^\top \nu^k - \gamma^k + \theta^k,$$

374 where  $\nu^k \in \mathbb{R}^m$ ,  $\gamma^k \geq 0$ , and  $\theta^k \geq 0$  are the KKT multipliers corresponding to the  
 375 optimality conditions (see (4.11)–(4.13) for the exact form) of the model NLPs (4.1)  
 376 – (4.2) under the general MFCQ for  $C$ . From these KKT conditions, we remark that  
 377  $\hat{R}(\tilde{x}^k) = -2\alpha_{k-1}(\tilde{x}^k - \tilde{x}^{k-1})$ , which makes this lagged stationarity both computationally  
 378 cheap and practical (since it does not require dual information from the numerical  
 379 solver used for the master problem). Therefore, the inner loop stopping criterion is

$$380 \quad (4.6) \quad \|\hat{R}(\tilde{x}^k)\|_\infty = 2\alpha_{k-1} \|\tilde{x}^k - \tilde{x}^{k-1}\|_\infty \leq \varepsilon_s.$$

381 Later, in Proposition 4.5 we show that the inner loop reaches (4.6) for  $\varepsilon_s > 0$  in a  
 382 finite number of iterations. Theorem 4.6 analyzes the infinite inner loop (without the  
 383 stopping criterion) and shows that the inner iterates  $\tilde{x}^k$  are bounded and any of their  
 384 accumulation or cluster points  $x^{*\lambda}$  are stationary points of (1.4)–(1.5).

385 Algorithm 4.1 repeats the sequential NLP steps discussed above for a given se-  
 386 quence  $\lambda_s \downarrow 0$ . It is shown in Theorem 4.11 that under mild conditions the outer  
 387 iterates  $x^s$  of Algorithm 4.1 are bounded and any of their accumulation points  $x^*$   
 388 satisfy stationarity conditions

$$389 \quad (4.7) \quad 0 \in \nabla f(x^*) + \sum_{i \in \mathcal{K}} \partial r_i(x^*) + \nabla h(x^*)^\top \nu^* - \gamma^* + \theta^*,$$

$$390 \quad (4.8) \quad h(x^*) = 0, \quad 0 \leq \gamma^* \perp (x^* - \ell) \geq 0, \quad 0 \leq \theta^* \perp (u - x^*) \geq 0,$$

391 for some multipliers  $\nu^* \in \mathbb{R}^m$ ,  $\gamma^* \geq 0$ , and  $\theta^* \geq 0$  under MFCQ for the feasible set  $C$ .  
 392 The limiting stationarity condition (4.7) targeted here is analogous in form to standard

393 first-order necessary optimality conditions, see, for example, [30]. Consequently, we  
 394 use the following basic CQ (BCQ) specific to non-Lipschitz functions, counterpart of  
 395 the horizon BCQ from [30, Theorem 8.15].

396 **BCQ for non-Lipschitz objectives.** If there exist  $v_i \in \partial^\infty r_i(x^*)$ ,  $i \in \mathcal{K}$ , such that  
 397  $-\sum_{i \in \mathcal{K}} v_i \in N_C(x^*)$ , then  $v_i = 0$  for all  $i \in \mathcal{K}$ . Under MFCQ, the BCQ is equivalent  
 398 to: if there exist  $v_i \in \partial^\infty r_i(x^*)$ ,  $i \in \mathcal{K}$ , and  $\nu \in \mathbb{R}^m$ ,  $\gamma \geq 0$ , and  $\theta \geq 0$  satisfying  
 399 complementarity conditions from (4.8) such that  $\sum_{i \in \mathcal{K}} v_i + \nabla h(x^*)^\top \nu - \gamma + \theta = 0$ ,  
 400 then  $v_i = 0$  for all  $i \in \mathcal{K}$ .

---

**Algorithm 4.1** Proximal oracle sequential NLP continuation algorithm

---

```

1: Initial point  $x^0 \in C$  and  $\lambda_0$  below the prox-thresholds of all  $r_i$ ; parameters  $\bar{\rho} \in$ 
   (0, 1],  $\sigma > 0$ , and  $\underline{\alpha} > 0$ ; sequences of tolerances  $\varepsilon_s \downarrow 0$  and regularizations  $\lambda_s \downarrow 0$ .
2:  $s \leftarrow 0$ 
3: while true do
4:   Inner loop initializations:  $k \leftarrow 0$ ,  $\tilde{x}^0 \leftarrow x^s$ ,  $\alpha_k \leftarrow \max\{\underline{\alpha}, \frac{|K|}{2\lambda_s}\}$ 
5:   while true do
6:     for  $i \in \mathcal{K}$  do
7:       Compute  $w_i^k \in \text{prox}_{\lambda_s} r_i(\tilde{x}^k)$  and  $g_i^k = (\tilde{x}^k - w_i^k)/\lambda_s$ 
8:     end for
9:      $g^k \leftarrow \sum_{i \in \mathcal{K}} g_i^k$ 
10:    Compute a stationary point  $\bar{x}^{k+1}$  of  $M(\cdot; \tilde{x}^k, g^k; \lambda_s, \alpha_k)$  from (4.1)–(4.2) sat-
    isfying (4.3)
11:    Compute  $\rho_k(\bar{x}^{k+1}) = \text{Ared}_k(\bar{x}^{k+1})/(\text{Pred}_k(\bar{x}^{k+1}) + \sigma \|\bar{x}^{k+1} - \tilde{x}^k\|^2)$  using (4.5)
    if  $\bar{x}^{k+1} \neq \tilde{x}^k$ , otherwise set  $\rho_k(\bar{x}^{k+1}) = 1$ 
12:    if  $\rho_k(\bar{x}^{k+1}) \geq \bar{\rho}$  then
13:       $\tilde{x}^{k+1} \leftarrow \bar{x}^{k+1}$ ,  $\alpha_{k+1} \leftarrow \alpha_k$ ,  $k \leftarrow k + 1$  {accept iterate}
14:    else
15:       $\tilde{x}^{k+1} \leftarrow \tilde{x}^k$ ,  $\alpha_{k+1} \leftarrow 2\alpha_k$  {reject iterate}
16:       $g^{k+1} \leftarrow g^k$ ,  $g_i^{k+1} \leftarrow g_i^k$ ,  $w_i^{k+1} \leftarrow w_i^k$ ,  $i \in \mathcal{K}$ ,  $k \leftarrow k + 1$ 
17:      Return to step 10
18:    end if
19:    if  $\|\hat{R}(\tilde{x}^k)\|_\infty \leq \varepsilon_s$  then
20:      break {stopping criterion for inner loop}
21:    end if
22:  end while
23:   $s \leftarrow s + 1$ ,  $x^s \leftarrow \tilde{x}^k$ ,  $g_i^s \leftarrow g_i^{k-1}$ ,  $w_i^s \leftarrow w_i^{k-1}$  for all  $i \in \mathcal{K}$ 
24: end while

```

---

401 BCQ rules out abnormal stationarity conditions and is a prerequisite for nec-  
 402 essary first-order optimality conditions (4.7)–(4.8) to hold for constrained problems  
 403 with nonsmooth and nonconvex objectives. In this respect, it can be viewed as the  
 404 counterpart of KKT conditions, which rule out Fritz-John points with zero objective  
 405 multipliers, in smooth constrained optimization. Geometrically, near  $x^*$ , limiting sub-  
 406 gradients can become arbitrarily large, and their normalized limiting directions belong  
 407 to the horizon subgradients  $\partial^\infty r_i(x^*)$  (for some  $i \in \mathcal{K}$ ). If one such direction exists  
 408 and is opposite to the normal cone of  $C$  at some boundary point, or, equivalently,  
 409 cancels out the multiplier terms in the BCQ above, then the constraint boundary  
 410 can absorb this unbounded tilt. In such cases general first-order stationarity condi-  
 411 tions such as [30, Theorem 8.15] may fail at minimizers as shown in Examples 4.1

412 and 4.2. The BCQ excludes exactly this pathology. We remark that the BCQ holds  
 413 automatically for locally Lipschitz functions  $r_i$ , since  $\partial^\infty r_i(x^*) = \{0\}$ .

414 *Example 4.1.* Consider  $\min_x \{r_1(x) := x^{1/3} \mid x \geq 0\}$  with the (global) minimizer  
 415  $x^* = 0$ . This problem fails BCQ above since there exists  $v = \frac{1}{3} \in \partial^\infty r_1(0)$  such that  
 416  $v - \gamma = 0$ , where  $0 \leq x \perp \gamma \geq 0$ . We remark that  $\partial r_1(0) = \emptyset$  and therefore the nec-  
 417 essary first-order stationary conditions (4.7)–(4.8) for  $f = 0$  or ones in [30][Theorem  
 418 8.15]) do not hold at this minimizer.

419 *Example 4.2.* The following example shows that (global) minimizers fail to satisfy  
 420 first-order stationarity when the BCQ does not hold even when the subdifferential is  
 421 nonempty and the feasible set satisfies MFCQ. Consider

$$422 \quad \min_{x \in \mathbb{R}^2} \{r_1(x_1, x_2) := x_1^2 + x_1 + \sqrt{\max(x_2, 0)} \mid x \in C\},$$

423 with  $C := \{x \in \mathbb{R}^2 : h(x) := x_1^2 - x_2 \leq 0\}$ . Since  $r_1(x_1, x_2) = x_1^2 + x_1 + \sqrt{x_2} \geq x_1^2 +$   
 424  $x_1 + |x_1|$  over  $x = (x_1, x_2) \in C$ , we see that  $x^* = (0, 0)$  is a strict local/global minimizer  
 425 over  $C$ . We remark that MFCQ holds at  $x^*$  and that  $\partial r_1(0, 0) = \partial r(0, 0) = \{1\} \times \mathbb{R}_+$ .

426 BCQ fails at  $x^*$  since  $0 \neq v = (0, 1) \in \partial^\infty r(0, 0)$  satisfies  $v + \nu(0, -1) = 0$   
 427 with  $0 \leq \nu \perp -h(0, 0) \geq 0$  (for  $\nu = 1$ ). Here we remark that  $\nu$  is signed since  
 428  $h(x)$  is involved in an inequality constraint and, also, the bound multipliers  $\gamma$  and  
 429  $\theta$  are not present. To see that  $(0, 1) \in \partial^\infty r_1(0, 0)$ , take  $x_\tau = (0, \tau)$  with  $\tau \downarrow 0$ . As  
 430  $x_\tau \rightarrow (0, 0)$ , we have  $r_1(x_\tau) = \sqrt{\tau} \rightarrow 0$  and  $v^\tau := \nabla r(0, \tau) = (1, \frac{1}{2\sqrt{\tau}}) \in \hat{\partial} r_1(x_\tau)$ . By  
 431 choosing  $t_\tau = 2\sqrt{\tau} \rightarrow 0$  we obtain  $t_\tau v^\tau = (2\sqrt{\tau}, 1) \rightarrow (0, 1)$ . Then  $(0, 1) \in \partial^\infty r_1(0, 0)$   
 432 by definition of horizon subdifferential from Section 2. In fact, one can prove that  
 433  $\partial^\infty r_1(0, 0) = \{0\} \times \mathbb{R}_+$ .

434 Even though  $(0, 0)$  is a minimizer the stationarity conditions do not hold since

$$435 \quad 0 \notin \partial r_1(0, 0) + \nu(0, -1) = \{1\} \times \mathbb{R}_+ + \nu(0, -1)$$

436 for any  $\nu$  such that  $0 \leq \nu \perp h(0, 0) \geq 0$ .

437 We use the concept of sum-stationarity using  $\sum_{i \in \mathcal{K}} \partial r_i$ , which comes naturally  
 438 in our composite objective, instead of  $\partial(\sum_{i \in \mathcal{K}} r_i)$  commonly used with weakly con-  
 439 vex summands. To the best of our knowledge, no inclusion is guaranteed between  
 440 the two sets in our general setup without additional qualifications. We show in  
 441 Lemma 4.12 that  $\sum_{i \in \mathcal{K}} \partial r_i \subseteq \partial(\sum_{i \in \mathcal{K}} r_i)$  for subdifferentially regular functions, which  
 442 cover weakly convex functions, hence the stationarity conditions used in this work  
 443 characterize the typical limiting stationarity for this class of problems.

444 **4.1. Inner-loop convergence.** We start with a majorization lemma specific to  
 445 upper- $C^2$  functions [33]. Here we provide a stronger version for Moreau envelopes with  
 446 an explicit bound in terms of the regularization  $\lambda$  and proximal oracle subgradients.

447 **LEMMA 4.3.** *For any  $\lambda > 0$  below the prox-thresholds of  $r_i(\cdot)$  and a given  $\bar{x}$ , the*  
 448 *following majorization bound holds: for any  $\bar{g}_i \in \frac{1}{\lambda}(\bar{x} - \text{prox}_\lambda r_i(\bar{x}))$ ,  $i \in \mathcal{K}$ ,*

$$449 \quad r_\lambda(x) \leq r_\lambda(\bar{x}) + \langle \bar{g}, x - \bar{x} \rangle + Q_\lambda \|x - \bar{x}\|^2, \quad \forall x, \quad \text{where } \bar{g} = \sum_{i \in \mathcal{K}} \bar{g}_i \text{ and } Q_\lambda = \frac{|\mathcal{K}|}{2\lambda}.$$

450 *Proof.* Fix  $i \in \mathcal{K}$  and consider  $\bar{g}_i \in \frac{1}{\lambda}(\bar{x} - \text{prox}_\lambda r_i(\bar{x}))$ , namely, there exists  $w_i \in$   
 451  $\text{prox}_\lambda r_i(\bar{x})$  such that  $\bar{g}_i = \frac{1}{\lambda}(\bar{x} - w_i)$ . The optimality of  $w_i$  with respect to the  
 452 minimization problem defining the Moreau envelope implies  $e_{\lambda_i} r_i(x) \leq r_i(w_i) + \frac{1}{2\lambda} \|x -$

453  $w_i\|^2$  for all  $x$ . Also, for  $x = \bar{x}$  we have  $e_{\lambda_i} r_i(\bar{x}) = r_i(w_i) + \frac{1}{2\lambda} \|\bar{x} - w_i\|^2$ . By subtracting  
 454 these two conditions, we obtain after some manipulations of the quadratic terms that

$$\begin{aligned}
 455 \quad e_{\lambda_i} r_i(x) &\leq e_{\lambda_i} r_i(\bar{x}) + \frac{1}{2\lambda} \langle x - \bar{x} + \bar{x} - w_i, x - \bar{x} + \bar{x} - w_i \rangle - \frac{1}{2\lambda} \|\bar{x} - w_i\|^2 \\
 456 \quad &= e_{\lambda_i} r_i(\bar{x}) + \left\langle \frac{1}{\lambda} (\bar{x} - w_i), x - \bar{x} \right\rangle + \frac{1}{2\lambda} \|x - \bar{x}\|^2. \quad \square
 \end{aligned}$$

457 Finally, summing over  $i \in \mathcal{K}$  gives the desired majorization bound.

458 The above lemma is the key to proving the convergence of the proposed algorithm.  
 459 First, we use it to prove that the candidate iterates  $\bar{x}^{k+1}$  are eventually accepted.

460 **LEMMA 4.4.** *Fix an outer iteration  $s$ . Then only finitely many inner steps can be*  
 461 *rejected in Algorithm 4.1. Consequently, the penalty parameters satisfy  $\alpha_k \in [\underline{\alpha}, \bar{\alpha}_s]$ ,*  
 462 *with  $\bar{\alpha}_s = \max\{\underline{\alpha}, 2(Q_\lambda + \bar{\rho}\sigma)\}$ , for all  $k$  greater than some iteration threshold index.*

463 *Proof.* We first show that whenever  $\alpha_k$  becomes larger than  $Q_\lambda + \bar{\rho}\sigma$ , all sub-  
 464 sequent iterations are accepted. Assume that at iterate  $k$ , the trial point  $\bar{x}^{k+1}$  is  
 465 computed with  $\alpha_k > Q_\lambda + \bar{\rho}\sigma$ . We use the majorization of Lemma (4.3) to write

$$\begin{aligned}
 466 \quad F_\lambda(\bar{x}^{k+1}) &= f(\bar{x}^{k+1}) + r_\lambda(\bar{x}^{k+1}) \leq f(\bar{x}^{k+1}) + r_\lambda(\tilde{x}^k) + \langle g^k, \bar{x}^{k+1} - \tilde{x}^k \rangle + \\
 467 \quad &Q_\lambda \|\bar{x}^{k+1} - \tilde{x}^k\|^2 = M_k(\bar{x}^{k+1}; \tilde{x}^k, g^k; \lambda, \alpha_k) + (Q_\lambda - \alpha_k) \|\bar{x}^{k+1} - \tilde{x}^k\|^2.
 \end{aligned}$$

468 By subtracting the above from  $F_\lambda(\tilde{x}^k) = M_k(\tilde{x}^k; \tilde{x}^k, g^k; \lambda, \alpha_k)$  and using our working  
 469 assumption that  $\alpha_k > Q_\lambda + \bar{\rho}\sigma$ , we obtain

$$\begin{aligned}
 470 \quad \text{Ared}_k(\bar{x}^{k+1}) &\geq \text{Pred}_k(\bar{x}^{k+1}) + (\alpha_k - Q_\lambda) \|\bar{x}^{k+1} - \tilde{x}^k\|^2 \\
 471 \quad &> \text{Pred}_k(\bar{x}^{k+1}) + \bar{\rho}\sigma \|\bar{x}^{k+1} - \tilde{x}^k\|^2 \geq \bar{\rho}(\text{Pred}_k(\bar{x}^{k+1}) + \sigma \|\bar{x}^{k+1} - \tilde{x}^k\|^2), \blacksquare
 \end{aligned}$$

472 with the last inequality due to the algorithmic parameter  $\bar{\rho} \leq 1$  and  $\text{Pred}_k(\bar{x}^{k+1}) \geq 0$ .  
 473 Therefore,  $\rho_k(\bar{x}^{k+1}) \geq \bar{\rho}$  and  $\bar{x}^{k+1}$  will be accepted in the acceptance test from step 13  
 474 of Algorithm 4.1, namely  $\tilde{x}^{k+1} = \bar{x}^{k+1}$ , and, also,  $\alpha_{k+1} = \alpha_k$ . Proceeding by induction  
 475 through the next iterations, it should be apparent that all iterations  $k + j$  with  $j \geq 1$   
 476 of Algorithm 4.1 accept  $\bar{x}^{k+j}$  and  $\alpha_{k+j} = \alpha_k$  for all  $j \geq 1$ .

477 We remark that the inner loop of Algorithm 4.1 may satisfy the termination  
 478 criteria in step 20 and exit after finite number of iterations without increasing the  
 479 penalty parameters  $\alpha_k$  above the  $Q_\lambda + \bar{\rho}\sigma$  threshold. In this case, it should be apparent  
 480 that  $\alpha_k \in [\underline{\alpha}, Q_\lambda + \bar{\rho}\sigma]$  for all  $k$ . If this is not the case, let  $k_0$  denote the smallest  
 481 iteration index such that  $\alpha_{k_0} > Q_\lambda + \bar{\rho}\sigma$ . This index must be finite because of the  
 482 geometric increase of penalties  $\alpha_k$  for the rejected iterates.

483 If  $k_0 = 0$ , then  $\alpha_{k_0} = \max\{\underline{\alpha}, Q_\lambda\} \in [\underline{\alpha}, \max\{\underline{\alpha}, Q_\lambda\}]$  due to initialization. If  
 484  $k_0 \geq 1$ , then we immediately see that the previous step must have been rejected,  
 485 therefore,  $\alpha_{k_0-1} = \frac{1}{2}\alpha_{k_0} < Q_\lambda + \bar{\rho}\sigma$ , which implies that  $\alpha_{k_0} < 2(Q_\lambda + \bar{\rho}\sigma)$ . Finally,  
 486 we remark that  $\alpha_{k_0} = \alpha_{k_0+j} = \dots < 2(Q_\lambda + \bar{\rho}\sigma)$  for all  $j \geq 1$ , as per the first part of  
 487 the proof.

488 The expression for  $\bar{\alpha}_s$  follows after collecting the lower and upper bounds for  $\alpha_k$   
 489 over the above three cases.  $\square$

490 **PROPOSITION 4.5.** *At any outer iteration  $s$ , the inner iterates  $\{\tilde{x}^k\}$  of Algo-*  
 491 *algorithm 4.1 are bounded and*

492 (i)  $\{F_\lambda(\tilde{x}^k)\}$  is decreasing monotonically:

$$493 \quad F_\lambda(\tilde{x}^{k+1}) - F_\lambda(\tilde{x}^k) \leq -\bar{\rho}\sigma \|\tilde{x}^{k+1} - \tilde{x}^k\|^2, \forall k \geq 0;$$

- 494 (ii) the updates  $\Delta^k = \tilde{x}^{k+1} - \tilde{x}^k$  are square-summable, that is,  $\sum_k \|\Delta_k\|^2 < \infty$ ,  
 495 and, thus,  $\Delta_k \rightarrow 0$ ;  
 496 (iii) for any tolerance  $\varepsilon_s > 0$ , the inner-loop stopping criterion (4.6) is reached in  
 497 a finite number of inner iterations  $k_s$ .

498 *Proof.* We first remark that the sequence  $\{\tilde{x}^k\}$  is bounded since it is generated  
 499 by the algorithm in the feasible set  $C$ , which is compact by the Assumption 2.5.

500 (i) By Lemma 4.4 all iterates  $k$  greater than some index  $k_0$  are accepted in step 13  
 501 of the algorithm, hence

$$502 \quad F_\lambda(\tilde{x}^k) - F_\lambda(\tilde{x}^{k+1}) \geq \bar{\rho} (\text{Pred}_k(\tilde{x}^{k+1}) + \sigma \|\tilde{x}^{k+1} - \tilde{x}^k\|^2) \geq \bar{\rho} \sigma \|\tilde{x}^{k+1} - \tilde{x}^k\|^2.$$

503 The iterates 0 to  $k_0$  also satisfy the monotonicity: the accepted iterates in the form  
 504 above, while the rejected iterates trivially, since  $\tilde{x}^{k+1} = \tilde{x}^k$ .

505 (ii) Let us sum over  $k \in [k_0, \bar{k}]$  the inequality from (i):

$$506 \quad \bar{\rho} \sigma \sum_{k=k_0}^{\bar{k}} \|\Delta_k\|^2 \leq F_\lambda(\tilde{x}^{k_0}) - F_\lambda(\tilde{x}^{\bar{k}}).$$

507 Since  $f(\cdot)$  is  $C^2$  and the Moreau envelopes defining  $r_\lambda(\cdot)$  are upper- $C^2$ , hence  
 508 Lipschitz continuous, the function  $F_\lambda(\cdot)$  is continuous. Therefore, boundedness of  
 509  $\{\tilde{x}^k\}$  implies that  $\{F_\lambda(\tilde{x}^k)\}$  is bounded. Taking the limit  $\bar{k} \rightarrow \infty$  we obtain that  
 510  $\bar{\rho} \sum_k \|\Delta_k\|^2 < \infty$ . This completes the proof of (ii) since  $\bar{\rho} \sigma > 0$ .

511 (iii) First observe that  $\|\hat{R}(\tilde{x}^k)\| = 2\alpha_k \|\tilde{x}^k - \tilde{x}^{k-1}\| \leq 2\bar{\alpha}_s \|\Delta_{k-1}\|$  by Lemma 4.4.  
 512 Using  $\Delta_k \rightarrow 0$  from (ii), stopping criterion (4.6) is going to be met for sufficiently  
 513 large  $k$ , so the inner loop terminates after finitely many iterations.  $\square$

514 The convergence behavior of the inner loop (for fixed  $\lambda$ ) of Algorithm 4.1 fits  
 515 under a standard subsequence convergence and only requires a standard CQ on the  
 516 feasible set of the target problem. We assume MFCQ.

517 **THEOREM 4.6.** *Assume MFCQ holds for the feasible set  $C$ , in addition to As-*  
 518 *sumptions 2.4 and 2.5. Fix an outer iteration  $s$  and consider the sequence of inner*  
 519 *iterates  $\{\tilde{x}^k\}$  generated by the inner loop of of Algorithm 4.1 without the stopping test*  
 520 *in step 20, i.e., the infinite loop consisting of steps 5–17. Then  $\{\tilde{x}^k\}$  is bounded and*  
 521 *all its cluster points  $x^{*\lambda}$  satisfy the proximal limiting stationarity conditions*

$$522 \quad (4.9) \quad 0 = \nabla f(x^{*\lambda}) + \sum_{i \in \mathcal{K}} g_i^{*\lambda} + \nabla h(x^{*\lambda})^\top \nu^{*\lambda} - \gamma^{*\lambda} + \theta^{*\lambda},$$

$$523 \quad (4.10) \quad h(x^{*\lambda}) = 0, \quad 0 \leq \gamma^{*\lambda} \perp (x^{*\lambda} - \ell) \geq 0, \quad 0 \leq \theta^{*\lambda} \perp (u - x^{*\lambda}) \geq 0.$$

524 for some multipliers  $\nu^{*\lambda} \in \mathbb{R}^m$ ,  $\gamma^{*\lambda} \geq 0$ , and  $\theta^{*\lambda} \geq 0$  and with  $g_i^{*\lambda} = \frac{1}{\lambda}(x^{*\lambda} - w_i^{*\lambda}) \in$   
 525  $\partial r_i(w_i^{*\lambda})$ , where  $w_i^{*\lambda} \in \text{prox}_\lambda r_i(x^{*\lambda})$ .

526 Moreover,  $x^{*\lambda}$  also satisfies Clarke stationarity conditions

$$527 \quad 0 \in \nabla f(x^{*\lambda}) + \partial^C r_\lambda(x^{*\lambda}) + \nabla h(x^{*\lambda})^\top \nu^{*\lambda} - \gamma^{*\lambda} + \theta^{*\lambda},$$

$$528 \quad h(x^{*\lambda}) = 0, \quad 0 \leq \gamma^{*\lambda} \perp (x^{*\lambda} - \ell) \geq 0, \quad 0 \leq \theta^{*\lambda} \perp (u - x^{*\lambda}) \geq 0.$$

529 for some multipliers  $\nu^{*\lambda}$ ,  $\gamma^{*\lambda} \geq 0$ , and  $\theta^{*\lambda} \geq 0$ .

530 *Proof.* The boundedness of  $\{\tilde{x}^k\}$  from Proposition 4.5 shows that cluster points of  
 531  $\{\tilde{x}^k\}$  exist. Let  $x^{*\lambda}$  be an arbitrary such point and, without re-indexing, let  $\tilde{x}^k \rightarrow x^{*\lambda}$ .

532 By Lemma 4.4,  $\tilde{x}^{k+1}$  becomes an accepted iterate for large enough  $k$  and therefore  
 533 satisfies the KKT conditions for the sequential NLP model (4.1)-(4.2), namely, there  
 534 exist multipliers  $\nu^{k+1}$ ,  $\gamma^{k+1}$ , and  $\theta^{k+1}$  such that

$$535 \quad (4.11) \quad \nabla f(\tilde{x}^{k+1}) + g^k + 2\alpha_k \Delta^k + \nabla h(\tilde{x}^{k+1})^\top \nu^{k+1} - \gamma^{k+1} + \theta^{k+1} = 0,$$

$$536 \quad (4.12) \quad h(\tilde{x}^{k+1}) = 0,$$

$$537 \quad (4.13) \quad 0 \leq \gamma^{k+1} \perp (\tilde{x}^{k+1} - \ell) \geq 0, \quad 0 \leq \theta^{k+1} \perp (u - \tilde{x}^{k+1}) \geq 0.$$

538 The main steps of the proof of limiting stationarity equations (4.9)-(4.10) are:

- 539 (a) show the sequence  $\{g^k\}_k$  is bounded and its cluster points  $g^{\lambda^*}$  satisfy  $g^{\lambda^*} =$   
 540  $\sum_{i \in \mathcal{K}} g_i^{\lambda^*}$  with  $g_i^{\lambda^*} = \frac{1}{\lambda}(x^{\lambda^*} - w_i^{\lambda^*}) \in \partial r_i(w_i^{\lambda^*})$  and  $w_i^{\lambda^*} \in \text{prox}_{\lambda} r_i(x^{\lambda^*})$ , for  
 541 all  $i \in \mathcal{K}$ .  
 542 (b) show the sequences  $\{\nu^{k+1}\}_k$ ,  $\{\gamma^{k+1}\}_k$ , and  $\{\theta^{k+1}\}_k$  of multipliers are bounded  
 543 as well;  
 544 (c) use (a) and (b) to take the limit  $k \rightarrow \infty$  in the above KKT system to obtain  
 545 the conclusion.

546 To prove (a), we first look at each  $g_i^k$  term from the expression of  $g^k$  in steps 7  
 547 and 9 of Algorithm 4.1 (see also (4.4)). We use a well-known proximal behavior,  
 548 namely, if  $\tilde{x}_k \rightarrow x^{\lambda^*}$  and  $w_i^k \in \text{prox}_{\lambda} r_i(\tilde{x}^k)$ , then the sequence  $\{w_i^k\}$  is bounded and  
 549 its cluster points  $w_i^{\lambda^*}$  lie in  $\text{prox}_{\lambda} r_i(x^{\lambda^*})$  [30, Proposition 1.25]. Hence the sequence  
 550 of  $g_i^k = \frac{1}{\lambda}(\tilde{x}^k - w_i^k)$  is bounded, and so is the sequence of their sum, which is  $g^k$ .  
 551 Similarly, cluster points satisfy the relationship from (i), with the observation that  
 552 the relationship  $g_i^{\lambda^*} \in \partial r_i(w_i^{\lambda^*})$  follows from (i) of Lemma 2.1.

553 We prove (b) by contradiction, namely assume there exists an unbounded subse-  
 554 quence (not relabeled) such that  $t_k := \|(\nu^k, \gamma^k, \theta^k)\| \rightarrow \infty$ . Since  $\|(\frac{\nu^k}{t_k}, \frac{\gamma^k}{t_k}, \frac{\theta^k}{t_k})\| = 1$   
 555 for all  $k$ , there must exist a cluster point and a subsequence (not relabeled) such that  
 556

$$557 \quad (4.14) \quad \left( \frac{\nu^k}{t_k}, \frac{\gamma^k}{t_k}, \frac{\theta^k}{t_k} \right) \rightarrow (\hat{\nu}, \hat{\gamma}, \hat{\theta}) \text{ as } k \rightarrow \infty \text{ with } \|(\hat{\nu}, \hat{\gamma}, \hat{\theta})\| = 1.$$

558 Next consider (4.11) rescaled by  $t_{k+1}$ , namely

$$559 \quad \frac{1}{t_{k+1}} (\nabla f(\tilde{x}^{k+1}) + g^k + 2\alpha_k \Delta^k) + \nabla h(\tilde{x}^{k+1})^\top \frac{\nu^{k+1}}{t_{k+1}} - \frac{\gamma^{k+1}}{t_{k+1}} + \frac{\theta^{k+1}}{t_{k+1}} = 0.$$

560 Recall that we work with a convergent (sub)sequence  $\{\tilde{x}^{k+1}\} \rightarrow x^{\lambda^*}$ , which makes the  
 561 term  $\nabla f(\tilde{x}^{k+1}) + g^k + 2\alpha_k \Delta^k$  bounded/convergent due to  $\nabla f$  being continuous (by  
 562 Assumption 2.4),  $\{g^k\}$  being bounded (by (i) above),  $\Delta^k \rightarrow 0$  (by Proposition 4.5),  
 563 and  $\alpha_k \leq \bar{\alpha}_s < \infty$  (by Lemma 4.4). Also  $\nabla h(\tilde{x}^{k+1}) \rightarrow \nabla h(x^{\lambda^*})$  since  $\nabla h$  is continuous  
 564 per Assumption 2.4. We see that by also using (4.14), we can take the limit to obtain

$$565 \quad \nabla h(x^{\lambda^*})^\top \hat{\nu} - \hat{\gamma} + \hat{\theta} = 0.$$

566 Similar continuity arguments can be applied to (4.13) to obtain that

$$567 \quad 0 \leq \hat{\gamma} \perp (x^{\lambda^*} - \ell) \geq 0, \quad 0 \leq \hat{\theta} \perp (u - x^{\lambda^*}) \geq 0.$$

568 By MFCQ, the only multipliers  $\hat{\nu}$ ,  $\hat{\gamma}$ ,  $\hat{\theta}$  satisfying the above relationships are the trivial  
 569 ones,  $(\hat{\nu}, \hat{\gamma}, \hat{\theta}) = (0, 0, 0)$ , which contradicts (4.14). Therefore,  $\{\nu^{k+1}\}_k$ ,  $\{\gamma^{k+1}\}_k$ , and  
 570  $\{\theta^{k+1}\}_k$  are bounded.

571 We now prove (c). First consider any cluster point  $(\nu^{\lambda^*}, \gamma^{\lambda^*}, \theta^{\lambda^*})$  of  $(\nu^k, \gamma^k, \theta^k)$   
572 and cluster point  $g^{\lambda^*}$  of  $g^k$ . Since  $\Delta_k \rightarrow 0$  by Proposition 4.5,  $\alpha_k \leq \bar{\alpha}_s \leq \infty$  by  
573 Lemma 4.4, and  $\nabla f$ ,  $\nabla h$ , and  $h$  are continuous by Assumption 2.4, we take the limit  
574 over any converging subsequence in (4.11)–(4.13) to obtain (4.9)–(4.10).

575 The Clarke stationarity follows from (4.9)–(4.10) since, as we next show,  $\sum_{i \in \mathcal{K}} g_i^{\lambda^*} \in$   
576  $\partial^C r_\lambda(x^{\lambda^*})$ . A general characterization [18][Proposition 3.1 and Theorem 3.2] for the  
577 Clarke subdifferentials of Moreau envelopes  $e_\lambda r_i(\cdot)$  is available via convex hulls in-  
578 volving proximal points, namely,

$$579 \quad \partial^C e_\lambda r_i(x) = \text{co} \left\{ \frac{1}{\lambda}(x - w) : w \in \text{prox}_\lambda r_i(x) \right\}.$$

580 Therefore,  $g_i^{\lambda^*}$  defined in (4.9)–(4.10) are Clarke subgradients of the Moreau envelopes,

$$581 \quad g_i^{\lambda^*} \in \partial^C e_{\lambda^*} f_i(x^{\lambda^*}), \forall i \in \mathcal{K}.$$

582 By Lemma A.1, we have that  $\sum_{i \in \mathcal{K}} g_i^{\lambda^*} \in \partial^C (\sum_{i \in \mathcal{K}} e_{\lambda^*} f_i(x^{\lambda^*})) = \partial^C r_{\lambda^*}(x^{\lambda^*})$ , which  
583 concludes the proof.  $\square$

584 *Remark 4.7.* Clarke stationarity conditions (ii) of Theorem 4.6 are weaker than  
585 the (limiting) sum-stationarity from (i) since the Clarke subdifferential is the convex  
586 hull of the limiting subdifferential in our general setup. This is a major improvement  
587 over our previous algorithm for upper- $C^2$  [33], mainly possible because we exploit  
588 the additional structure that Moreau envelopes have over general upper- $C^2$  functions,  
589 namely, the relationships between proximal points (minimizers of the Moreau enve-  
590 lope) and limiting subdifferentials from Lemma 2.1. In this respect, (4.9)–(4.10) of  
591 Theorem 4.6 can be viewed as stationarity conditions involving limiting subgradients  
592 at “nearby” points. More important, they are instrumental in proving the (sub-  
593 sequence) convergence of the outer iterates in the limit of Moreau parameter  $\lambda \downarrow 0$ .  
594 Clarke stationarity conditions (ii) are too weak to obtain the same convergence result.

595 **4.2. Outer-loop convergence.** In this section we elucidate the convergence  
596 behavior of the outer loop of Algorithm 4.1. We start with a technical lemma showing  
597 that proximal points of convergent sequences are also convergent.

598 **LEMMA 4.8.** *Consider a sequence  $x^s$  that converges to  $x^*$  and a sequence  $\lambda_s \downarrow 0$ .  
599 For any choice of proximal points  $w_i^s \in \text{prox}_{\lambda_s} r_i(x^s)$ , the sequences  $\{w_i^s\}$  converge to  
600  $x^*$  for all  $i \in \mathcal{K}$ .*

601 *Proof.* The proof is common to all  $i \in \mathcal{K}$ , therefore let  $r := r_i$  and  $w_s := w_i^s$ . Prox-  
602 boundedness of  $r$  implies that there exists a fixed  $\bar{\lambda} > 0$  such that  $\bar{m} := e_{\bar{\lambda}}(x^*) > -\infty$ .  
603 We then use the definition (1.3) of the Moreau envelope to write  $r(w^s) \geq \bar{m} - \frac{1}{2\bar{\lambda}} \|w^s -$   
604  $x^*\|^2$ . By the choice of  $w^s \in \text{prox}_{\lambda_s} r(x^s)$ , for  $\lambda_s < \bar{\lambda}/2$ , we can write

$$605 \quad (4.15) \quad r(w^s) + \frac{1}{2\lambda_s} \|w^s - x^s\|^2 \leq r(x^*) + \frac{1}{2\lambda_s} \|x^s - x^*\|^2.$$

606 The previous two inequalities imply that

$$607 \quad \bar{m} - \frac{1}{2\bar{\lambda}} \|w^s - x^*\|^2 + \frac{1}{2\lambda_s} \|w^s - x^s\|^2 \leq r(x^*) + \frac{1}{2\lambda_s} \|x^s - x^*\|^2.$$

608 Now use  $\|w^s - x^s\|^2 = \|(w^s - x^*) - (x^s - x^*)\|^2 \geq \frac{1}{2} \|w^s - x^*\|^2 - \|x^s - x^*\|^2$  and bring  
609 the terms involving  $\|x^s - x^*\|^2$  to the left to obtain

$$610 \quad \bar{m} + \left( \frac{1}{4\lambda_s} - \frac{1}{2\bar{\lambda}} \right) \|w^s - x^*\|^2 \leq r(x^*) + \frac{1}{\lambda_s} \|x^s - x^*\|^2.$$

611 Consequently,  $\|w^s - x^*\|^2 \leq \frac{4\bar{\lambda}}{\bar{\lambda} - 2\lambda_s} [\lambda_s(r(x^*) - \bar{m}) + \|x^s - x^*\|^2]$ , which proves  $w^s \rightarrow$   
612  $x^*$ .  $\square$

613 One would expect the function values  $r_i(w_i^s) \rightarrow r_i(x^*)$  as well; however, this is  
614 not necessarily the case when  $r_i(\cdot)$  has a jump discontinuity at  $x^*$ . For example,  
615  $r(x) = 1$  for  $x \neq 0$  and  $r(0) = 0$ , the Moreau envelope is  $e_\lambda(x) = \min\{1, \frac{|x|^2}{2\lambda}\}$ ; also,  
616  $\text{prox}_\lambda r(x) = \{0\}$  for  $|x| < \sqrt{2\lambda}$ ,  $\text{prox}_\lambda r(x) = \{0, x\}$  for  $|x| = \sqrt{2\lambda}$ , and  $\text{prox}_\lambda r(x) =$   
617  $\{x\}$  for  $|x| > \sqrt{2\lambda}$ . For simplicity, let us work with only one term,  $\mathcal{K} = \{1\}$ , and  
618 set  $r_i(x) = r(x)$ . Assume that Algorithm 4.1 generates  $x^s = \sqrt{2\lambda_s}$  and that the  
619 proximal oracle always returns  $w_i^s = \sqrt{2\lambda_s} \in \text{prox}_{\lambda_s} r(x) = \{0, \sqrt{2\lambda_s}\}$ . We observe  
620 that  $r_i(w_i^s) \rightarrow 1 \neq r_i(0)$ , even though  $w_i^s \rightarrow x^* = 0$ . This pathology has previously  
621 been recognized in [5], where an attentive convergence condition is used (see their  
622 assumption (H3)) to prove convergence of their PPA framework, namely the condition  
623 below.

624 *Condition 4.9.* For convergent sequences  $x^s \rightarrow x^*$  and corresponding  $w_i^s \in \text{prox}_{\lambda_s} r_i(x^s)$ ,  
625 attentive convergence holds along  $\{w_i^s\}$ , namely  $r_i(w_i^s) \rightarrow r_i(x^*)$ , for all  $i \in \mathcal{K}$ .  $\blacksquare$

626 This condition is quite general and is satisfied, for example, when the functions  
627  $r_i(\cdot)$  are continuous. It should also be noted that it holds in the lsc case whenever  
628 the outer iterates  $x^s$ , as returned by the inner loop of Algorithm 4.1, are exact min-  
629 imizers of the regularized problem (1.4)–(1.5) as  $\lambda_s \downarrow 0$ . The latter fact is proved in  
630 Lemma 4.15.

631 Given two tolerances  $\varepsilon > 0$  and  $\delta > 0$ , a point  $x \in C$  is  $(\varepsilon, \delta)$ -approximately sum-  
632 stationary if there exist points  $w_i$ , subgradients  $g_i \in \partial r_i(w_i)$ , and a normal vector  
633  $\eta \in N_C(x)$  such that

$$634 \quad \|w_i - x\|_\infty \leq \delta, \quad \forall i \in \mathcal{K}, \quad \text{and} \quad \|\nabla f(x) + \sum_{i \in \mathcal{K}} g_i + \eta\|_\infty \leq \varepsilon.$$

635 In words,  $x \in C$  is  $(\varepsilon, \delta)$ -approximately sum-stationary if it is stationary up to a  
636 residual norm  $\varepsilon$  with subgradients sampled at points no further than  $\delta$  from  $x$ .

637 Our main and general convergence results are provided by the following proposi-  
638 tion. In the following, let  $k_s$  be the number of inner iterations at each outer iteration  
639  $s$  in Algorithm 4.1. We recall that  $x^s = \tilde{x}^{k_s}$  according to our notation convention;  
640 also, let  $x_-^s = \tilde{x}^{k_s-1}$  and  $\Delta^s = x^s - x_-^s$ .

641 **PROPOSITION 4.10.** *Consider the iterates sequence  $\{x^s\}$  generated by Algorithm 4.1*  
642 *under a sequence of regularization parameters  $\lambda_s \downarrow 0$  and inner-loop tolerances  $\varepsilon_s \downarrow 0$ .*  
643 *Assume that Condition 4.9 holds along subsequences  $x^s \rightarrow x^*$ , that the conditions of*  
644 *Theorem 4.6 hold, that BCQ holds at cluster points  $x^*$ , and that Assumptions 2.4*  
645 *and 2.5 are satisfied. Then:*

646 (i) *Each outer-loop iterate  $x^s$  satisfies*

$$647 \quad \|\nabla f(x^s) + \sum_{i \in \mathcal{K}} g_i^s + \nabla h(x^s)^\top \nu^s - \gamma^s + \theta^s\|_\infty \leq \varepsilon_s,$$

648 *together with the KKT feasibility  $x^s \in C$  and KKT complementarity condi-*  
649 *tions. As a result  $x^s$  is  $(\varepsilon_s, \delta_s)$ -approximately sum-stationary, with*

$$650 \quad \delta_s := \max_{i \in \mathcal{K}} \|x^s - w_i^s\|_\infty \leq \frac{\varepsilon_s}{2\alpha} + \lambda_s \max_{i \in \mathcal{K}} \|g_i^s\|_\infty.$$

651 (ii) *The sequences  $\{g_i^s\}$  and  $\{\eta^s\}$ , where  $\eta^s := \nabla h(x^s)^\top \nu^s - \gamma^s + \theta^s$ , are bounded.*  
652 *Consequently,  $\delta_s \leq \frac{\varepsilon_s}{2\alpha} + M_g \lambda_s$  with  $M_g := \sup_s \max_{i \in \mathcal{K}} \|g_i^s\|_\infty < \infty$  and,*  
653 *therefore  $\delta_s \downarrow 0$ .*

654 (iii) For any prescribed tolerances  $\lambda_{tol} > 0$  and  $\varepsilon_{tol} > 0$ , Algorithm 4.1 with  
 655 stopping criteria  $\lambda_s \leq \lambda_{tol}$  and  $\varepsilon_s \leq \varepsilon_{tol}$  terminates after finitely many iter-  
 656 ations and returns a point  $x^s$  that is  $(\varepsilon_s, \delta_s)$ -approximately sum-stationary,  
 657 with  $\delta_s \leq \frac{\varepsilon_{tol}}{2\alpha} + M_g \lambda_{tol}$ .

658 *Proof.* (i) Fix an outer iteration  $s$ . At the terminating inner iteration index  $k_s$ ,  
 659 which is finite per Proposition 4.5, the KKT conditions of the model NLP at this  
 660 iteration (see (4.11)–(4.13)) give

$$661 \quad \nabla f(x^s) + \sum_{i \in \mathcal{K}} g_i^s + \eta^s = -2\alpha_{k_s-1}(x^s - x_-^s),$$

662 where  $\eta^s$  is the shorthand defined in (ii) above. We also remark that  $\hat{R}(x^s) = \nabla f(x^s) +$   
 663  $\sum_{i \in \mathcal{K}} g_i^s + \eta^s$  simply by the form of  $\eta^s$  and definition of  $\hat{R}(x^s)$  (see above (4.6)).  
 664 By the algorithm's inner loop stopping criterion (4.6) at iteration  $k_s$ , we have that  
 665  $\|\hat{R}(x^s)\|_\infty \leq \varepsilon_s$ , yielding the stationarity inequality from (i). The feasibility and  
 666 complementarity conditions follow directly from the aforementioned KKT conditions  
 667 of the model NLP.

668 The inner loop stopping criterion  $\|\hat{R}(x^s)\|_\infty \leq \varepsilon_s$  also implies

$$669 \quad \|x^s - x_-^s\|_\infty \leq \frac{\varepsilon_s}{2\alpha_{k_s-1}} \leq \frac{\varepsilon_s}{2\alpha},$$

670 where the last inequality is due to the lower bound for penalty parameters  $\alpha_k$  from  
 671 Lemma 4.4.

672 On the other hand, since  $x_-^s - w_i^s = \lambda_s g_i^s$  for all  $i \in \mathcal{K}$ , we can write

$$673 \quad \|x^s - w_i^s\|_\infty \leq \|x^s - x_-^s\|_\infty + \|x_-^s - w_i^s\|_\infty = \|x^s - x_-^s\|_\infty + \lambda_s \|g_i^s\|_\infty,$$

674 which proves the  $(\varepsilon_s, \delta_s)$ -approximate sum-stationarity of  $x^s$ .

675 (ii) Since  $\{x^s\} \in C$ , Assumption 2.5 ensures  $\{x^s\}$  is bounded. Let

$$676 \quad e^s = \nabla f(x^s) + \sum_{i \in \mathcal{K}} g_i^s + \eta^s.$$

677 Suppose by contradiction that the family  $\{g_i^s\}_s$ ,  $i \in \mathcal{K}$  and  $\{\eta^s\}_s$  is unbounded  
 678 as  $\lambda_s \rightarrow 0$ . Passing to an unbounded subsequence and relabeling if necessary, we may  
 679 assume  $\sum_{i \in \mathcal{K}} \|g_i^s\| + \|\eta^s\| \rightarrow +\infty$ . Define the scaling

$$680 \quad t_s = \frac{1}{1 + \sum_{i \in \mathcal{K}} \|g_i^s\| + \|\eta^s\|} \downarrow 0,$$

681 and multiply stationarity identity above by  $t_s$  to obtain

$$682 \quad t_s \nabla f(x^s) + \sum_{i \in \mathcal{K}} t_s g_i^s + t_s \eta^s = t_s e^s.$$

683 Since  $x^s$  is bounded and  $\nabla f(\cdot)$  is continuous,  $\{\nabla f(x^s)\}$  must be bounded, hence  
 684  $t_s \nabla f(x^s) \rightarrow 0$ . Similarly, since  $\|e^s\|_\infty = \|\hat{R}(x^s)\|_\infty \leq \varepsilon_s$ , we must have  $t_s e^s \rightarrow 0$ .  
 685 Hence,

$$686 \quad (4.16) \quad \sum_{i \in \mathcal{K}} t_s g_i^s + t_s \eta^s \rightarrow 0.$$

687 On the other hand, simply by the definition of  $t_s$  and the unboundedness hypoth-  
688 esis, we observe that

$$689 \quad \sum_{i \in \mathcal{K}} \|t_s g_i^s\| + \|t_s \eta^s\| = \frac{\sum_{i \in \mathcal{K}} \|g_i^s\| + \|\eta^s\|}{1 + \sum_{i \in \mathcal{K}} \|g_i^s\| + \|\eta^s\|} \rightarrow 1,$$

690 and, therefore,  $\{(t_s g_i^s)_{i \in \mathcal{K}}, t_s \eta^s\}$  has a nonzero cluster point. Passing to a subsequence  
691 if needed, we obtain

$$692 \quad t_s g_i^s \rightarrow v_i^\infty \quad \forall i \in \mathcal{K} \text{ and } t_s \eta^s \rightarrow \eta^\infty,$$

693 with at least one of the  $v_i^\infty$  and  $\eta^\infty$  being nonzero.

694 Consider a convergent subsequence, still not relabeled, such that  $x^s \rightarrow x^*$ . From  
695 part (i),  $\|x^s - x_-^s\| \leq \frac{\varepsilon_s}{2\alpha}$ , hence  $x_-^s \rightarrow x^*$  as well. Lemma 4.8 applied to  $\{x_-^s\}$   
696 implies  $w_i^s \rightarrow x^*$ ; also  $r_i(w_i^s) \rightarrow r_i(x^*)$  under the  $r_i$ -attentive convergence of  $\{w_i^s\}$   
697 from Condition 4.9.

698 We next prove  $v_i^\infty \in \partial^\infty r_i(x^*)$ . We first pass to regular subgradients required by  
699 the definition of horizon subdifferentials. Since  $g_i^s \in \partial r_i(w_i^s)$ , definition of the limiting  
700 subdifferential yields points  $\hat{w}_i^s$  and  $\hat{g}_i^s \in \hat{\partial} r_i(\hat{w}_i^s)$  for all  $s$  such that

$$701 \quad \|\hat{w}_i^s - w_i^s\| \leq t_s, \quad |r_i(\hat{w}_i^s) - r_i(w_i^s)| \leq t_s, \quad \text{and} \quad \|\hat{g}_i^s - g_i^s\| \leq t_s.$$

702 Then  $\hat{w}_i^s \rightarrow x^*$ ,  $r_i(\hat{w}_i^s) \rightarrow r_i(x^*)$ , and

$$703 \quad \|t_s \hat{g}_i^s - v_i^\infty\| \leq t_s \|\hat{g}_i^s - g_i^s\| + \|t_s g_i^s - v_i^\infty\| \rightarrow 0.$$

704 With these regular subgradients, the definition of the horizon subdifferential provides

$$705 \quad v_i^\infty \in \partial^\infty r_i(x^*), \quad \forall i \in \mathcal{K}.$$

706 Similarly, since each  $N_C(x^s)$  is a cone and  $\eta^s \in N_C(x^s)$ , we have  $t_s \eta^s \in N_C(x^s)$ .  
707 By outer semicontinuity of the limiting normal cone [30, Proposition 6.6],  $x^s \rightarrow x^*$   
708 and  $t_s \eta^s \rightarrow \eta^\infty$  imply  $\eta^\infty \in N_C(x^*)$ . Taking the limit  $t_s \downarrow 0$  in (4.16), we obtain  
709  $\sum_{i \in \mathcal{K}} v_i^\infty + \eta^\infty = 0$ . BCQ at  $x^*$  then forces  $v_i^\infty = 0$  for all  $i$  and  $\eta^\infty = 0$ , contradicting  
710 that at least one cluster is nonzero. Therefore  $\{g_i^s\}$  and  $\{\eta^s\}_s$  are bounded.

711 (iii) Since  $\lambda_s \downarrow 0$  and  $\varepsilon_s \downarrow 0$ , there exists a finite outer index  $\bar{s}$  with  $\lambda_{\bar{s}} \leq \lambda_{tol}$   
712 and  $\varepsilon_{\bar{s}} \leq \varepsilon_{tol}$ . By Proposition 4.5 for all outer iterations  $s \leq \bar{s}$  the inner iterations  
713 terminate finitely, therefore  $x^{\bar{s}}$  is reached in a finite number of iterations. The  $(\varepsilon_{\bar{s}}, \delta_{\bar{s}})$ -  
714 approximate stationarity follows from (i) and (ii) above.  $\square$

715 An explicit coupling between inner-loop tolerance  $\varepsilon_s$  and regularization  $\lambda_s$  from  
716 Proposition 4.10 should be used by a numerical method based on Algorithm 4.1 to  
717 properly balance accuracy of the inner and outer solves. The key bound is given by  
718 Proposition 4.10:

$$719 \quad \delta_s \leq \frac{\varepsilon_s}{2\alpha} + \lambda_s \max_{i \in \mathcal{K}} \|g_i^s\|_\infty = \frac{\varepsilon_s}{2\alpha} + \max_{i \in \mathcal{K}} \|x_-^s - w_i^s\|_\infty.$$

720 A natural choice is to set  $\varepsilon_s$  the same order as the proximal displacement  $B_s :=$   
721  $\max_{i \in \mathcal{K}} \|x_-^s - w_i^s\|_\infty$ . Otherwise, if  $\varepsilon_s \ll B_s$ , the regularized subproblem (1.4)–  
722 (1.5) is oversolved during the inner iterations; and, if  $B_s \ll \varepsilon_s$ , the overall error is  
723 dominated by the inner-loop accuracy. A practical, a posteriori rule for the inner-loop  
724 termination (step 20 in the algorithm) would therefore be

$$725 \quad \|\hat{R}(\tilde{x}^k)\|_\infty \leq B \max_{i \in \mathcal{K}} \|\tilde{x}^k - w_i^k\|_\infty \quad \text{or, equivalently,} \quad \|\hat{R}(\tilde{x}^k)\|_\infty \leq B \lambda_s \max_{i \in \mathcal{K}} \|g_i^k\|_\infty,$$

726 for some prescribed  $B > 0$ .

727 THEOREM 4.11. Consider Algorithm 4.1 under a sequence of regularizations  $\lambda_s \downarrow$   
728 0 and inner-loop tolerances  $\varepsilon_s \downarrow 0$ . Under the same assumptions as in Proposi-  
729 tion 4.10, any cluster point  $x^*$  of the bounded sequence  $\{x^s\}$  satisfies the sum-stationarity  
730 conditions (4.7)–(4.8) for the target problem (1.1)–(1.2).

731 *Proof.* The sequence  $\{x^s\}$  is bounded since  $C$  is compact by Assumption 2.5. We  
732 take an arbitrary accumulation point  $x^*$  and a convergent subsequence (not relabeled)  
733  $x^s \rightarrow x^*$ . By (ii) of Proposition 4.10, we can take subsequences, not relabeled, such  
734 that

$$735 \quad g_i^s \rightarrow g_i^*, \forall i \in \mathcal{K}, \text{ and } \eta^s \rightarrow \eta^*.$$

736 Possibly after a relabeling, for  $x^s \rightarrow x^*$  we have  $\nabla f(x^s) \rightarrow \nabla f(x^*)$  by continuity of  
737  $\nabla f$  given by Assumption 2.4. Also by closedness of normal cones  $\eta^* \in N_C(x^*)$ .

738 Using a similar argument as in the proof of (ii) of Proposition 4.10, one can prove  
739 that  $x_-^s \rightarrow x^*$  as well. Lemma 4.8 applied to the sequence  $\{x_-^s\}$  implies  $w_i^s \rightarrow x^*$ .  
740 Since  $g_i^s \in \partial r_i(w_i^s)$  and  $g_i^s \rightarrow g_i^*$ , under the  $r_i$ -attentive convergence of  $\{w_i^s\}$  from  
741 Condition 4.9, we must have  $g_i^* \in \partial r_i(x^*)$  for each  $i \in \mathcal{K}$ , by the outer semicontinuity  
742 of limiting subdifferential [30, Proposition 8.7].

743 With the above convergence properties, we can take the limit of the stationarity  
744 condition from (i) of Proposition 4.10 to obtain

$$745 \quad \nabla f(x^*) + \sum_{i \in \mathcal{K}} g_i^* + \eta^* = 0, \text{ with } g_i^* \in \partial r_i(x^*) \text{ and } \eta^* \in N_C(x^*).$$

746 Under the MFCQ for the feasible set  $C$  of the master problem (see Theorem 4.6),  
747  $\eta^* \in N_C(x^*)$  has the representation  $\eta^* = \nabla h(x^*)^\top \nu^* - \gamma^* + \theta^*$  with the multipliers  
748  $\nu^*$ ,  $\gamma^*$ , and  $\theta^*$  satisfying sign/complementarity conditions from (4.8). Finally, the  
749 stationarity equation above implies the stationarity condition (4.7), which completes  
750 the proof.  $\square$

751 A sharper iteration-complexity analysis via Kurdyka–Łojasiewicz (KL) assump-  
752 tion and arguments is not immediate in our framework, since our stationarity measure  
753 is based on proximal oracle subgradients, whereas the standard KL machinery is typi-  
754 cally formulated in terms of the limiting subdifferential and its associated slope.  
755 Relating these two notions in a way that yields clean rate complexity estimates would  
756 very likely require additional regularity assumptions, so we defer such a KL-based  
757 analysis to future work.

758 **4.3. Stronger convergence guarantees under additional regularity.** For  
759 several classes of objective functions  $r_i$  covered by Assumption 2.4, including lower-  
760  $C^1$ , lower- $C^2$ , and (lsc) convex functions, the sum-stationarity obtained above can be  
761 strengthened to the more familiar limiting subgradient stationarity. Moreover, the  
762 convergence prerequisites from the previous section hold automatically for broader  
763 classes of nonsmooth objectives, including locally Lipschitz functions (which encom-  
764 pass all of the aforementioned classes) and semismooth functions.

765 LEMMA 4.12. If the functions  $r_i$  satisfying Assumption 2.4 are also subdifferen-  
766 tially regular in the sense of Rockafellar and Wets [30, Definition 7.25] (also known  
767 as lower regular [24]), the sum-stationarity condition (4.7) from Theorem 4.11 takes  
768 the more commonly used form

$$769 \quad 0 \in \nabla f(x^*) + \partial \left( \sum_{i \in \mathcal{K}} r_i(x^*) \right) + \nabla h(x^*)^\top \nu^* - \gamma^* + \theta^*.$$

770 *Proof.* The stationarity condition follows from

$$771 \quad \sum_{i \in \mathcal{K}} \partial r_i(x^*) = \sum_{i \in \mathcal{K}} \hat{\partial} r_i(x^*) \subseteq \hat{\partial} \left( \sum_{i \in \mathcal{K}} r_i(x^*) \right) \subseteq \partial \left( \sum_{i \in \mathcal{K}} r_i(x^*) \right),$$

772 where the first inclusion is given by the sum rule for regular subgradients [30, Corollary  
773 10.9], the second inclusion is due to the regular subdifferential being a subset of the  
774 limiting subdifferential, and the equality is due to  $r_i$  being subdifferentially regular.  $\square$

775 We remark that if  $r_i$  are also locally Lipschitz, then the horizon qualification from  
776 Corollary 10.9 of [30] implies the stronger relationship  $\sum_{i \in \mathcal{K}} \partial r_i(x^*) = \partial \left( \sum_{i \in \mathcal{K}} r_i(x^*) \right)$ .  $\blacksquare$

777 Under the Lipschitzian property for objective terms  $r_i$ , the conditions needed for  
778 the finite termination and subsequence convergence of Proposition 4.10 and Theo-  
779 rem 4.11 can be considerably relaxed, as shown in the following proposition.

780 **PROPOSITION 4.13.** *If the functions  $r_i$  are locally Lipschitz, the finite convergence*  
781 *properties of Proposition 4.10 and the asymptotic convergence of Theorem 4.11 hold*  
782 *only under Assumptions 2.4 and 2.5 and MFCQ of Theorem 4.6 since Condition 4.9*  
783 *and BCQ hold automatically.*

784 *Proof.* The proof is immediate once we observe that (i)  $\partial^\infty r_i(\bar{x}) = \{0\}$  by local  
785 Lipschitz continuity [25], hence BCQ is satisfied, and (ii) Condition 4.9 also holds due  
786 to continuity of  $r_i$ .  $\square$

787 *Remark 4.14.* Proposition 4.13 applies to semismooth functions [22], lower- $C^1$ ,  
788 lower- $C^2$ , and (lsc) convex functions, which are locally Lipschitz.

789 The following lemma shows that Condition 4.9 holds for the general case ( $r_i$   
790 only satisfy Assumption 2.4 and are possibly discontinuous) if the inner iterations  
791 of Algorithm 4.1 finish with (approximate) minimizers  $x^s$  of the inner, fixed- $\lambda$  prob-  
792 lem (1.4)–(1.5) instead of (approximate) stationary points (in the sense of (i) of Propo-  
793 sition 4.10). This occurs, for example, when the feasible set  $C$  and the objective  $f$  in  
794 the master problem are convex and it may also occur as well when a global numerical  
795 solver is used for the inner problem.

796 **LEMMA 4.15.** *Consider Algorithm 4.1 under a sequence of regularization param-*  
797 *eters  $\lambda_s \downarrow 0$ . Suppose that, for each outer iteration  $s$ , the inner loop terminates after*  
798 *a finite number  $k_s$  of inner iterations and returns the outer iterate  $x^s := \tilde{x}^{k_s} \in C$ ,*  
799 *satisfying*

$$800 \quad F_{\lambda_s}(x^s) \leq \inf\{F_{\lambda_s}(x) : x \in C\} + \tau_s,$$

801 *for some sequence  $\tau_s \downarrow 0$ . Let  $x^s \rightarrow x^*$  be a convergent subsequence of such outer*  
802 *iterates. Then Condition 4.9 holds along this subsequence.*

803 *Proof.* Let  $w_i^s \in \text{prox}_{\lambda_s} r_i(x^s)$  and define

$$804 \quad t_i^s := \frac{1}{2\lambda_s} \|x^s - w_i^s\|^2 \geq 0.$$

805 From the choice of  $w_i^s$ , we have  $e_{\lambda_s} r_i(x^s) = r_i(w_i^s) + t_i^s$ , and therefore  $F_{\lambda_s}(x^s) =$   
806  $f(x^s) + \sum_{i \in \mathcal{K}} (r_i(w_i^s) + t_i^s)$ .

807 Since  $x^* \in C$  and  $x^s$  is  $\tau_s$ -suboptimal for  $F_{\lambda_s}$  over  $C$ ,

$$808 \quad F_{\lambda_s}(x^s) \leq \inf\{F_{\lambda_s}(x) : x \in C\} + \tau_s \leq F_{\lambda_s}(x^*) + \tau_s.$$

809 Also,  $e_{\lambda_s} r_i(x^*) \leq r_i(x^*)$  by the monotonicity of the Moreau envelope (see Theo-  
 810 rem 3.1), hence

$$811 \quad f(x^s) + \sum_{i \in K} (r_i(w_i^s) + t_i^s) \leq f(x^*) + \sum_{i \in K} r_i(x^*) + \tau_s.$$

812 Since  $f$  is continuous,  $x^s \rightarrow x^*$ , and  $\tau_s \rightarrow 0$ , we obtain

$$813 \quad (4.17) \quad \limsup_{s \rightarrow \infty} \sum_{i \in K} (r_i(w_i^s) + t_i^s) \leq \sum_{i \in K} r_i(x^*).$$

814 On the other hand, Lemma 4.8 gives  $w_i^s \rightarrow x^*$ , and lower semicontinuity of  $r_i$  implies

$$815 \quad r_i(x^*) \leq \liminf_{s \rightarrow \infty} r_i(w_i^s), \quad \forall i \in K.$$

816 Summing over  $i \in K$  and using  $t_i^s \geq 0$ , we get

$$817 \quad \sum_{i \in K} r_i(x^*) \leq \sum_{i \in K} \liminf_{s \rightarrow \infty} r_i(w_i^s) \leq \liminf_{s \rightarrow \infty} \sum_{i \in K} r_i(w_i^s) \leq \liminf_{s \rightarrow \infty} \sum_{i \in K} (r_i(w_i^s) + t_i^s),$$

818 where the second inequality uses properties of  $\liminf$  (i.e., Fatou's Lemma). Combin-  
 819 ing this with (4.17) yields

$$820 \quad \lim_{s \rightarrow \infty} \sum_{i \in K} (r_i(w_i^s) + t_i^s) = \sum_{i \in K} r_i(x^*).$$

821 Finally, for each  $i \in K$ , the lower semicontinuity of  $r_i$ , nonnegativity of  $t_i^s$ , and  
 822  $w_i^s \rightarrow x^*$  imply

$$823 \quad \liminf_{s \rightarrow \infty} (r_i(w_i^s) + t_i^s) \geq \liminf_{s \rightarrow \infty} r_i(w_i^s) \geq r_i(x^*).$$

824 The last two relations imply  $r_i(w_i^s) + t_i^s \rightarrow r_i(x^*)$ ,  $\forall i \in K$ . This and  $w_i^s \rightarrow x^*$  imply

$$825 \quad 0 \leq \limsup_{s \rightarrow \infty} t_i^s = \limsup_{s \rightarrow \infty} [r_i(w_i^s) + t_i^s - r_i(w_i^s)] \leq \limsup_{s \rightarrow \infty} [r_i(w_i^s) + t_i^s] - \liminf_{s \rightarrow \infty} r_i(w_i^s)$$

$$826 \quad = r_i(x^*) - \liminf_{s \rightarrow \infty} r_i(w_i^s) \leq 0,$$

827 therefore  $t_i^s \rightarrow 0$  and, consequently,  $r_i(w_i^s) \rightarrow r_i(x^*)$ , proving Condition 4.9.  $\square$

828 **5. A case study in parametric optimization with applications to elec-**  
 829 **trical power grids.** As mentioned in the introduction, a main motivation for the  
 830 present framework comes from optimization models in which the nonsmooth terms  $r_i$   
 831 are scenario-wise optimal value functions. A representative form is

$$832 \quad (5.1) \quad r_i(x) := \min_{y_1, y_2} c_i(y_1, y_2)$$

$$\text{s.t. } x - \Delta_i \leq y_1 \leq x + \Delta_i, \quad y_\ell \leq y_2 \leq y_u, \quad h_i(y_1, y_2) = 0,$$

833 where  $x$  denotes the master decision,  $y_1$  collects scenario variables constrained to re-  
 834 main close to  $x$ ,  $y_2$  collects scenario-local state variables,  $\Delta_i > 0$  prescribes admissible  
 835 post-scenario deviations, and  $i \in K$  indexes scenarios, e.g., system contingencies or  
 836 uncertainty realizations. We use this bound-coupling form only for notational sim-  
 837 plicity. The same discussion applies when  $x$  also enters the objective or the general  
 838 constraints of the scenario problem.

839 This is precisely the setting in which Algorithm 4.1 is useful. In many applications,  
 840 evaluating  $r_i(x)$  already requires solving a large NLP. The main difficulty is therefore  
 841 not the function evaluation itself, but obtaining reliable first-order information for the  
 842 value function  $r_i$ . Even when the scenario data are  $C^2$ , the marginal function  $r_i$  is,  
 843 in general, only lsc [30]. Under additional constraint qualifications such as MFCQ, it  
 844 becomes locally Lipschitz [13, 14]. However, even under MFCQ, both Clarke and lim-  
 845 iting subdifferentials may fail to admit multiplier-based representations in degenerate  
 846 or unstable nonconvex cases, see for example, Example 1 in [13] and Theorem 4 in [23],  
 847 respectively. Algorithm 4.1 avoids differentiating through the scenario NLP (5.1). In-  
 848 stead of relying on potentially numerically unavailable subgradients, it requires only  
 849 a proximal point of  $r_i$ , which can be obtained by solving essentially the same scenario  
 850 problem with an additional convex quadratic term. Concretely, one needs to solve

$$\begin{aligned}
 851 \quad (5.2) \quad e_\lambda r_i(x) &= \min_{y_1, y_2, w} c_i(y_1, y_2) + \frac{1}{2\lambda} \|w - x\|^2 \\
 &\text{s.t. } w - \Delta_i \leq y_1 \leq w + \Delta_i, \quad y_\ell \leq y_2 \leq y_u, \quad h_i(y_1, y_2) = 0,
 \end{aligned}$$

852 and any minimizer  $w_i^\lambda$  yields a proximal oracle subgradient  $g_i^\lambda = \frac{1}{\lambda}(x - w_i^\lambda)$ . Thus,  
 853 once a solver for the scenario NLP is available, the oracle required by Algorithm 4.1  
 854 can often be implemented with minimal structural changes, without differentiating  
 855 through KKT multipliers or relying on sensitivity routines.

856 This viewpoint is particularly natural for preventive security-constrained optimal  
 857 power flow models (SC-OPF) and related stochastic programming models relevant in  
 858 power dispatch operations under stochastic renewable generation, loads, and operat-  
 859 ing conditions [3]. The master variable  $x$  represents base-case power dispatch and  
 860 control settings. Each scenario  $i \in K$  corresponds to an  $N - 1$  contingency, such as  
 861 the loss of a generator, line, or transformer. For fixed  $x$ , the scenario subproblem (5.1)  
 862 computes a post-contingency operating point subject to power flow equations, nodal  
 863 active/reactive balance, generator capability limits, emergency voltage and thermal  
 864 limits, and contingency transition rules that describe the post-contingency response  
 865 and grid state via the constraints in (5.1). The scenario objective  $c_i$  is often con-  
 866 vex, as it captures quadratic penalties associated with the post-contingency response.  
 867 The NLP (5.1) is generally convex for direct-current (DC) OPF but nonconvex for  
 868 alternating-current (AC) OPF. For DC models, computing numerical proximal sub-  
 869 gradients with local optimizers applied to the proximal problem (5.2) is therefore  
 870 straightforward; for AC models, nonconvexity can be problematic. In our extensive  
 871 computational experience with scenario AC NLPs (5.1), local interior-point solvers  
 872 consistently return high-quality solutions that are stable across a range of initial-  
 873 izations. Since the proximal NLP (5.2) preserves the same scenario structure while  
 874 adding convex quadratic regularization, this evidence suggests that such solvers are  
 875 capable of recovering minimizers relevant to the Moreau envelope, and thereby provide  
 876 valid proximal oracle points and subgradients. A systematic numerical investigation  
 877 is deferred to future work.

878 Algorithm 4.1 matches the operational structure of grid models in several ways.  
 879 First, the master constraints (1.2) are enforced exactly, so every accepted iterate is a  
 880 feasible base-case operating point. This is especially attractive in real-time grid oper-  
 881 ations, which are subject to stringent time-to-solution requirements, where preserving  
 882 feasibility and security is often more important than aggressively pursuing marginal  
 883 cost decreases. Second, the proximal oracles are scenario-separable since for a given  $x$ ,  
 884 each  $e_\lambda r_i$  can be evaluated independently, resulting in a computational pattern with

885 exceptional parallel efficiency at large scale [20]. Third, the method returns first-order  
 886 information in a form that is immediately usable by the master algorithm, namely the  
 887 proximal oracle subgradients  $g_i^\lambda$ , without requiring the delicate post-processing of dual  
 888 information from the inner NLP solves. In [27] we successfully employed derivative-  
 889 free convex proximal stabilizations for  $r_i$ , constructed using power systems engineering  
 890 insight, to address the pervasive difficulties associated with derivative evaluations of  
 891  $r_i$ . This strategy, however, entails a loss of robustness and limits the method’s appli-  
 892 cability to other classes of optimization problems. In this sense, Algorithm 4.1 turns  
 893 the proximal stabilization intuition behind decomposition-based SC-OPF heuristics  
 894 from [27] into a principled continuation method with explicit stationarity guarantees.

895 The general regularity assumptions and results of Section 4 apply in this setting.  
 896 In addition to being lsc, the value functions  $r_i$  are proper and finite-valued whenever  
 897 the scenario problems are feasible for all relevant  $x$ ; furthermore, since the scenario  
 898 objective  $c_i$  from (5.1) is generally bounded from below,  $r_i$  is also prox-bounded and  
 899 therefore falls under Assumption 2.4. When MFCQ holds for the scenario NLP (5.1),  
 900  $r_i$  is locally Lipschitz [14], so the stronger Proposition 4.13 applies (since Condition 4.9  
 901 and BCQ hold automatically). Thus, the parametric NLP setting provides a concrete  
 902 class of problems compatible with the convergence guarantees presented above.

903 **6. Conclusions and future work.** We introduced a proximal oracle sequential  
 904 NLP framework for nonsmooth constrained optimization problems in which the non-  
 905 smooth objective is a finite sum of proper lower semicontinuous nonconvex terms ac-  
 906 cessed through proximal oracles. The method applies Moreau regularization termwise  
 907 to preserve rich decomposition properties of the original problem and builds exact-  
 908 constraint NLP models from the proximal points  $w_i \in \text{prox}_\lambda r_i(x)$  and the associated  
 909 proximal oracle subgradients  $(x - w_i)/\lambda$ . This construction avoids differentiating  
 910 through scenario subproblems and does not rely on multiplier sensitivity calcula-  
 911 tions. We proved consistency of the distributed Moreau-envelope continuation at the  
 912 level of global minimizers. For fixed  $\lambda$ , the inner sequential NLP iteration has clus-  
 913 ter points satisfying proximal limiting stationarity and Clarke stationarity for the  
 914 regularized problem. Along the outer continuation, cluster points satisfy first-order  
 915 sum-stationarity for the original problem under the stated attentive-convergence and  
 916 CQ assumptions; when the terms  $r_i$  are locally Lipschitz, these additional conver-  
 917 gence prerequisites hold automatically. The iterative algorithm terminates finitely  
 918 for prescribed positive regularization and stationarity tolerances with an explicit er-  
 919 ror bound. We also show how the required proximal oracles can be computed with  
 920 minimal modifications of the scenario NLP for the cases when  $r_i$  are optimal value  
 921 functions of parametric problems.

922 Several directions remain for future work. A primary direction is the development  
 923 of a robust implementation of the proposed algorithm, including treatment of inexact  
 924 proximal oracles, adaptive choices of regularization and inner-loop tolerances, and  
 925 safeguards for nonconvex NLP oracles when solved by local numerical solvers. Such an  
 926 implementation would enable systematic numerical evaluation on large-scale problems  
 927 and would clarify the tradeoffs among the various algorithmic and computational  
 928 choices that affect scalability, robustness, and applicability. Further work will also  
 929 investigate sharper complexity guarantees, for example under a KL assumption. This  
 930 would require relating the proximal oracle subgradient stationarity used here to a  
 931 limiting subdifferential slope [5], possibly under additional structure such as prox-  
 932 regularity [28] or a KL property [4, 5].

933 **Acknowledgments.** This work was performed under the auspices of the U.S.  
 934 Department of Energy by Lawrence Livermore National Laboratory under contract  
 935 DE-AC52-07NA27344.

936 **Appendix A. Sum rule for Clarke subdifferentials.** The following lemma  
 937 establishes that the Clarke subdifferential satisfies the exact sum rule for upper- $C^2$   
 938 functions. We include a proof because we were not able to find it established previously  
 939 in the literature.

940 LEMMA A.1. *Let  $g_1, g_2 : \mathbb{R}^n \rightarrow \mathbb{R}$  be upper- $C^2$  in a neighborhood of  $x$ . Then*

941 
$$\partial^C(g_1 + g_2)(x) = \partial^C g_1(x) + \partial^C g_2(x).$$

942 *Proof.* Since  $g_1$  and  $g_2$  are upper- $C^2$  near  $x$ , by the local representation theorem  
 943 for upper- $C^2$  functions [30, Theorem 10.33], there exist a neighborhood  $U$  of  $x$ ,  $C^2$   
 944 functions  $\phi_i : U \rightarrow \mathbb{R}$ , and convex functions  $q_i : U \rightarrow \mathbb{R}$  such that  $g_i = \phi_i - q_i$  on  $U$ ,  
 945  $i = \{1, 2\}$ . Therefore  $g_1 + g_2 = (\phi_1 + \phi_2) - (q_1 + q_2)$  on  $U$ , with  $\phi_1 + \phi_2$  being  $C^2$  and  
 946  $q_1 + q_2$  being convex.

947 For any proper convex function  $q$ , the Clarke subdifferential coincides with the  
 948 convex subdifferential [7], that is,  $\partial^C q(x) = \partial q(x)$ . Moreover, for a  $C^1$  function  $\phi$   
 949 and a locally Lipschitz function  $\psi$ , the Clarke sum rule [7] gives  $\partial^C(\phi + \psi)(x) =$   
 950  $\nabla\phi(x) + \partial^C\psi(x)$ . Applying this with  $\psi = -q$  and using  $\partial^C(-q)(x) = -\partial^C q(x)$ , we  
 951 obtain  $\partial^C(\phi - q)(x) = \nabla\phi(x) - \partial q(x)$ . Therefore,

952 
$$\partial^C(g_1 + g_2)(x) = \nabla(\phi_1 + \phi_2)(x) - \partial(q_1 + q_2)(x).$$

953 Since  $q_1$  and  $q_2$  are convex and finite on  $U$ , the convex subdifferential sum rule [30]  
 954 yields  $\partial(q_1 + q_2)(x) = \partial q_1(x) + \partial q_2(x)$ . Consequently,

955 
$$\begin{aligned} \partial^C(g_1 + g_2)(x) &= \nabla\phi_1(x) + \nabla\phi_2(x) - (\partial q_1(x) + \partial q_2(x)) \\ 956 &= (\nabla\phi_1(x) - \partial q_1(x)) + (\nabla\phi_2(x) - \partial q_2(x)) = \partial^C g_1(x) + \partial^C g_2(x). \end{aligned}$$

957 This proves the result. □

958 REFERENCES

959 [1] F. J. ARAGÓN ARTACHO AND J. M. BORWEIN, *Global convergence of a non-convex Douglas-*  
 960 *Rachford iteration*, Journal of Global Optimization, 57 (2013), pp. 753–769.  
 961 [2] F. J. ARAGÓN ARTACHO, J. M. BORWEIN, AND M. K. TAM, *Global behavior of the Douglas-*  
 962 *Rachford method for a nonconvex feasibility problem*, Journal of Global Optimization, 65  
 963 (2016), pp. 309–327.  
 964 [3] I. ARAVENA, D. K. MOLZAHN, S. ZHANG, C. G. PETRA, F. E. CURTIS, S. TU, A. WÄCHTER,  
 965 E. WEI, E. WONG, A. GHOLAMI, K. SUN, X. A. SUN, S. T. ELBERT, J. T. HOLZER,  
 966 AND A. VEERAMANY, *Recent developments in security-constrained AC optimal power flow:*  
 967 *Overview of Challenge 1 in the ARPA-E Grid Optimization Competition*, Operations Re-  
 968 search, 71 (2023), pp. 1997–2014.  
 969 [4] H. ATTOUCH AND J. BOLTE, *On the convergence of the proximal algorithm for nonsmooth*  
 970 *functions involving analytic features*, Mathematical Programming, 116 (2009), pp. 5–16.  
 971 [5] H. ATTOUCH, J. BOLTE, AND B. F. SVAITER, *Convergence of descent methods for semi-algebraic*  
 972 *and tame problems: proximal algorithms, forward-backward splitting, and regularized*  
 973 *Gauss–Seidel methods*, Mathematical Programming, 137 (2013), pp. 91–129.  
 974 [6] J. V. BURKE, F. E. CURTIS, A. S. LEWIS, M. L. OVERTON, AND L. E. A. SIMÕES, *Gradient*  
 975 *sampling methods for nonsmooth optimization*, in Numerical Nonsmooth Optimization:  
 976 State of the Art Algorithms, A. M. Bagirov, M. Gaudioso, N. Karimitsa, M. M. Mäkelä,  
 977 and S. Taheri, eds., Springer, Cham, 2020, pp. 201–225.

- 978 [7] Y. CUI AND J.-S. PANG, *Modern Nonconvex Nondifferentiable Optimization*, vol. 29 of MOS-  
979 SIAM Series on Optimization, Mathematical Optimization Society (MOS) and Society for  
980 Industrial and Applied Mathematics (SIAM), Philadelphia, PA, Dec 2021.
- 981 [8] F. E. CURTIS, T. MITCHELL, AND M. L. OVERTON, *A BFGS-SQP method for nonsmooth, non-*  
982 *convex, constrained optimization and its evaluation using relative minimization profiles*,  
983 *Optimiz. Meth. and Software*, 32 (2017), pp. 148–181.
- 984 [9] F. E. CURTIS AND M. L. OVERTON, *A sequential quadratic programming algorithm for noncon-*  
985 *vex, nonsmooth constrained optimization*, *SIAM J. on Optimization*, 22 (2012), pp. 474–  
986 500.
- 987 [10] F. E. CURTIS, X. QU, AND D. P. ROBINSON, *A proximal-gradient method for solving regularized*  
988 *optimization problems with general constraints*, 2025, <https://arxiv.org/abs/2512.23166>.  
989 Lehigh ISE Technical Report 25T-023; arXiv version revised 15 January 2026.
- 990 [11] D. DAVIS AND D. DRUSVYATSKIY, *Stochastic model-based minimization of weakly convex func-*  
991 *tions*, *SIAM J. on Optimization*, 29 (2019), pp. 207–239.
- 992 [12] A. DE MARCHI, X. JIA, C. KANZOW, AND P. MEHLITZ, *Constrained composite optimization*  
993 *and augmented Lagrangian methods*, *Mathematical Programming*, (2023).
- 994 [13] J. GAUVIN, *The generalized gradient of a marginal function in mathematical programming*,  
995 *Mathematics of Operations Research*, 4 (1979), pp. 458–463.
- 996 [14] J. GAUVIN AND F. DUBEAU, *Differential properties of the marginal function in mathemati-*  
997 *cal programming*, in *Optimality and Stability in Mathematical Programming*, vol. 19 of  
998 *Mathematical Programming Studies*, Springer, 1982, pp. 101–119.
- 999 [15] K. GUO, D. HAN, D. Z. W. WANG, AND T. WU, *Convergence of ADMM for multi-block*  
1000 *nonconvex separable optimization models*, *Frontiers of Mathematics in China*, 12 (2017),  
1001 pp. 1139–1162.
- 1002 [16] R. HESSE AND D. R. LUKE, *Nonconvex notions of regularity and convergence of fundamental*  
1003 *algorithms for feasibility problems*, *SIAM J. on Optimization*, 23 (2013), pp. 2397–2419.
- 1004 [17] M. HONG, Z.-Q. LUO, AND M. RAZAVIYAYN, *Convergence analysis of alternating direction*  
1005 *method of multipliers for a family of nonconvex problems*, *SIAM J. on Optimization*, 26  
1006 (2016), pp. 337–364.
- 1007 [18] A. JOURANI, L. THIBAUT, AND D. ZAGRODNY, *Differential properties of the Moreau envelope*,  
1008 *Journal of Functional Analysis*, 266 (2014), pp. 1185–1237.
- 1009 [19] P. D. KHANH, V. V. H. KHOA, B. S. MORDUKHOVICH, AND V. T. PHAT, *Local minimizers of*  
1010 *nonconvex functions in Banach spaces via Moreau envelopes*, *Vietnam Journal of Mathe-*  
1011 *matics*, 53 (2025), pp. 803–813.
- 1012 [20] LAWRENCE LIVERMORE NATL. LAB., *LLNL-Developed Software Leverages Exascale Power to*  
1013 *Optimize National Grid’s Emergency Response*. HPCwire, Off the Wire Press Releases,  
1014 Aug. 2023, [https://www.hpcwire.com/off-the-wire/llnl-developed-software-leverages-exa-](https://www.hpcwire.com/off-the-wire/llnl-developed-software-leverages-exa-scale-power-to-optimize-national-grids-emergency-response/)  
1015 [scale-power-to-optimize-national-grids-emergency-response/](https://www.hpcwire.com/off-the-wire/llnl-developed-software-leverages-exa-scale-power-to-optimize-national-grids-emergency-response/). Accessed: 2026-05-08.
- 1016 [21] G. LI AND T. K. PONG, *Douglas–Rachford splitting for nonconvex optimization with application*  
1017 *to nonconvex feasibility problems*, *Mathematical Programming*, 159 (2016), pp. 371–401.  
1018 Series A.
- 1019 [22] R. MIFFLIN, *Semismooth and semiconvex functions in constrained optimization*, *SIAM J. on*  
1020 *Control and Optimization*, 15 (1977), pp. 959–972.
- 1021 [23] B. S. MORDUKHOVICH, N. M. NAM, AND N. D. YEN, *Subgradients of marginal functions in*  
1022 *parametric mathematical programming*, *Mathematical Programming*, 116 (2009), pp. 369–  
1023 396.
- 1024 [24] B. S. MORDUKHOVICH AND T. T. A. NGHIA, *Subdifferentials of nonconvex supremum functions*  
1025 *and their applications to semi-infinite and infinite programs with lipschitzian data*, *SIAM*  
1026 *Journal on Optimization*, 23 (2013), pp. 406–431.
- 1027 [25] B. S. MORDUKHOVICH AND R. T. ROCKAFELLAR, *Second-order subdifferential calculus with*  
1028 *applications to tilt stability in optimization*, *SIAM Journal on Optimization*, 22 (2012),  
1029 pp. 953–986.
- 1030 [26] N. PARIKH AND S. BOYD, *Proximal algorithms*, *Foundations and Trends® in Optimization*, 1  
1031 (2014), pp. 127–239.
- 1032 [27] C. G. PETRA AND I. ARAVENA, *A surrogate-based asynchronous decomposition technique for*  
1033 *realistic security-constrained optimal power flow problems*, *Operations Research*, 71 (2023),  
1034 pp. 2015–2030.
- 1035 [28] R. A. POLIQUIN AND R. T. ROCKAFELLAR, *Prox-regular functions in variational analysis*, *Trans-*  
1036 *actions of the American Mathematical Society*, 348 (1996), pp. 1805–1838.
- 1037 [29] Z. QI, Y. CUI, Y. LIU, AND J. PANG, *Asymptotic properties of stationary solutions of nonconvex*  
1038 *nonsmooth empirical risk minimization*, *Math. Oper. Res.*, 47 (2022), pp. 2034–2064.
- 1039 [30] R. T. ROCKAFELLAR AND R. J.-B. WETS, *Variational Analysis*, Springer-Verlag, Berlin Hei-

- 1040 delberg, 1998.
- 1041 [31] A. THEMELIS AND P. PATRINOS, *Douglas–Rachford splitting and ADMM for nonconvex opti-*  
1042 *mization: Tight convergence results*, SIAM J. on Optimization, 30 (2020), pp. 149–181.
- 1043 [32] A. THEMELIS, L. STELLA, AND P. PATRINOS, *Douglas–Rachford splitting and ADMM for non-*  
1044 *convex optimization: accelerated and Newton-type linesearch algorithms*, Computational  
1045 Optimization and Applications, 82 (2022), pp. 395–440.
- 1046 [33] J. WANG AND C. G. PETRA, *A sequential quadratic programming algorithm for nonsmooth*  
1047 *problems with upper- $C^2$  objective*, SIAM J. on Optimization, 33 (2023), pp. 2379–2405.
- 1048 [34] Y. WANG, W. YIN, AND J. ZENG, *Global convergence of ADMM in nonconvex nonsmooth*  
1049 *optimization*, Journal of Scientific Computing, 78 (2019), pp. 29–63.
- 1050 [35] G. D. YALCIN AND F. E. CURTIS, *Incremental quasi-Newton algorithms for solving a nonconvex,*  
1051 *nonsmooth, finite-sum optimization*, Optimiz. Meth. and Software, 39 (2024), pp. 345–367.