

**Titre:** A proximal modified quasi-newton method for nonsmooth  
Title: regularized optimization

**Auteurs:** Youssef Diouane, Mohamed L. Habiboullah, & Dominique Orban  
Authors:

**Date:** 2026

**Type:** Article de revue / Article

**Référence:** Diouane, Y., Habiboullah, M. L., & Orban, D. (2026). A proximal modified quasi-  
Citation: newton method for nonsmooth regularized optimization. SIAM Journal on  
Optimization, 36(2), 534-563. <https://doi.org/10.1137/24m169761x>

 **Document en libre accès dans PolyPublie**  
Open Access document in PolyPublie

**URL de PolyPublie:** <https://publications.polymtl.ca/75896/>  
PolyPublie URL:

**Version:** Version finale avant publication / Accepted version  
Révisé par les pairs / Refereed

**Conditions d'utilisation:** Creative Commons Attribution 4.0 International (CC BY)  
Terms of Use:

 **Document publié chez l'éditeur officiel**  
Document issued by the official publisher

**Titre de la revue:** SIAM Journal on Optimization (vol. 36, no. 2)  
Journal Title:

**Maison d'édition:** Society for Industrial and Applied Mathematics  
Publisher:

**URL officiel:** <https://doi.org/10.1137/24m169761x>  
Official URL:

**Mention légale:**  
Legal notice:

# A PROXIMAL MODIFIED QUASI-NEWTON METHOD FOR NONSMOOTH REGULARIZED OPTIMIZATION

YOUSSEF DIOUANE\*, MOHAMED L. HABIBOULLAH†, AND DOMINIQUE ORBAN‡

**Abstract.** We develop R2N, a modified quasi-Newton method for minimizing the sum of a  $\mathcal{C}^1$  function  $f$  and a lower semi-continuous prox-bounded  $h$ . Both  $f$  and  $h$  may be nonconvex. At each iteration, our method computes a step by minimizing the sum of a quadratic model of  $f$ , a model of  $h$ , and an adaptive quadratic regularization term. A step may be computed by way of a variant of the proximal-gradient method. An advantage of R2N over competing trust-region methods is that proximal operators do not involve an extra trust-region indicator. We also develop the variant R2DH, in which the model Hessian is diagonal, which allows us to compute a step without relying on a subproblem solver when  $h$  is separable. R2DH can be used as standalone solver, but also as subproblem solver inside R2N. We describe non-monotone variants of both R2N and R2DH. Global convergence of a first-order stationarity measure to zero holds without relying on local Lipschitz continuity of  $\nabla f$ , while allowing model Hessians to grow unbounded, an assumption particularly relevant to quasi-Newton models. Under Lipschitz-continuity of  $\nabla f$ , we establish a tight worst-case evaluation complexity bound of  $O(1/\epsilon^{2/(1-p)})$  to bring said measure below  $\epsilon > 0$ , where  $0 \leq p < 1$  controls the growth of model Hessians. Specifically, the latter must not diverge faster than  $|S_k|^p$ , where  $S_k$  is the set of successful iterations up to iteration  $k$ . When  $p = 1$ , we establish the tight exponential complexity bound  $O(\exp(c\epsilon^{-2}))$  where  $c > 0$  is a constant. We describe our Julia implementation and report numerical experience on a classic basis-pursuit problem, an image denoising problem, a minimum-rank matrix completion problem, a nonlinear support vector machine and an inverse nonlinear problem.

**Key words.** Nonsmooth optimization; Nonconvex optimization; Regularized optimization; Composite optimization; Modified quasi-Newton method; Proximal quasi-Newton method; Proximal gradient method

**AMS subject classifications.** 90C30, 90C53,

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

$$(1.1) \quad \underset{x \in \mathbb{R}^n}{\text{minimize}} \quad f(x) + h(x),$$

where  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is  $\mathcal{C}^1$  on  $\mathbb{R}^n$ , and  $h : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$  is lower semi-continuous (lsc). Both  $f$  and  $h$  may be nonconvex. Our motivation is to develop a modified Newton variant of the trust-region algorithms of Aravkin et al. [3] and Leconte and Orban [25] because the proximal operators used in the subproblems should be easier to derive. For instance, when  $h$  is the rank function, the proximal operator with a trust-region indicator is not known analytically at this time.

We introduce method R2N, at each iteration of which the sum of a quadratic model of  $f$ , a model of  $h$ , and an adaptive quadratic regularization term, is approximately minimized. Both models may be nonconvex. The Hessian of the quadratic model of  $f$  may be that of  $f$  if it exists, or an approximation such as those derived from quasi-Newton updates. We establish global convergence of R2N under the assumption that the models of  $h$  are prox-bounded and approximate  $h(x + s)$  as  $o(\|s\|)$ —an assumption

---

\*GERAD and Department of Mathematics and Industrial Engineering, Polytechnique Montréal. E-mail: [youssef.diouane@polymtl.ca](mailto:youssef.diouane@polymtl.ca). Research partially supported by an NSERC Discovery Grant.

†GERAD and Department of Mathematics and Industrial Engineering, Polytechnique Montréal. E-mail: [mohamed-laghdaf-2.habiboullah@polymtl.ca](mailto:mohamed-laghdaf-2.habiboullah@polymtl.ca). Research partially supported by an FRQNT scholarship.

‡GERAD and Department of Mathematics and Industrial Engineering, Polytechnique Montréal. E-mail: [dominique.orban@gerad.ca](mailto:dominique.orban@gerad.ca). Research partially supported by an NSERC Discovery Grant.

that covers composite terms with Hölder Jacobian, see [Model Assumption 5.1](#) for details. No assumption on local Lipschitz continuity of  $\nabla f$  is required, nor is boundedness of the model Hessians, provided they do not diverge too fast. Specifically, if  $B_k$  is the model Hessian at iteration  $k$ , we require that the series with general term  $1/(1 + \max_{0 \leq j \leq k} \|B_j\|)$  diverge—an assumption similar to that used in trust-region methods [[14](#), §8.4]. Our assumptions are significantly weaker than assumptions commonly found in the analysis of competing methods, and, consequently, the applicability of R2N is significantly more general—see the related research section below for details.

R2N specializes to method R2DH when  $B_k$  is diagonal, as did the solver of [Leconte and Orban](#) [[25](#)]. For a number of choices of separable  $h$  that are relevant in applications, steps can be computed explicitly without resort to an iterative subproblem solver. R2DH can be used as standalone solver or as subproblem solver inside R2N.

We also develop complexity results inspired from those of [Leconte and Orban](#) [[24](#)] and [Diouane et al.](#) [[18](#)], that account for potentially unbounded model Hessians. Specifically, we require that either  $\|B_k\| = O(|\mathcal{S}_k|^p)$  for some  $0 \leq p \leq 1$ , where  $\mathcal{S}_k$  is the set of successful iterations up to iteration  $k$ . When  $0 \leq p < 1$ , we establish a tight  $O(\epsilon^{-2/(1-p)})$  complexity, and when  $p = 1$ , we establish a tight exponential complexity, i.e., a bound in  $O(\exp(c\epsilon^{-2}))$  where  $c > 0$  is a constant. Though the latter bound is tight, it is not known if it is attained for a quasi-Newton update.

We provide efficient implementations of R2N and R2DH. The latter can use one of several diagonal quasi-Newton updates. Both have non-monotone variants that preserve their convergence and complexity properties. Our open-source Julia implementations are available from [[5](#)]. In [Section 8](#), we illustrate the performance of R2N and R2DH on challenging problems, including minimum-rank problems for which the trust-region methods of [[3](#), [4](#)] are impractical.

**Contributions and related research.** The proximal-gradient method [[20](#), [29](#)] is the prototypical first-order method for [\(1.1\)](#). Vast amounts of literature consider variants but restrict  $f$  and/or  $h$  to be convex, impose that  $f$  have (locally) Lipschitz-continuous gradient, or that  $h$  be Lipschitz continuous. For instance, [[27](#)] develop a proximal Newton method that requires  $f$  and  $h$  convex, a positive semi-definite Hessian, and solve the subproblem via the proximal-gradient method. [Cartis et al.](#) [[11](#)] require  $h$  to be globally Lipschitz continuous. [Kanzow and Lechner](#) [[22](#)], [Liu et al.](#) [[30](#)] develop an approach closely related to ours, but for convex  $h$ . Others dispense with convexity but require coercivity of  $f + h$  [[28](#)].

We are aware of several references that allow both  $f$  and  $h$  to be nonconvex. For instance, [Bolte et al.](#) [[9](#)] propose PALM, an alternating first-order method for problems with partitioned variables. They assume that  $f$ , which acts on both sets of variables  $x$  and  $y$ , has a gradient that is Lipschitz continuous with respect to  $x$  and  $y$  separately, and is Lipschitz continuous in  $(x, y)$  on a bounded set. They also assume that the sequence generated by their method is bounded, which is a strong requirement. Under those assumptions, they prove that every accumulation point of the iterates is a first-order stationary point. Moreover, if the KL property holds, then the entire sequence converges to a first-order stationary point. No second-order information is used in their method.

[Boj et al.](#) [[10](#)] study a proximal-gradient algorithm with momentum for solving [\(1.1\)](#). They assume that  $\nabla f$  is Lipschitz continuous, that  $h$  is bounded below, and that  $f + h$  is coercive. Under those conditions, they show that every accumulation point of the iterates is a first-order stationary point of [\(1.1\)](#). If, in addition, the function  $H(x, y) = f(x) + h(x) + M\|x - y\|^2$ , where  $M > 0$ , is a KL function, then

the entire sequence converges to a first-order stationary point. Their method does not use second-order information, and the coercivity assumption can be restrictive.

[Themelis et al. \[36\]](#) propose ZeroFPR, a non-monotone line-search proximal quasi-Newton method based on the forward-backward envelope. They assume that  $\nabla f$  is globally Lipschitz continuous—although they note that local Lipschitz continuity suffices if the domain of  $h$  is bounded and the search directions remain bounded—and that  $h$  is prox-bounded. Under these assumptions, any accumulation point is a stationary point of (1.1). Furthermore, if the iterates remain bounded, the forward-backward envelope satisfies the KL property, and  $f$  is twice continuously differentiable, then the entire sequence converges to a first-order stationary point. With a suitable desingularization function, they also obtain R-linear convergence of  $\{x_k\}_{k \in \mathbb{N}}$ . Although ZeroFPR allows quasi-Newton approximations, the model Hessians are assumed to be uniformly bounded—a condition that may be difficult to guarantee, as we will discuss later. It furthermore requires an estimate of the Lipschitz constant of  $\nabla f$ , and a preliminary loop to compute it.

[Stella et al. \[35\]](#) propose PANOC, a line-search limited-memory BFGS method, in the context of optimal control problems. They assume that  $\nabla f$  is Lipschitz continuous, although local Lipschitz continuity suffices if the domain of  $h$  is bounded and the search directions remain bounded using the similar arguments as in [36]. However, if the Lipschitz constant of  $\nabla f$  is unknown, it must be estimated in a preliminary loop. In addition, they require  $h$  to be bounded below, which is a strong assumption. Under these conditions, any accumulation point of the iterates is stationary for (1.1). Moreover, if the iterates converge to a strong local minimum of  $f + h$ , the forward-backward envelope is twice continuously differentiable, the proximal operator of  $h$  is strictly differentiable, and the model Hessians  $B_k$  satisfy the Denis–Moré condition, then the entire sequence converges at a superlinear rate.

One may observe that in the methods of [9, 10, 36], the KL property together with the Lipschitz continuity of  $\nabla f$  is used to establish convergence of the entire sequence of iterates to a first-order stationary point. Moreover, with a suitable desingularization function, R-linear convergence of the iterates is shown in [10, 36]. By contrast, our work does not rely on such strong assumptions. We avoid both the KL property and the Lipschitz continuity of  $\nabla f$ , and aim instead for generality.

More recently, some works have derived convergence analyses without assuming Lipschitz continuity of  $\nabla f$  or the KL property. For example, [Kanzow and Mehlitz \[23\]](#) analyze monotone and non-monotone first-order proximal-gradient methods under the assumptions that  $f$  is continuously differentiable and  $h$  is lower semicontinuous and bounded below by an affine function. [De Marchi \[15\]](#) extends their setting to allow  $h$  to be prox-bounded. They establish that for any convergent subsequence  $\{x_k\}_{k \in \mathcal{K}} \rightarrow x^*$ , a stationarity measure converges to zero along  $\mathcal{K}$ . Additional strong assumptions—such as local Lipschitz continuity of  $\nabla f$  or continuity of  $h$ —then ensure that each accumulation point is a first-order stationary point. Theirs are first-order methods and do not incorporate second-order information. Furthermore, their convergence results do not cover cases where the iterates may be unbounded.

In this work, we establish the global convergence of a first-order stationarity measure to zero under the minimal assumptions that  $f$  is continuously differentiable and  $h$  is lower semicontinuous and prox-bounded. We further guarantee that, for any  $\epsilon > 0$ , our stationarity measure falls below  $\epsilon$  in a finite number of iterations, even if the sequence of iterates is unbounded.

Our work follows the scheme laid out by [Aravkin et al. \[3\]](#); a trust-region framework applicable to nonconvex  $f$  and/or  $h$ , and that does not require coercivity or KL

assumptions. However, their analysis relies on the Lipschitz continuity of  $\nabla f$  in a neighborhood of the iterates. They also describe a method named R2 that amounts to a proximal-gradient method with adaptive step size, and that may be viewed as R2N where  $B_k$  is set to zero at each iteration, effectively reducing to a first-order method. Aravkin et al. [4] specialize their trust-region method to problems where  $f$  has a least-squares structure, and develop a Levenberg-Marquardt variant named LM that may also be viewed as a special case of R2N for least-squares  $f$ . If  $J_k$  is the least-squares residual's Jacobian at  $x_k$ , their model of  $f$  uses  $B_k = J_k^T J_k$ . Leconte and Orban [25] devise variants of the trust-region method of [3] for separable  $h$  in which the model Hessian is a diagonal quasi-Newton approximation. They also devise non-monotone schemes that are shown to significantly improve performance in certain cases. All of [3, 4, 25] assume uniformly bounded second-order information in the model of  $f$ .

Leconte and Orban [24] revisit the trust-region method of [3] but allow for unbounded model Hessians. They establish global convergence and a worst-case complexity bound of  $O(\epsilon^{-2/(1-p)})$  provided  $\|B_k\| = O(|\mathcal{S}_k|^p)$  with  $0 \leq p < 1$ . To the best of our knowledge, they were the first to use that assumption and to obtain a complexity bound in the presence of unbounded model Hessians. Unfortunately, their analysis does not generalize to  $p = 1$ . Under additional assumptions detailed in Section 5, they also prove the existence of a subsequence of iterates that converges to a first-order stationary point of (1.1). This result allows us to derive a similar convergence guarantee for both R2N and R2DH; nevertheless, their analysis still requires Lipschitz continuity of  $\nabla f$  in a neighborhood of the iterates.

Potentially unbounded model Hessians are a relevant assumption in several contexts, including quasi-Newton methods. Conn et al. [14, §8.4.1.2] show that the SR1 approximation satisfies  $\|B_k\| = O(|\mathcal{S}_k|)$ , and a similar bound for BFGS when  $f$  is convex. Powell [32] establishes a similar bound for his PSB update. Even though it is not currently known whether those bounds are tight, the case  $p = 1$  covers them.

Diouane et al. [18] generalized the results of [24] to  $p = 1$  and provided tighter complexity constants when  $0 \leq p < 1$  in the context of trust-region methods for smooth optimization, i.e.,  $h = 0$ . Our complexity analysis draws from [18, 24].

**Notation.** Unless otherwise noted, if  $x$  is a vector,  $\|x\|$  denotes its Euclidean norm and if  $A$  is a matrix,  $\|A\|$  denotes its spectral norm. For positive sequences  $\{a_k\}$  and  $\{b_k\}$ , we say that  $a_k = o(b_k)$  if and only if  $\limsup_k a_k/b_k = 0$ . The cardinality of a finite set  $\mathcal{A}$  is denoted  $|\mathcal{A}|$ . We denote  $\mathbb{N}_0$  the set of positive integers.

**2. Background.** We recall relevant concepts of variational analysis, e.g., [33].

The domain of  $h$  is  $\text{dom } h := \{x \in \mathbb{R}^n \mid h(x) < \infty\}$ . Because  $h$  is proper,  $\text{dom } h \neq \emptyset$ . If  $P : \mathbb{R}^n \rightarrow \mathbb{R}^n$  is a set-valued function,  $\text{dom } P = \{x \in \mathbb{R}^n \mid P(x) \neq \emptyset\}$ .

**DEFINITION 2.1.** (*Limiting subdifferential*) Consider  $\phi : \mathbb{R}^n \rightarrow \bar{\mathbb{R}}$  and  $\bar{x} \in \mathbb{R}^n$  such that  $\phi(\bar{x}) < +\infty$ . We say that  $v \in \mathbb{R}^n$  is a regular subgradient of  $\phi$  at  $\bar{x}$  if

$$\liminf_{x \rightarrow \bar{x}} \frac{\phi(x) - \phi(\bar{x}) - v^T(x - \bar{x})}{\|x - \bar{x}\|} \geq 0.$$

The set  $\hat{\partial}\phi(\bar{x})$  of all regular subgradients of  $\phi$  at  $\bar{x}$  is called the Fréchet subdifferential.

The limiting subdifferential of  $\phi$  at  $\bar{x}$  is the set  $\partial\phi(\bar{x})$  of all  $v \in \mathbb{R}^n$  such that there is  $\{x_k\} \rightarrow \bar{x}$  with  $\{\phi(x_k)\} \rightarrow \phi(\bar{x})$  and  $\{v_k\} \rightarrow v$  with  $v_k \in \hat{\partial}\phi(x_k)$  for all  $k$ .

If  $\phi = f + h$  with  $f$  continuously differentiable and  $h$  lower semi-continuous, then  $\partial\phi(x) = \nabla f(x) + \partial h(x)$  [33, Theorem 10.1].

DEFINITION 2.2. (*Proximal Operator*) Let  $h : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$  be proper and lower semi-continuous. The proximal operator of  $h$  with step length  $\nu > 0$  is

$$\text{prox}(x) := \underset{y}{\text{argmin}}_{\nu h} h(y) + \frac{1}{2}\nu^{-1}\|y - x\|^2.$$

Without further assumptions on  $h$ , the proximal operator might be empty, or contain one or more elements.

By [33, Exercise 8.8c],  $\bar{x}$  is first-order stationary for (1.1) if  $0 \in \nabla f(\bar{x}) + \partial h(\bar{x})$ .

DEFINITION 2.3. (*Outer limit*) Let  $\mathcal{K} \subseteq \mathbb{N}$  and a sequence  $\{S_k\}_{k \in \mathcal{K}}$  of subsets of  $\mathbb{R}^n$ . The outer limit of  $\{S_k\}_{k \in \mathcal{K}}$  is

$$\limsup_{k \in \mathcal{K}} S_k := \{v \in \mathbb{R}^n \mid \exists \mathcal{K}' \subseteq \mathcal{K}, \exists v_k \in S_k \text{ for all } k \in \mathcal{K}' \text{ with } \lim_{k \in \mathcal{K}'} v_k = v\}.$$

Let  $\mathcal{K} \subseteq \mathbb{N}$  and  $\{x_k\}_{k \in \mathcal{K}}$  such that  $\lim_{k \in \mathcal{K}} x_k = \bar{x}$  and  $\lim_{k \in \mathcal{K}} h(x_k) = h(\bar{x})$ . By [33, Proposition 8.7],  $\limsup_{k \in \mathcal{K}} \partial h(x_k) \subseteq \partial h(\bar{x})$ .

**3. Models.** For  $\sigma \geq 0$ ,  $x \in \mathbb{R}^n$ , and  $B(x) = B(x)^T \in \mathbb{R}^{n \times n}$ , consider the models

$$(3.1a) \quad \varphi(s; x) := f(x) + \nabla f(x)^T s + \frac{1}{2} s^T B(x) s$$

$$(3.1b) \quad \psi(s; x) \approx h(x + s)$$

$$(3.1c) \quad m(s; x, \sigma) := \varphi(s; x) + \frac{1}{2} \sigma \|s\|^2 + \psi(s; x).$$

Note that (3.1c) represents a regularized second-order model of the objective of (1.1), where  $f$  and  $h$  are modeled separately. More details, on the use of such model to solve (1.1), will be given in Section 4. By construction,  $\varphi(0; x) = f(x)$  and  $\nabla \varphi(0; x) = \nabla f(x)$ . We make the following assumption on (3.1b).

MODEL ASSUMPTION 3.1. For any  $x \in \mathbb{R}^n$ ,  $\psi(\cdot; x)$  is proper, lower semi continuous and prox-bounded with threshold  $\lambda_x \in \mathbb{R}_+ \cup \{+\infty\}$  [33, Definition 1.23]. In addition,  $\psi(0; x) = h(x)$ , and  $\partial \psi(0; x) = \partial h(x)$ .

We make the following additional assumption and say that  $\{\psi(\cdot; x)\}$  is uniformly prox-bounded.

MODEL ASSUMPTION 3.2. There is  $\lambda \in \mathbb{R}_+ \cup \{+\infty\}$  such that  $\lambda_x \geq \lambda$  for all  $x \in \mathbb{R}^n$ .

Model Assumption 3.2 is satisfied if  $h$  itself is prox-bounded and we select  $\psi(s; x) := h(x + s)$  for all  $x$ . Let

$$(3.2a) \quad p(x, \sigma) := \min_s m(s; x, \sigma) \leq m(0; x, \sigma) = f(x) + h(x)$$

$$(3.2b) \quad P(x, \sigma) := \underset{s}{\text{argmin}} m(s; x, \sigma),$$

be the value function and the set of minimizers of (3.1c), respectively.

For  $x \in \mathbb{R}^n$ ,  $s \in P(x, \sigma) \implies 0 \in \nabla \varphi(s; x) + \sigma s + \partial \psi(s; x)$ . Our first result states properties of the domain of  $p$  and  $P$  as given in (3.2a) and (3.2b).

LEMMA 3.1. Let Model Assumption 3.1 be satisfied and  $B(x) = B(x)^T$  for all  $x \in \mathbb{R}^n$ . Then,  $\text{dom } p = \mathbb{R}^n \times \mathbb{R}$ . In addition, if Model Assumption 3.2 holds,  $\text{dom } P \supseteq \{(x, \sigma) \mid \sigma > \max(\lambda^{-1} - \lambda_{\min}(B(x)), \lambda^{-1})\}$ , where  $\lambda_{\min}(B(x))$  is the smallest eigenvalue of  $B(x)$ .

*Proof.* By definition of the domain and [Model Assumption 3.1](#),

$$\begin{aligned} \text{dom } p &= \{(x, \sigma) \mid \inf_s m(s; x, \sigma) < +\infty\} = \{(x, \sigma) \mid \exists s m(s; x, \sigma) < +\infty\} \\ &= \{(x, \sigma) \mid \exists s \psi(s; x) < +\infty\} = \mathbb{R}^n \times \mathbb{R}, \end{aligned}$$

because  $\psi(\cdot; x)$  is proper. Moreover,

$$\text{dom } P = \{(x, \sigma) \mid \exists s(x, \sigma) \in \mathbb{R}^n, m(s(x, \sigma); x, \sigma) = \inf_s m(s; x, \sigma)\}.$$

Write

$$m(s; x, \sigma) = \varphi(s; x) + \frac{1}{2}(\sigma - \lambda^{-1})\|s\|^2 + \psi(s; x) + \frac{1}{2}\lambda^{-1}\|s\|^2.$$

By [Model Assumption 3.2](#) and [33, Exercise 1.24(c)], there is  $b \in \mathbb{R}$  such that  $\psi(s; x) + \frac{1}{2}\lambda^{-1}\|s\|^2 \geq b$  for all  $s \in \mathbb{R}^n$ . Let  $a \in \mathbb{R}$ . The above and (3.1a) imply that the level set  $\{s \in \mathbb{R}^n \mid m(s; x, \sigma) \leq a\}$  is contained in

$$\{s \in \mathbb{R}^n \mid \nabla f(x)^T s + \frac{1}{2}s^T(B(x) + (\sigma - \lambda^{-1})I)s \leq a - b - f(x)\},$$

which is a bounded set for  $\sigma > \lambda^{-1} - \lambda_{\min}(B(x))$ , i.e.,  $m(\cdot; x, \sigma)$  is level-bounded. Thus, [33, Theorem 1.9] implies that  $\inf_s m(s; x, \sigma)$  is attained, i.e., that  $P(x, \sigma) \neq \emptyset$ .

In [Lemma 3.1](#),  $\text{dom } P = \{(x, \sigma) \mid \sigma > \max(\lambda^{-1} - \lambda_{\min}(B(x)), \lambda^{-1})\}$  does not hold in general. Consider for example a situation where  $\psi(s; x)$  is bounded below for all  $x \in \mathbb{R}^n$ , i.e., each  $\lambda_x = +\infty$ . We can choose  $\lambda = +\infty$ . Assume also that, for a given  $x \in \mathbb{R}^n$ ,  $\varphi(s; x) = 0$  for all  $s$ , and  $\psi(s; x)$  level-bounded. Then,  $\lambda_{\min}(B(x)) = 0$ , and for  $\sigma = 0 = \lambda^{-1}$ ,  $m(s; x, \sigma) = \psi(s; x)$ . Therefore,  $P(x, \sigma) \neq \emptyset$ .

For a given  $s \in P(x, \sigma)$ , we define

$$(3.3) \quad \xi(s; x, \sigma) := f(x) + h(x) - (\varphi(s; x) + \psi(s; x)).$$

The next result relates (3.3) to first-order stationary for (1.1) and (3.1c).

**LEMMA 3.2.** *Let [Model Assumption 3.1](#) be satisfied, and  $x \in \mathbb{R}^n$  and  $\sigma \geq 0$  be given. Then, for  $s \in P(x, \sigma)$ ,  $\xi(s; x, \sigma) = 0 \implies s = 0 \implies x$  is first-order stationary for (1.1).*

*Proof.* For  $s \in P(x, \sigma)$ , if  $\xi(s; x, \sigma) = 0$ , then

$$0 = \xi(s; x, \sigma) = f(x) + h(x) - (\varphi(s; x) + \psi(s; x)) \geq \frac{1}{2}\sigma\|s\|^2,$$

which implies that  $s = 0$ , and therefore  $0 \in P(x, \sigma)$ . Therefore,  $0 \in \partial m(0; x, \sigma) = \nabla \varphi(0; x) + \partial \psi(0; x) = \nabla f(x) + \partial h(x)$ , and  $x$  is first-order stationary for (1.1).  $\square$

The following proposition states some properties of (3.2a) and (3.2b).

**PROPOSITION 3.3.** *Let [Model Assumptions 3.1](#) and [3.2](#) be satisfied. Assume also that  $\nabla f$  is bounded over  $\mathbb{R}^n$ . Let  $\epsilon > 0$ . Then,*

1. *at any  $(x, \sigma)$  such that  $\sigma \geq \lambda^{-1} - \lambda_{\min}(B(x)) + \epsilon$ ,  $p$  is finite and lsc, and  $P(x, \sigma)$  is nonempty and compact;*
2. *if  $\{(x_k, \sigma_k)\} \rightarrow (\bar{x}, \bar{\sigma})$  with  $\sigma_k \geq \lambda^{-1} - \lambda_{\min}(B(x_k)) + \epsilon$  for all  $k$  in such a way that  $\{p(x_k, \sigma_k)\} \rightarrow p(\bar{x}, \bar{\sigma})$ , and for each  $k$ ,  $s_k \in P(x_k, \sigma_k)$ , then  $\{s_k\}$  is bounded and all its limit points are in  $P(\bar{x}, \bar{\sigma})$ ;*
3. *for any  $x \in \mathbb{R}^n$ ,  $p(\bar{x}, \cdot)$  is continuous at any  $\bar{\sigma} \geq \lambda^{-1} - \lambda_{\min}(B(\bar{x})) + \epsilon$  and  $\{p(x_k, \sigma_k)\} \rightarrow p(\bar{x}, \bar{\sigma})$  holds in part 2.*

*Proof.* The proof consists in establishing that (3.1c) is level-bounded in  $s$  locally uniformly in  $(x, \sigma)$  [33, Definition 1.16] for  $\sigma \geq \lambda^{-1} - \lambda_{\min}(B(x)) + \epsilon$  and applying [33, Theorem 1.17]. It is nearly identical to that of [4, Proposition 3.2] and is omitted.  $\square$

Even though model (3.1c) is natural for incorporating second-order information, it is generally difficult to compute an exact minimizer of it. We proceed as Aravkin et al. [3, 4] and consider a simpler first-order model that will allow us to define an implementable stationary measure, to set minimal requirements steps computed in the course of the iterations of the algorithm of Section 4, and to derive convergence properties. This first-order model generalizes the concept of *Cauchy point* (“cp”) when solving (1.1). For fixed  $\nu > 0$  and  $x \in \mathbb{R}^n$ , define

$$(3.4a) \quad \varphi_{\text{cp}}(s; x) := f(x) + \nabla f(x)^T s$$

$$(3.4b) \quad m_{\text{cp}}(s; x, \nu^{-1}) := \varphi_{\text{cp}}(s; x) + \frac{1}{2}\nu^{-1}\|s\|^2 + \psi(s; x)$$

$$(3.4c) \quad p_{\text{cp}}(x, \nu^{-1}) := \min_s m_{\text{cp}}(s; x, \nu^{-1}) \leq m_{\text{cp}}(0; x, \nu^{-1}) = f(x) + h(x)$$

$$(3.4d) \quad P_{\text{cp}}(x, \nu^{-1}) := \operatorname{argmin}_s m_{\text{cp}}(s; x, \nu^{-1})$$

$$(3.4e) \quad \xi_{\text{cp}}(s_{\text{cp}}; x, \nu^{-1}) := f(x) + h(x) - (\varphi_{\text{cp}}(s_{\text{cp}}; x) + \psi(s_{\text{cp}}; x)),$$

where  $s_{\text{cp}} \in P_{\text{cp}}(x, \nu^{-1})$ . By [9, Lemma 2],

$$(3.5) \quad \xi_{\text{cp}}(s_{\text{cp}}; x, \nu^{-1}) \geq \frac{1}{2}\nu^{-1}\|s_{\text{cp}}(x, \nu^{-1})\|^2 \geq 0.$$

In the smooth case, i.e.,  $h = 0$  and  $\psi = 0$ ,  $s_{\text{cp}} = -\nu\nabla f(x)$ , so that

$$\xi_{\text{cp}}(s_{\text{cp}}; x, \nu^{-1}) = f(x) - (f(x) + \nabla f(x)^T s_{\text{cp}}) = \nu\|\nabla f(x)\|^2,$$

which suggests  $\nu^{-1/2}\xi_{\text{cp}}(s_{\text{cp}}; x, \nu^{-1})^{1/2}$  as a stationarity measure that generalizes the norm of the gradient to the nonsmooth setting.

Furthermore, this choice can be naturally interpreted in terms of the subdifferential of the model  $m_{\text{cp}}$  (3.4b) at  $s_{\text{cp}}$ . Specifically, from (3.5), we have

$$(3.6) \quad \sqrt{\nu^{-1}\xi_{\text{cp}}(s_{\text{cp}}; x, \nu^{-1})} \geq \frac{1}{\sqrt{2}}\nu^{-1}\|s_{\text{cp}}\|.$$

Moreover, by definition of  $s_{\text{cp}} \in P_{\text{cp}}(x, \nu^{-1})$ , the first-order optimality condition yields

$$-\nu^{-1}s_{\text{cp}} \in \nabla f(x) + \partial\psi(s_{\text{cp}}; x) = \partial m_{\text{cp}}(s_{\text{cp}}; x, \nu^{-1}).$$

Hence,

$$\operatorname{dist}(0, \partial m_{\text{cp}}(s_{\text{cp}}; x, \nu^{-1})) \leq \nu^{-1}\|s_{\text{cp}}\| \leq \sqrt{2\nu^{-1}\xi_{\text{cp}}(x, \nu^{-1})},$$

which directly links our stationarity measure to the distance to the subdifferential of the first-order model  $m_{\text{cp}}$  at  $s_{\text{cp}}$ .

It is worth noting that other authors, such as Kanzow and Mehlitz [23], adopt a different stationarity measure by directly considering  $\nu^{-1}\|s_{\text{cp}}\|$ .

The next results establish corresponding properties of  $p_{\text{cp}}$  and  $P_{\text{cp}}$ . The proofs are similar to those of Lemmas 3.1 and 3.2 and Proposition 3.3 and are omitted.

**LEMMA 3.4.** *Let Model Assumption 3.1 be satisfied. Then,  $\operatorname{dom} p_{\text{cp}} = \mathbb{R}^n \times \mathbb{R}$ . If Model Assumption 3.2 holds,  $\operatorname{dom} P_{\text{cp}} \supseteq \{(x, \nu^{-1}) \mid \nu > \max(\lambda^{-1} - \lambda_{\min}(B(x)), \lambda^{-1})\}$ .*

The next result characterizes first-order stationarity for (1.1).

LEMMA 3.5. *Let **Model Assumption 3.1** be satisfied and  $\nu > 0$ . Then, for  $s_{\text{cp}} \in P_{\text{cp}}(x, \nu^{-1})$ ,  $\xi_{\text{cp}}(s_{\text{cp}}; x, \nu^{-1}) = 0 \implies s_{\text{cp}} = 0 \implies x$  is first-order stationary for (1.1).*

The following result states properties of (3.4c) and (3.4d).

PROPOSITION 3.6. *Let **Model Assumptions 3.1** and **3.2** be satisfied and  $\nabla f(x)$  be bounded over  $\mathbb{R}^n$ . Let  $\epsilon > 0$ . Then,*

1. *at any  $(x, \nu^{-1})$  with  $\nu^{-1} \geq \lambda^{-1} + \epsilon$ ,  $p_{\text{cp}}$  is finite and lsc, and  $P_{\text{cp}}(x, \nu^{-1})$  is nonempty and compact;*
2. *if  $\{(x_k, \nu_k^{-1})\} \rightarrow (\bar{x}, \bar{\nu}^{-1})$  with  $\nu_k^{-1} \geq \lambda^{-1} + \epsilon$  for all  $k$  in such a way that  $\{p_{\text{cp}}(x_k, \nu_k^{-1})\} \rightarrow p_{\text{cp}}(\bar{x}, \bar{\nu}^{-1})$ , and for each  $k$ ,  $s_k \in P_{\text{cp}}(x_k, \nu_k^{-1})$ , then  $\{s_k\}$  is bounded and all its limit points are in  $P_{\text{cp}}(\bar{x}, \bar{\nu}^{-1})$ ;*
3. *for any  $\bar{x} \in \mathbb{R}^n$  and any  $\bar{\nu}^{-1} \geq \lambda^{-1} + \epsilon$ ,  $p_{\text{cp}}(\bar{x}, \cdot)$  is continuous at  $\bar{\nu}$  and  $\{p_{\text{cp}}(x_k, \nu_k^{-1})\} \rightarrow p_{\text{cp}}(\bar{x}, \bar{\nu}^{-1})$  holds in part 2.*

The main idea of the algorithm proposed in Section 4 is that (3.1c) is approximately minimized at each iteration. In order to establish convergence, the step  $s$  thus computed is required to satisfy *Cauchy decrease*, which we define as in [3, 4]:

$$(3.7) \quad \varphi(0; x) + \psi(0; x) - (\varphi(s; x) + \psi(s; x)) \geq (1 - \theta_1)\xi_{\text{cp}}(s_{\text{cp}}; x, \nu^{-1}),$$

for a preset value of  $\theta_1 \in (0, 1)$ . In other words,  $s$  must result in a decrease in  $\varphi(\cdot; x) + \psi(\cdot; x)$  that is at least a fraction of the decrease of the Cauchy model  $\varphi_{\text{cp}}(\cdot; x) + \psi(\cdot; x)$  obtained with the Cauchy step  $s_{\text{cp}}$  and a well-chosen step length  $\nu$ .

The following result parallels [24, Proposition 4] and establishes that if a step  $s$  reduces (3.1c) at least as much as  $s_{\text{cp}}$  does, Cauchy decrease holds. This observation is important because the first step of the proximal-gradient method from  $s = 0$  applied to (3.4b) and to (3.1c) with step length  $\nu$  is the same, and that step is precisely  $s_{\text{cp}}$ . Therefore, a step  $s$  may be obtained by continuing the proximal-gradient iterations on (3.1c) from  $s_{\text{cp}}$ .

PROPOSITION 3.7. *Let **Model Assumption 3.1** be satisfied. Let  $x \in \mathbb{R}^n$ ,  $\theta_1 \in (0, 1)$ ,  $\sigma > 0$  and let  $s_{\text{cp}}$  be computed with  $\nu = \theta_1 / (\|B(x)\| + \sigma)$ . Assume  $s \in \mathbb{R}^n$  is such that  $m(s; x, \sigma) \leq m(s_{\text{cp}}; x, \sigma)$ . Then,  $s$  satisfies (3.7).*

*Proof.* Let  $x \in \mathbb{R}^n$ ,  $\sigma > 0$ , and  $s \in \mathbb{R}^n$ , such that  $m(s; x, \sigma) \leq m(s_{\text{cp}}; x, \sigma)$ . Then,

$$\begin{aligned} \varphi(s; x) + \psi(s; x) + \frac{1}{2}\sigma\|s\|^2 &\leq \varphi(s_{\text{cp}}; x) + \psi(s_{\text{cp}}; x) + \frac{1}{2}\sigma\|s_{\text{cp}}\|^2 \\ &= \varphi_{\text{cp}}(s_{\text{cp}}; x) + \psi(s_{\text{cp}}; x) + \frac{1}{2}s_{\text{cp}}^T B(x)s_{\text{cp}} + \frac{1}{2}\sigma\|s_{\text{cp}}\|^2. \end{aligned}$$

The Cauchy-Schwarz inequality  $s_{\text{cp}}^T B(x)s_{\text{cp}} \leq \|B(x)\|\|s_{\text{cp}}\|^2$ , the identity  $\varphi(0; x) = \varphi_{\text{cp}}(0; x)$  and (3.5) yield

$$\begin{aligned} (\varphi + \psi)(0; x) - (\varphi + \psi)(s; x) &\geq \xi_{\text{cp}}(s_{\text{cp}}; x, \nu^{-1}) - \frac{1}{2}(\|B(x)\| + \sigma)\|s_{\text{cp}}\|^2 + \frac{1}{2}\sigma\|s\|^2 \\ &\geq \xi_{\text{cp}}(s_{\text{cp}}; x, \nu^{-1}) - \frac{1}{2}(\|B(x)\| + \sigma)\|s_{\text{cp}}\|^2 \\ &\geq \xi_{\text{cp}}(s_{\text{cp}}; x, \nu^{-1}) - (\|B(x)\| + \sigma)\nu\xi_{\text{cp}}(s_{\text{cp}}; x, \nu^{-1}) \\ &= (1 - \theta_1)\xi_{\text{cp}}(s_{\text{cp}}; x, \nu^{-1}). \quad \square \end{aligned}$$

Computing  $\|B(x)\|$  in the spectral norm comes at a cost. However, as we now illustrate, an inexact computation is sufficient in order to ensure (3.7). Assume that

we are able to compute  $\beta(x) \approx \|B(x)\|$  such that  $\beta(x) \geq \mu\|B(x)\|$  for  $0 < \mu < 1$ , and set  $\nu = \theta_1/(\beta(x) + \sigma)$ . The proof of [Proposition 3.7](#) continues to apply unchanged until the very last line, which becomes

$$\begin{aligned} \varphi(0; x) + \psi(0; x) - (\varphi(s; x) + \psi(s; x)) &\geq (1 - (\|B(x)\| + \sigma)\nu)\xi_{\text{cp}}(s_{\text{cp}}; x, \nu^{-1}) \\ &= \left(1 - \theta_1 \frac{\|B(x)\| + \sigma}{\beta(x) + \sigma}\right) \xi_{\text{cp}}(s_{\text{cp}}; x, \nu^{-1}). \end{aligned}$$

If  $\beta(x) \leq \|B(x)\|$ ,  $(\|B(x)\| + \sigma)/(\beta(x) + \sigma) \leq \|B(x)\|/\beta(x) \leq 1/\mu$ , so that

$$\left(1 - \theta_1 \frac{\|B(x)\| + \sigma}{\beta(x) + \sigma}\right) \xi_{\text{cp}}(s_{\text{cp}}; x, \nu^{-1}) \geq (1 - \theta_1/\mu)\xi_{\text{cp}}(s_{\text{cp}}; x, \nu^{-1}).$$

Thus, as long as  $\theta_1 < \mu$ , [\(3.7\)](#) is satisfied with  $\theta_1$  replaced with  $\theta_1/\mu$ .

If, on the other hand,  $\beta(x) \geq \|B(x)\|$ , then  $(\|B(x)\| + \sigma)/(\beta(x) + \sigma) \leq 1$ , and

$$\left(1 - \theta_1 \frac{\|B(x)\| + \sigma}{\beta(x) + \sigma}\right) \xi_{\text{cp}}(s_{\text{cp}}; x, \nu^{-1}) \geq (1 - \theta_1)\xi_{\text{cp}}(s_{\text{cp}}; x, \nu^{-1}),$$

and [\(3.7\)](#) holds unchanged.

The above observation also allows us to replace  $\|B(x)\|$  in the denominator of  $\nu$  with, e.g.,  $\|B(x)\|_1$ ,  $\|B(x)\|_\infty$  or  $\|B(x)\|_F$  if  $B(x)$  is available as an explicit matrix, or indeed with any other norm of  $B(x)$ .

**4. A modified quasi-Newton method for nonsmooth optimization.** We are in position to describe a modified quasi-Newton method to solve [\(1.1\)](#) named R2N. By contrast with trust-region-based approaches [\[3, 25\]](#), proximal operators are easier to evaluate in the R2N subproblem as they do not include a trust-region indicator.

At iteration  $k$ , we choose a step length  $\nu_k > 0$  based on the regularization parameter  $\sigma_k > 0$  and the norm of the model Hessian  $B(x_k)$  at the current iterate  $x_k \in \mathbb{R}^n$  as in [Proposition 3.7](#). We then compute the Cauchy step  $s_{k,\text{cp}}$  as a minimizer of [\(3.4b\)](#). A step  $s_k$  is subsequently computed that satisfies the assumptions of [Proposition 3.7](#).

The rest of the algorithm is standard. The decrease in  $f + h$  at  $x_k + s_k$  is compared to the decrease predicted by the model. If both are in sufficient agreement,  $x_k + s_k$  becomes the new iterate, and  $\sigma_k$  is possibly reduced. If the model turns out to predict poorly the actual decrease, the trial point is rejected and  $\sigma_k$  is increased. [Algorithm 4.1](#) states the whole procedure.

The interaction between  $\sigma_k$  and the unknown threshold  $\lambda_{x_k}$  works as in [\[3, Algorithm 6.1\]](#) and [\[4\]](#). If  $\sigma_k \leq \lambda_{x_k}^{-1}$ ,  $\psi(s_k; x_k) = -\infty$ , and according to the rules of extended arithmetic, which state that  $\pm\infty \cdot 0 = 0 \cdot (\pm\infty) = (\pm\infty)/(\pm\infty) := 0$  [\[33\]](#),  $\rho_k = 0$ . Consequently,  $s_k$  will be rejected at [Line 10](#), and  $\sigma_{k+1}$  will be set larger than  $\sigma_k$  at [Line 11](#). After a finite number of such increments,  $\sigma_k$  will surpass  $\lambda_{x_k}^{-1}$ , resulting in a step with finite  $\psi(s_k; x_k)$ . In effect, [Model Assumption 3.2](#) is only required to hold at the iterates generated by the algorithm.

Importantly, R2N does not require  $B_k \succeq 0$ , which may be useful in practice in order to capture natural problem curvature. In addition, we allow  $\{B_k\}$  to be unbounded. In [Section 5](#), we establish convergence provided it does not diverge too fast, using an assumption similar to that used in trust-region methods [\[14, §8.4\]](#). In [Section 6](#), we study the effect of using different bounds on  $\{\|B_k\|\}$  on worst-case evaluation complexity. The complexity results are obtained by adapting results from [\[18\]](#) (in the context of trust-region methods for smooth optimization) to R2N.

Our main working assumption is the following.

---

**Algorithm 4.1** R2N: A proximal modified Quasi-Newton method.

---

- 1: Choose constants  $0 < \theta_1 < 1 < \theta_2$ ,  $0 < \eta_1 \leq \eta_2 < 1$  and  $0 < \gamma_3 \leq 1 < \gamma_1 \leq \gamma_2$ .
- 2: Choose  $\sigma_0 > 0$  and  $x_0 \in \mathbb{R}^n$  where  $h$  is finite.
- 3: **for**  $k = 0, 1, \dots$  **do**
- 4:   Choose  $B_k := B(x_k) \in \mathbb{R}^{n \times n}$  such that  $B_k = B_k^T$ .
- 5:   Compute  $\nu_k := \theta_1 / (\|B_k\| + \sigma_k)$ .
- 6:   Compute  $s_{k,\text{cp}} \in \operatorname{argmin}_s m_{\text{cp}}(s; x_k, \nu_k^{-1})$  and  $\xi_{\text{cp}}(s_{k,\text{cp}}, x_k, \nu_k^{-1})$  as defined in (3.4e).
- 7:   Compute a step  $s_k$  such that  $m(s_k; x_k, \sigma_k) \leq m(s_{k,\text{cp}}; x_k, \sigma_k)$ .
- 8:   If  $\|s_k\| > \theta_2 \|s_{k,\text{cp}}\|$ , reset  $s_k = s_{k,\text{cp}}$ .
- 9:   Compute the ratio

$$\rho_k := \frac{f(x_k) + h(x_k) - (f(x_k + s_k) + h(x_k + s_k))}{\varphi(0; x_k) + \psi(0; x_k) - (\varphi(s_k; x_k) + \psi(s_k; x_k))}.$$

- 10:   If  $\rho_k \geq \eta_1$ , set  $x_{k+1} = x_k + s_k$ . Otherwise, set  $x_{k+1} = x_k$ .
- 11:   Update the regularization parameter according to

$$\sigma_{k+1} \in \begin{cases} [\gamma_3 \sigma_k, \sigma_k] & \text{if } \rho_k \geq \eta_2, & \text{very successful iteration} \\ [\sigma_k, \gamma_1 \sigma_k] & \text{if } \eta_1 \leq \rho_k < \eta_2, & \text{successful iteration} \\ [\gamma_1 \sigma_k, \gamma_2 \sigma_k] & \text{if } \rho_k < \eta_1. & \text{unsuccessful iteration} \end{cases}$$

12: **end for**

---

**PROBLEM ASSUMPTION 4.1.** *The function  $f$  is continuously differentiable over the set  $\{x \in \mathbb{R}^n \mid (f+h)(x) \leq (f+h)(x_0)\}$  and  $h$  is proper and lower semi-continuous.*

**Problem Assumption 4.1** is very mild as one does not require boundedness nor Lipschitz continuity of  $f$  or  $\nabla f$ , in contrast with [4, Problem Assumption 4.1] or the assumptions of [Kanzow and Lechner \[22\]](#). For instance, our analysis includes cases where  $f$  is continuously differentiable, but whose gradient is not locally Lipschitz continuous at  $x = 0$ , e.g.,  $f(x) = |x|^{\frac{3}{2}}$ .

In the next sections, we derive convergence and worst-case complexity analysis for [Algorithm 4.1](#). We will repeatedly use the notation

- (4.1a)  $\mathcal{S} := \{i \in \mathbb{N} \mid \rho_i \geq \eta_1\}$  (all successful iterations)
- (4.1b)  $\mathcal{S}_k := \{i \in \mathcal{S} \mid i \leq k\}$  (successful iterations until iteration  $k$ )
- (4.1c)  $\mathcal{U} := \{i \in \mathbb{N} \mid \rho_i < \eta_1\}$  (all unsuccessful iterations)
- (4.1d)  $\mathcal{U}_k := \{i \in \mathbb{N} \mid i \notin \mathcal{S}, i \leq k\}$  (unsuccessful iterations until iteration  $k$ ).

**5. Convergence analysis of Algorithm 4.1.** In this section, we investigate the convergence properties of [Algorithm 4.1](#) under [Problem Assumption 4.1](#). For notational convenience, we denote  $\xi_{\text{cp}}(s_{k,\text{cp}}, x_k, \nu_k^{-1})$  by  $\xi_{k,\text{cp}}$ . We then show that  $\liminf_{k \rightarrow \infty} \nu_k^{-\frac{1}{2}} \xi_{k,\text{cp}}^{\frac{1}{2}} = 0$ . We stress that the obtained convergence properties of [Algorithm 4.1](#) are more general than those of [4, 22, 23], and do not require boundedness of the model Hessians nor (local) Lipschitz continuity of  $\nabla f$ .

We first establish lower bounds on  $\xi_{k,\text{cp}}$  in terms of  $\|s_k\|$ .

**LEMMA 5.1.** *For all  $k \in \mathbb{N}$ ,*

$$(5.1) \quad \xi_{k,\text{cp}} \geq \frac{1}{2\theta_2} \nu_k^{-1} \|s_k\|^2.$$

Additionally, for any  $\alpha > 0$ ,

$$(5.2) \quad \nu_k^{-\frac{1}{2}} \xi_{k,\text{cp}}^{\frac{1}{2}} \geq \alpha \quad \Rightarrow \quad \xi_{\text{cp}}(x_k; \nu_k^{-1}) \geq \frac{\alpha}{\theta_2 \sqrt{2}} \|s_k\|.$$

*Proof.* From [Algorithm 4.1](#), we have  $\|s_k\| \leq \theta_2 \|s_{k,\text{cp}}\|$ . Hence,

$$\xi_{k,\text{cp}} \geq \frac{1}{2} \nu_k^{-1} \|s_{k,\text{cp}}\|^2 \geq \frac{1}{2\theta_2^2} \nu_k^{-1} \|s_k\|^2.$$

If  $\nu_k^{-\frac{1}{2}} \xi_{k,\text{cp}}^{\frac{1}{2}} \geq \alpha$ ,

$$\xi_{k,\text{cp}} \geq \alpha \nu_k^{\frac{1}{2}} \xi_{k,\text{cp}}^{\frac{1}{2}} \geq \alpha \nu_k^{\frac{1}{2}} \left( \frac{1}{2\theta_2^2} \nu_k^{-1} \|s_k\|^2 \right)^{\frac{1}{2}} = \frac{\alpha}{\theta_2 \sqrt{2}} \|s_k\|. \quad \square$$

The next lemma shows that the convergence of  $\{x_k\}_{k \in \mathbb{N}}$  holds if the objective is bounded below, the algorithm generates infinitely many successful iterations and the stationarity measure  $\nu_k^{-1/2} \xi_{k,\text{cp}}^{1/2}$  is bounded away from zero.

**LEMMA 5.2.** *Assume that [Algorithm 4.1](#) generates infinitely many successful iterations and that there is  $(f+h)_{\text{low}} \in \mathbb{R}$  such that  $(f+h)(x_k) \geq (f+h)_{\text{low}}$  for all  $k \in \mathbb{N}$ . Additionally, assume, that there is  $\alpha > 0$  such that for all  $k \in \mathbb{N}$ ,  $\nu_k^{-1/2} \xi_{k,\text{cp}}^{1/2} \geq \alpha$ . Then,  $\{x_k\}_{k \in \mathbb{N}}$  is a Cauchy sequence, and hence converges.*

*Proof.* For all  $k \in \mathcal{S}$ , using (5.2) from [Lemma 5.1](#), we have

$$\begin{aligned} f(x_k) + h(x_k) - f(x_{k+1}) - h(x_{k+1}) &\geq \eta_1(1 - \theta_1) \xi_{k,\text{cp}} \\ &\geq \frac{\eta_1(1 - \theta_1)\alpha}{\theta_2 \sqrt{2}} \|s_k\| = \frac{\eta_1(1 - \theta_1)\alpha}{\theta_2 \sqrt{2}} \|x_{k+1} - x_k\|. \end{aligned}$$

Summing over all successful iterations from 1 to  $k$ , we obtain

$$\begin{aligned} f(x_0) + h(x_0) - (f+h)_{\text{low}} &\geq \sum_{j \in \mathcal{S}_k} f(x_j) + h(x_j) - f(x_{j+1}) - h(x_{j+1}) \\ &\geq \frac{\eta_1(1 - \theta_1)\alpha}{\theta_2 \sqrt{2}} \sum_{j \in \mathcal{S}_k} \|x_{j+1} - x_j\| \\ &= \frac{\eta_1(1 - \theta_1)\alpha}{\theta_2 \sqrt{2}} \sum_{j=0}^k \|x_{j+1} - x_j\|. \end{aligned}$$

Thus,  $\sum_{j \in \mathbb{N}} \|x_{j+1} - x_j\| < +\infty$ . Hence,  $\{x_k\}_{k \in \mathbb{N}}$  is a Cauchy sequence, and converges.  $\square$

The following lemma shows that when  $\{x_k\}$  converges and  $\{\sigma_k\}$  diverges along common subsequences, the corresponding subsequence of  $\{s_k\}$  converges to zero.

**LEMMA 5.3.** *Let [Problem Assumption 4.1](#) and [Model Assumptions 3.1](#) and [3.2](#) be satisfied and assume that there is an index set  $\mathcal{K} \subseteq \mathbb{N}$  such that  $\lim_{k \in \mathcal{K}} \sigma_k = +\infty$  and  $\{x_k\}_{k \in \mathcal{K}}$  is bounded. Then,  $\lim_{k \in \mathcal{K}} s_k = \lim_{k \in \mathcal{K}} s_{k,\text{cp}} = 0$ .*

*Proof.* By contradiction, assume that there is an index set  $\mathcal{K}' \subseteq \mathcal{K}$  and  $\alpha > 0$  such that  $\|s_{k,\text{cp}}\| \geq \alpha$  for all  $k \in \mathcal{K}'$ . By definition of  $s_{k,\text{cp}}$  and [Model Assumption 3.1](#),

$f(x_k) + h(x_k) = m_{\text{cp}}(0; x_k, \nu_k^{-1}) \geq m_{\text{cp}}(s_{k,\text{cp}}; x_k, \nu_k^{-1})$ . Hence,

$$\begin{aligned} f(x_k) + h(x_k) &\geq \varphi_{\text{cp}}(s_{k,\text{cp}}; x_k) + \psi(s_{k,\text{cp}}; x_k) + \frac{1}{2}\nu_k^{-1}\|s_{k,\text{cp}}\|^2 \\ &= f(x_k) + \nabla f(x_k)^T s_{k,\text{cp}} + \frac{1}{2\theta_1}(\|B_k\| + \sigma_k)\|s_{k,\text{cp}}\|^2 + \psi(s_{k,\text{cp}}; x_k) \\ &\geq f(x_k) + \nabla f(x_k)^T s_{k,\text{cp}} + \frac{1}{2\theta_1}\sigma_k\|s_{k,\text{cp}}\|^2 + \psi(s_{k,\text{cp}}; x_k) \\ &\geq f(x_k) - \|\nabla f(x_k)\|\|s_{k,\text{cp}}\| + \frac{1}{2}\left(\frac{\sigma_k}{\theta_1} - \lambda^{-1}\right)\|s_{k,\text{cp}}\|^2 \\ &\quad + \psi(s_{k,\text{cp}}; x_k) + \frac{1}{2}\lambda^{-1}\|s_{k,\text{cp}}\|^2. \end{aligned}$$

By [Model Assumption 3.2](#) and [33, Exercise 1.24(c)], there is  $b_h \in \mathbb{R}$  such that  $\psi(s; x) + \frac{1}{2}\lambda^{-1}\|s\|^2 \geq b_h$  for all  $s$  and  $x$ . Hence, for all sufficiently large  $k \in \mathcal{K}'$ ,  $\sigma_k > \lambda^{-1}$  and

$$\begin{aligned} f(x_k) + h(x_k) &\geq f(x_k) - \|\nabla f(x_k)\|\|s_{k,\text{cp}}\| + \frac{1}{2}\left(\frac{\sigma_k}{\theta_1} - \lambda^{-1}\right)\|s_{k,\text{cp}}\|^2 + b_h \\ &\geq f(x_k) - \|\nabla f(x_k)\|\|s_{k,\text{cp}}\| + \frac{1}{2}\alpha\left(\frac{\sigma_k}{\theta_1} - \lambda^{-1}\right)\|s_{k,\text{cp}}\| + b_h \\ (5.3) \quad &= f(x_k) + \left(\frac{1}{2}\alpha\left(\frac{\sigma_k}{\theta_1} - \lambda^{-1}\right) - \|\nabla f(x_k)\|\right)\|s_{k,\text{cp}}\| + b_h. \end{aligned}$$

Since  $\{x_k\}_{k \in \mathcal{K}'}$  is bounded, so are  $\{f(x_k)\}_{k \in \mathcal{K}'}$  and  $\{\nabla f(x_k)\}_{k \in \mathcal{K}'}$  by [Problem Assumption 4.1](#). Let  $b_f := \min_{k \in \mathcal{K}'} f(x_k) > -\infty$  and  $b_{f'} = \max_{k \in \mathcal{K}'} \|\nabla f(x_k)\| < \infty$ . Because  $\{f(x_k) + h(x_k)\}$  is nonincreasing, (5.3) yields

$$(5.4) \quad f(x_0) + h(x_0) \geq f(x_k) + h(x_k) \geq b_f + \left(\frac{1}{2}\alpha\left(\frac{\sigma_k}{\theta_1} - \lambda^{-1}\right) - b_{f'}\right)\|s_{k,\text{cp}}\| + b_h.$$

As  $\lim_{k \in \mathcal{K}'} \sigma_k = +\infty$ , for  $k$  sufficiently large,  $\frac{1}{2}\alpha\left(\frac{\sigma_k}{\theta_1} - \lambda^{-1}\right) > b_{f'}$ . Thus, for all sufficiently large  $k \in \mathcal{K}'$ , (5.4) combines with  $\|s_{k,\text{cp}}\| \geq \alpha$  to give

$$f(x_0) + h(x_0) \geq b_f + \left(\frac{1}{2}\alpha\left(\frac{\sigma_k}{\theta_1} - \lambda^{-1}\right) - b_{f'}\right)\alpha + b_h,$$

which is a contradiction because the right-hand side diverges. Thus,  $\lim_{k \in \mathcal{K}} \|s_{k,\text{cp}}\| = 0$ . Finally, since  $\|s_k\| \leq \theta_2\|s_{k,\text{cp}}\|$ , we get also  $\lim_{k \in \mathcal{K}} \|s_k\| = 0$ .  $\square$

For the remainder of this section, we need the following assumption.

**MODEL ASSUMPTION 5.1.** *For all  $k \in \mathbb{N}$ , the model function  $\psi(\cdot, x_k)$  satisfies*

$$(5.5) \quad |h(x_k + s_k) - \psi(s_k; x_k)| = o(\|s_k\|) \quad \text{as } s_k \rightarrow 0.$$

[Model Assumption 5.1](#) is trivially satisfied if, at each iteration  $k$ , we set  $\psi(s; x_k) = h(x_k + s)$ , which is what [Kanzow and Lechner \[22\]](#) do. However, the assumption also holds when  $h(x) = g(c(x))$ , where  $c : \mathbb{R}^n \rightarrow \mathbb{R}^m$  has Lipschitz-continuous or  $\alpha_h$ -Hölder-continuous Jacobian,  $g : \mathbb{R}^m \rightarrow \mathbb{R}$  is  $L$ -Lipschitz continuous, and we choose  $\psi(s; x_k) := g(c(x_k) + \nabla c(x_k)^T s)$ . Indeed, there exists  $M > 0$  such that  $|h(x_k + s) - \psi(s; x_k)| \leq L\|c(x_k + s) - c(x_k) - \nabla c(x_k)^T s\| \leq LM\|s\|^{1+\alpha_h} = o(\|s\|)$ .

**THEOREM 5.4.** *Let [Problem Assumption 4.1](#) and [Model Assumptions 3.1, 3.2](#) and [5.1](#) be satisfied. Assume that there is an index set  $\mathcal{K} \subseteq \mathbb{N}$  so that (i) there is  $\alpha > 0$  such that  $\nu_k^{-1/2}\xi_{k,\text{cp}}^{1/2} \geq \alpha$  for all  $k \in \mathcal{K}$ , (ii)  $\{\sigma_k(1 + \|B_k\|)^{-1}\}_{k \in \mathcal{K}}$  is unbounded and (iii)  $\{x_k\}_{k \in \mathcal{K}}$  is bounded. Then, there is an index set  $\mathcal{K}' \subseteq \mathcal{K}$  such that for all  $k \in \mathcal{K}'$  sufficiently large,  $k$  is a very successful iteration.*

*Proof.* By Assumption (ii), there is an index set  $\mathcal{K}' \subset \mathcal{K}$  such that  $\lim_{k \in \mathcal{K}'} \sigma_k(1 + \|B_k\|)^{-1} = \infty$ . Since  $\sigma_k \geq \sigma_k(1 + \|B_k\|)^{-1}$ , we also have  $\lim_{k \in \mathcal{K}'} \sigma_k = \infty$ . [Lemma 5.3](#) then implies  $\lim_{k \in \mathcal{K}'} \|s_k\| = 0$ . For all  $k \in \mathcal{K}'$ , [Model Assumption 5.1](#) combines with [\(3.7\)](#) and a Taylor expansion of  $f$  about  $x_k$  to give

$$\begin{aligned}
|\rho_k - 1| &= \left| \frac{(f+h)(x_k+s_k) - (\varphi(s_k; x_k) + \psi(s_k; x_k))}{\varphi(0; x_k) + \psi(0; x_k) - (\varphi(s_k; x_k) + \psi(s_k; x_k))} \right| \\
&= \left| \frac{(f+h)(x_k+s_k) - (f(x_k) + \nabla f(x_k)^T s_k + \frac{1}{2} s_k^T B_k s_k + \psi(s_k; x_k))}{\varphi(0; x_k) + \psi(0; x_k) - (\varphi(s_k; x_k) + \psi(s_k; x_k))} \right| \\
&\leq \frac{|f(x_k+s_k) - f(x_k) - \nabla f(x_k)^T s_k|}{(1-\theta_1)\xi_{k,\text{cp}}} + \frac{\|B_k\| \|s_k\|^2}{2(1-\theta_1)\xi_{k,\text{cp}}} \\
&\quad + \frac{|h(x_k+s_k) - \psi(s_k; x_k)|}{(1-\theta_1)\xi_{k,\text{cp}}} \\
&= \frac{o(\|s_k\|)}{(1-\theta_1)\xi_{k,\text{cp}}} + \frac{\|B_k\| \|s_k\|^2}{2(1-\theta_1)\xi_{k,\text{cp}}} + \frac{o(\|s_k\|)}{(1-\theta_1)\xi_{k,\text{cp}}} \\
(5.6) \quad &\leq \frac{o(\|s_k\|)}{\xi_{k,\text{cp}}} + \frac{(1 + \|B_k\|) \|s_k\|^2}{2(1-\theta_1)\xi_{k,\text{cp}}}.
\end{aligned}$$

By Assumption (i), [Lemma 5.1](#) implies  $\xi_{\text{cp}}(x_k; \nu_k^{-1}) \geq \frac{\alpha}{\theta_2 \sqrt{2}} \|s_k\|$  for all  $k \in \mathcal{K}'$ , which we apply to the first term in the right-hand side of [\(5.6\)](#). Similarly, [\(5.1\)](#) implies

$$\xi_{k,\text{cp}} \geq \frac{1}{2\theta_2^2} \nu_k^{-1} \|s_k\|^2 = \frac{1}{2\theta_1 \theta_2^2} (\|B_k\| + \sigma_k) \|s_k\|^2 \geq \frac{1}{2\theta_1 \theta_2^2} \sigma_k \|s_k\|^2,$$

which we apply to the second term in the right-hand side of [\(5.6\)](#). Hence, [\(5.6\)](#) simplifies to

$$(5.7) \quad |\rho_k - 1| \leq \frac{o(\|s_k\|)}{\frac{\alpha}{\theta_2 \sqrt{2}} \|s_k\|} + \frac{(1 + \|B_k\|) \|s_k\|^2}{\frac{(1-\theta_1)}{\theta_1 \theta_2^2} \sigma_k \|s_k\|^2} = \frac{o(\|s_k\|)}{\|s_k\|} + \frac{\theta_1 \theta_2^2}{(1-\theta_1) \sigma_k (1 + \|B_k\|)^{-1}}.$$

By Assumption (ii), the right-hand side of [\(5.7\)](#) converges to zero. Thus, for all sufficiently large  $k \in \mathcal{K}'$ ,  $|\rho_k - 1| \leq 1 - \eta_2$ , which implies that  $\rho_k \geq \eta_2$ .  $\square$

[Theorem 5.4](#) shares similarities with [\[4, Theorem 4.1\]](#) but uses weaker assumptions. In particular, compared to [\[4, Theorem 4.1\]](#), we do not use the Lipschitz continuity of  $\nabla f$  nor do we require model Hessians to be uniformly bounded.

**LEMMA 5.5.** *Let [Problem Assumption 4.1](#) and [Model Assumptions 3.1, 3.2](#) and [5.1](#) be satisfied. Assume that  $\{x_k\}_{k \in \mathbb{N}}$  is bounded and that there is  $\alpha > 0$  such that for all  $k \in \mathbb{N}$ ,  $\nu_k^{-1/2} \xi_{k,\text{cp}}^{1/2} \geq \alpha$ . Then,  $\{\sigma_k(1 + \|B_k\|)^{-1}\}_{k \in \mathbb{N}}$  is bounded.*

*Proof.* Assume, by contradiction, that  $\{\sigma_k(1 + \|B_k\|)^{-1}\}_{k \in \mathbb{N}}$  is unbounded. Since  $\{\sigma_k\}_{k \in \mathbb{N}}$  increases only on unsuccessful iterations and  $\{(1 + \|B_k\|)^{-1}\}_{k \in \mathbb{N}}$  is bounded,  $\{\sigma_k(1 + \|B_k\|)^{-1}\}_{k \in \mathcal{U}}$  must be unbounded, where  $\mathcal{U}$  is defined in [\(4.1\)](#). Hence, using [Theorem 5.4](#), we deduce that there is an index set  $\mathcal{U}' \subseteq \mathcal{U}$  such that for all  $k \in \mathcal{U}'$  sufficiently large,  $k$  is a very successful iteration, i.e.,  $k \in \mathcal{S}$ , which is absurd.  $\square$

Consider the following assumption

MODEL ASSUMPTION 5.2. *The sequence  $\{B_k\}_{k \in \mathbb{N}}$  satisfies:*

$$\sum_{k \in \mathbb{N}} \frac{1}{r_k} = +\infty, \quad r_k := \max_{0 \leq j \leq k} \|B_j\| + 1.$$

The next theorem examines the case where [Algorithm 4.1](#) generates only a finite number of successful iterations. The proof of the second part of the theorem, which establishes stationarity of the limit point, follows ideas similar to those in [\[23, Lemma 3.3\]](#).

**THEOREM 5.6.** *Let [Problem Assumption 4.1](#) and [Model Assumptions 3.1, 3.2, 5.1](#) and [5.2](#) be satisfied. If [Algorithm 4.1](#) generates finitely many successful iterations, then there is  $x^* \in \mathbb{R}^n$  such that  $x_k = x^*$  for all sufficiently large  $k$ , and  $\liminf_{k \rightarrow \infty} \nu_k^{-\frac{1}{2}} \xi_{\text{cp}}(x_k, \nu_k^{-1})^{\frac{1}{2}} = 0$ . Moreover,  $x^*$  is stationary for [\(1.1\)](#).*

*Proof.* First, we prove that  $\liminf_{k \rightarrow \infty} \nu_k^{-\frac{1}{2}} \xi_{\text{cp}}(x_k, \nu_k^{-1})^{\frac{1}{2}} = 0$ . Assume, by contradiction, that there is  $\alpha > 0$  such that  $\nu_k^{-\frac{1}{2}} \xi_{\text{cp}}(x_k, \nu_k^{-1})^{\frac{1}{2}} \geq \alpha$  for all  $k \in \mathbb{N}$  and let  $k_f$  be the last successful iteration. Hence,  $x_k = x_{k_f}$  for all  $k \geq k_f$  and  $\lim_k x_k = x_{k_f} = x^*$ . Using [Lemma 5.5](#),  $\{\sigma_k(1 + \|B_k\|)^{-1}\}_{k \in \mathbb{N}}$  is bounded by a constant  $b_\sigma > 0$ . This, implies, that for all  $k > k_f$ ,

$$\frac{1}{r_k} = \frac{1}{1 + \max_{0 \leq j \leq k} \|B_j\|} \leq \frac{1}{1 + \|B_k\|} \leq \frac{b_\sigma}{\sigma_k}.$$

Thus, [Model Assumption 5.2](#) implies

$$(5.8) \quad \sum_{k=k_f+1}^{\infty} \frac{1}{\sigma_k} = +\infty.$$

On the other hand, all  $k > k_f$ ,  $k$  is an unsuccessful iteration. The mechanism of [Algorithm 4.1](#) then ensures  $\frac{\sigma_k}{\sigma_{k+1}} \leq \frac{1}{\gamma_1} < 1$  for all  $k > k_f$ . But this implies that  $\sum_{k=k_f+1}^{\infty} \frac{1}{\sigma_k}$  converges, which contradicts [\(5.8\)](#). Hence,  $\liminf_{k \rightarrow \infty} \nu_k^{-\frac{1}{2}} \xi_{\text{cp}}(x_k, \nu_k^{-1})^{\frac{1}{2}} = 0$ .

The first part of the theorem implies existence of  $\mathcal{K} \subseteq \mathbb{N}$  such that

$$\lim_{k \in \mathcal{K}} \nu_k^{-\frac{1}{2}} \xi_{\text{cp}}(x_k, \nu_k^{-1})^{\frac{1}{2}} = 0.$$

Let  $\mathcal{K}_f := \{k \in \mathcal{K} \mid k > k_f\}$ . By definition of  $s_{k,\text{cp}}$ ,

$$(5.9) \quad \nabla f(x^*)^T s_{k,\text{cp}} + \frac{1}{2} \nu_k^{-1} \|s_{k,\text{cp}}\|^2 + \psi(s_{k,\text{cp}}; x^*) \leq \psi(0; x^*) \quad (k \in \mathcal{K}_f).$$

Since  $\mathcal{K}_f \subseteq \mathcal{U}$  by assumption,  $\lim_{k \in \mathcal{K}_f} \sigma_k = +\infty$ , and  $\{x_k\}_{k \in \mathcal{K}_f}$  is constant, hence bounded. Thus, [Lemma 5.3](#) implies that

$$(5.10) \quad \lim_{k \in \mathcal{K}_f} s_{k,\text{cp}} = 0.$$

Because  $\lim_{k \in \mathcal{K}_f} \nu_k^{-\frac{1}{2}} \xi_{\text{cp}}(x_k, \nu_k^{-1})^{\frac{1}{2}} = 0$ , [\(3.5\)](#) implies that

$$(5.11) \quad \lim_{k \in \mathcal{K}_f} \nu_k^{-1} s_{k,\text{cp}} = 0.$$

By taking the limit superior in (5.9) over  $\mathcal{K}_f$ , and using (5.10) and (5.11),

$$(5.12) \quad \limsup_{k \in \mathcal{K}_f} \psi(s_{k,\text{cp}}; x^*) \leq \psi(0; x^*).$$

Using (5.10) and the lower semi-continuity of  $\psi(\cdot; x^*)$  at 0 (from [Model Assumption 3.1](#)), we have  $\liminf_{k \in \mathcal{K}_f} \psi(s_{k,\text{cp}}; x^*) \geq \psi(0; x^*)$ . Hence,

$$(5.13) \quad \lim_{k \in \mathcal{K}_f} \psi(s_{k,\text{cp}}; x^*) = \psi(0; x^*).$$

On the other hand, since  $s_{k,\text{cp}} \in P(x_k, \nu_k^{-1})$ ,

$$-\nu_k^{-1} s_{k,\text{cp}} \in \nabla f(x^*) + \partial\psi(s_{k,\text{cp}}; x^*) \quad (k \in \mathcal{K}_f).$$

Thus, by using (5.11),

$$(5.14) \quad 0 \in \nabla f(x^*) + \limsup_{k \in \mathcal{K}_f} \partial\psi(s_{k,\text{cp}}; x^*).$$

Using (5.10) and (5.13), [33, Proposition 8.7] implies that

$$\limsup_{k \in \mathcal{K}_f} \partial\psi(s_{k,\text{cp}}; x^*) \subseteq \partial\psi(0; x^*),$$

which, combined with (5.14), gives

$$0 \in \nabla f(x^*) + \partial\psi(0; x^*) = \nabla f(x^*) + \partial h(x^*).$$

In other words,  $x^*$  is stationary for (1.1).  $\square$

To the best of our knowledge, in the case of a finite number of successful iterations, [Theorem 5.6](#) is the first convergence result of a regularized or trust-region method that does not rely on the boundedness of the regularization parameter or trust-region radius. Remarkably, even in the absence of such boundedness, the theorem establishes the stationarity of the limit point  $x^*$ , despite the fact that the only a subsequence of the stationarity measure  $\nu_k^{-1/2} \xi_{\text{cp}}(x_k, \nu_k^{-1})^{1/2}$  converges to zero.

Now we consider the case where the number of successful iterations is infinite. Let  $\tau \in \mathbb{N}_0$  and define, as in [18],

$$(5.15a) \quad \mathcal{T}_k^\tau = \{j = 0, \dots, k \mid j < \tau |\mathcal{S}_j|\},$$

$$(5.15b) \quad \mathcal{W}_k^\tau = \{j = 0, \dots, k \mid j \geq \tau |\mathcal{S}_j|\}.$$

The sets (5.15) are the sets defined in [18, Equation (52)] with  $\lambda := 0$ . For readability, we drop the superscript  $\lambda$  in the following. The next lemma provides a series comparison result that will be used in the proof of the main theorem.

**LEMMA 5.7** (18, Lemma 10). *Let  $\{r_j\}_{j \in \mathbb{N}}$  be a non-decreasing positive real sequence. For any  $k \in \mathbb{N}$ ,*

$$\tau \sum_{j \in \mathcal{S}_k} \frac{1}{r_j} \geq \sum_{j \in \mathcal{T}_k^\tau} \frac{1}{r_j} = \sum_{j=0}^k \frac{1}{r_j} - \sum_{j \in \mathcal{W}_k^\tau} \frac{1}{r_j},$$

where  $\mathcal{T}_k^\tau$  and  $\mathcal{W}_k^\tau$  are defined in (5.15).

The following lemma plays a key role in deriving a convergence result in the case where the number of successful iterations is infinite.

LEMMA 5.8. *Let **Problem Assumption 4.1** and **Model Assumptions 3.1, 3.2, 5.1** and **5.2** be satisfied. Assume that (i)  $\tau \in \mathbb{N}_0$  is chosen so that  $\gamma_3 \gamma_1^{\tau-1} > 1$ , (ii) there is  $\alpha > 0$  such that  $\nu_k^{-1/2} \xi_{k,\text{cp}}^{1/2} \geq \alpha$  for all  $k \in \mathbb{N}$ , and (iii)  $\{x_k\}_{k \in \mathbb{N}}$  is bounded. Then,  $\left\{ \sum_{j \in \mathcal{W}_k^\tau} \frac{1}{r_j} \right\}_{k \in \mathbb{N}}$  is bounded, where  $r_j$  is as in **Model Assumption 5.2**.*

*Proof.* For any  $j \geq 0$ , **Lemma 5.5** and the update mechanism of **Algorithm 4.1** imply that there is  $b_\sigma > 0$  such that

$$\frac{1}{r_j} \leq \frac{1}{1 + \|B_j\|} \leq \frac{b_\sigma}{\sigma_j} \leq \frac{b_\sigma}{\gamma_3^{|\mathcal{S}_j|} \gamma_1^{|\mathcal{U}_j|} \sigma_0} = \frac{b_\sigma}{\gamma_3^{|\mathcal{S}_j|} \gamma_1^{j-|\mathcal{S}_j|} \sigma_0}.$$

Consider now  $k \geq 0$  and  $j \in \mathcal{W}_k^\tau$ . Then,  $j > \tau|\mathcal{S}_j|$ , which, together with the fact that  $\gamma_1 > 1$  and  $0 < \gamma_3 \leq 1$  leads to

$$\frac{1}{r_j} \leq \frac{b_\sigma}{\gamma_3^{|\mathcal{S}_j|} \gamma_1^{j-|\mathcal{S}_j|} \sigma_0} < \frac{b_\sigma}{\gamma_3^{j/\tau} \gamma_1^{j-j/\tau} \sigma_0} = \frac{b_\sigma}{(\gamma_3 \gamma_1^{\tau-1})^{j/\tau} \sigma_0}.$$

We sum the above inequalities over  $j \in \mathcal{W}_k^\tau$  and use the fact that  $\gamma_3 \gamma_1^{\tau-1} > 1$  to obtain

$$\sum_{j \in \mathcal{W}_k^\tau} \frac{1}{r_j} < \frac{b_\sigma}{\sigma_0} \sum_{j \in \mathcal{W}_k^\tau} \frac{1}{(\gamma_3 \gamma_1^{\tau-1})^{j/\tau}} \leq \frac{b_\sigma}{\sigma_0} \sum_{j \in \mathbb{N}} \frac{1}{(\gamma_3 \gamma_1^{\tau-1})^{j/\tau}} < \infty. \quad \square$$

We state now our main convergence result.

THEOREM 5.9. *Let **Problem Assumption 4.1** and **Model Assumptions 3.1, 3.2, 5.1** and **5.2** be satisfied. Assume that **Algorithm 4.1** generates infinitely many successful iterations and that there is  $(f+h)_{\text{low}} \in \mathbb{R}$  such that  $(f+h)(x_k) \geq (f+h)_{\text{low}}$  for all  $k \in \mathbb{N}$ . Then,  $\liminf_{k \rightarrow +\infty} \nu_k^{-1/2} \xi_{k,\text{cp}}^{1/2} = 0$ .*

*Proof.* By contradiction, assume that there is  $\alpha > 0$  such that for all  $k \in \mathbb{N}$ ,  $\nu_k^{-1/2} \xi_{k,\text{cp}}^{1/2} \geq \alpha$ . **Lemma 5.2** shows that  $\{x_k\}_{k \in \mathbb{N}}$  is convergent, hence bounded. **Lemma 5.5** then implies that  $\{\sigma_k(1 + \|B_k\|)^{-1}\}_{k \in \mathbb{N}}$  is bounded, say by  $b_\sigma > 0$ . Equivalently,  $\sigma_k \leq b_\sigma(1 + \|B_k\|)$ . For any  $j \in \mathcal{S}$ ,  $\rho_j \geq \eta_1$ , and

$$\begin{aligned} f(x_j) + h(x_j) - f(x_{j+1}) - h(x_{j+1}) &\geq \eta_1(1 - \theta_1) \xi_{\text{cp}}(x_j, \nu_j^{-1}) \\ &\geq \eta_1(1 - \theta_1) \nu_j \alpha^2 \\ &= \frac{\eta_1(1 - \theta_1) \theta_1 \alpha^2}{\sigma_j + \|B_j\|} \\ &\geq \frac{\eta_1(1 - \theta_1) \theta_1 \alpha^2}{(b_\sigma(1 + \|B_j\|) + \|B_j\|)} \\ &\geq \frac{\eta_1(1 - \theta_1) \theta_1 \alpha^2}{(1 + b_\sigma)(1 + \|B_j\|)} \\ &\geq \frac{\eta_1(1 - \theta_1) \theta_1 \alpha^2}{1 + b_\sigma} \frac{1}{r_j}, \end{aligned}$$

where  $r_j$  is defined in [Model Assumption 5.2](#). Let  $k \in \mathcal{S}$ . We sum the above inequalities over all  $j \in \mathcal{S}_k$ , and obtain

$$f(x_0) + h(x_0) - f(x_{k+1}) - h(x_{k+1}) \geq \frac{\eta_1(1-\theta_1)\theta_1\alpha^2}{1+b_\sigma} \sum_{j \in \mathcal{S}_k} \frac{1}{r_j}.$$

Since  $f + h$  is bounded below, it follows that  $\sum_{k \in \mathcal{S}} \frac{1}{r_k} < \infty$ . Let  $\tau \in \mathbb{N}_0$  be chosen so that  $\gamma_3\gamma_1^{\tau-1} > 1$ . By [Lemma 5.8](#),  $\sum_{j \in \mathcal{W}_k^\tau} \frac{1}{r_j}$  is uniformly bounded for all  $k \in \mathbb{N}$ . However, [Lemma 5.7](#) yields that for all  $k \in \mathbb{N}$ ,

$$\sum_{j=0}^k \frac{1}{r_j} \leq \tau \sum_{j \in \mathcal{S}_k} \frac{1}{r_j} + \sum_{j \in \mathcal{W}_k^\tau} \frac{1}{r_j},$$

which implies that  $\sum_{k=0}^{\infty} \frac{1}{r_k}$  converges, and contradicts [Model Assumption 5.2](#).  $\square$

Note that the assumptions involved in [Theorem 5.9](#) are weaker compared to existing methods in the literature—see [Section 1](#).

Moreover, our result is more general in that it establishes the existence of a subsequence  $(x_k)_{k \in \mathcal{K}}$  such that  $\lim_{k \in \mathcal{K}} \nu_k^{-1/2} \xi_{\text{cp}}(s_{k,\text{cp}}; x_k, \nu_k^{-1})^{1/2} = \lim_{k \in \mathcal{K}} \nu_k^{-1} \|s_{k,\text{cp}}\| = 0$ . That is in contrast with [Kanzow and Mehrlitz \[23\]](#), who prove that  $\lim_{k \in \mathcal{K}} \nu_k^{-1} \|s_{k,\text{cp}}\| = 0$  only for subsequences  $\mathcal{K}$  along which  $(x_k)_{k \in \mathcal{K}}$  converges. Furthermore, their analysis relies on the assumption that  $(\nu_k^{-1})_{k \in \mathbb{N}}$  is bounded away from zero, which is not required in our case.

Regarding the stationarity of accumulation points of [Algorithm 4.1](#), [[24](#), [Theorem 5](#)] directly implies that if there exists a subsequence  $\{x_k\}_{k \in \mathcal{K}} \rightarrow x^*$  such that  $\lim_{k \in \mathcal{K}} \nu_k^{-1/2} \xi_{\text{cp}}(s_{k,\text{cp}}; x_k, \nu_k^{-1})^{1/2} = 0$ , and such that there exists a limiting model  $\psi(\cdot; x^*)$  that satisfies [[24](#), [Model Assumption 3.1](#)], then  $x^*$  is a stationary point of [\(1.1\)](#), provided that the sequence  $(\sigma_k)_{k \in \mathbb{N}}$  is bounded from below.

The analysis of [[24](#)] does not establish whether all accumulation points of  $\{x_k\}_{k \in \mathbb{N}}$  are stationary for [\(1.1\)](#). A complete characterization of the stationarity of all accumulation points of [Algorithm 4.1](#) under similar or weaker assumptions, in particular the removal of the boundedness assumption on the sequence  $(\sigma_k)_{k \in \mathbb{N}}$ , is left for future work.

**6. Complexity analysis of [Algorithm 4.1](#).** In this section, we study the evaluation complexity of [Algorithm 4.1](#) in the case where the model Hessians are allowed to be unbounded. We replace [Model Assumption 5.1](#) with the following.

**MODEL ASSUMPTION 6.1.** *There is  $\kappa_m > 0$  such that for all  $k \in \mathbb{N}$ ,*

$$(6.1) \quad |(f+h)(x_k + s_k) - (\varphi + \psi)(s_k; x_k)| \leq \kappa_m(1 + \|B_k\|)\|s_k\|^2.$$

Note that if  $\nabla f$  is Lipschitz continuous and  $\psi(\cdot; x)$  satisfies [Model Assumption 5.1](#) then [Model Assumption 6.1](#) is satisfied as discussed by [Leconte and Orban \[24\]](#).

The next lemma will allow us to show that  $\{\sigma_k(1 + \|B_k\|^{-1})\}$  is bounded

**LEMMA 6.1.** *Let [Problem Assumption 4.1](#) and [Model Assumptions 3.1](#), [3.2](#) and [6.1](#) be satisfied. Define*

$$b_{\text{succ}} := \frac{2\kappa_m}{1 - \eta_2} > 0.$$

If  $x_k$  is not first-order stationary and  $\sigma_k(1 + \max_{0 \leq j \leq k} \|B_j\|)^{-1} \geq b_{\text{succ}}$ , iteration  $k$  is very successful and  $\sigma_{k+1} \leq \sigma_k$ .

*Proof.* Because  $x_k$  is not first-order stationary,  $s_k \neq 0$ . By definition of  $s_k$ ,  $m(0; x_k, \sigma_k) \geq m(s_k; x_k, \sigma_k)$ . Hence,

$$(6.2) \quad \varphi_k(0; x_k) + \psi_k(0; x_k) \geq \varphi_k(s_k; x_k) + \psi_k(s_k; x_k) + \frac{1}{2}\sigma_k \|s_k\|^2.$$

Thus, [Model Assumption 6.1](#) yields

$$\begin{aligned} |\rho_k - 1| &= \left| \frac{(f+h)(x_k + s_k) - (\varphi_k(s_k; x_k) + \psi_k(s_k; x_k))}{\varphi_k(0; x_k) + \psi_k(0; x_k) - (\varphi_k(s_k; x_k) + \psi_k(s_k; x_k))} \right| \\ &\leq \frac{\kappa_m(1 + \|B_k\|)\|s_k\|^2}{\frac{1}{2}\sigma_k \|s_k\|^2} \\ &\leq \frac{2\kappa_m}{\sigma_k(1 + \max_{0 \leq j \leq k} \|B_j\|)^{-1}} \leq \frac{2\kappa_m}{b_{\text{succ}}} = 1 - \eta_2. \end{aligned}$$

Thus, we obtain  $\rho_k \geq \eta_2$ , meaning that the iteration  $k$  is very successful.  $\square$

The next theorem shows that  $\{\sigma_k(1 + \|B_k\|^{-1})\}$  is bounded.

**THEOREM 6.2.** *Let [Problem Assumption 4.1](#) and [Model Assumptions 3.1, 3.2](#) and [6.1](#) be satisfied. For all  $k \in \mathbb{N}$ , if  $x_k$  is not stationary,*

$$\sigma_k(1 + \max_{0 \leq j \leq k} \|B_j\|)^{-1} \leq b_{\text{max}} := \min \{ \sigma_0(1 + \|B_0\|)^{-1}, \gamma_2 b_{\text{succ}} \} > 0.$$

*Proof.* Set  $b_k := \sigma_k(1 + \max_{0 \leq j \leq k} \|B_j\|)^{-1}$  for all  $k$ . We proceed by induction. For  $k = 0$ ,  $\sigma_0(1 + \|B_0\|)^{-1} \leq b_{\text{max}}$  by definition. Assume that  $b_k \leq b_{\text{max}}$  for  $k \geq 0$ .

Assume first that  $b_k < b_{\text{succ}}$ . Because  $\{(1 + \max_{0 \leq j \leq k} \|B_j\|)^{-1}\}$  is non-increasing, the update of  $\sigma_k$  in [Algorithm 4.1](#) ensures that

$$b_{k+1} = (1 + \max_{0 \leq j \leq k+1} \|B_j\|)^{-1} \sigma_{k+1} \leq (1 + \max_{0 \leq j \leq k} \|B_j\|)^{-1} \gamma_2 \sigma_k = \gamma_2 b_k < \gamma_2 b_{\text{succ}} \leq b_{\text{max}}.$$

Now, assume conversely that  $b_k \geq b_{\text{succ}}$ . [Lemma 6.1](#) implies that iteration  $k$  is very successful, and  $\sigma_{k+1} \leq \sigma_k$ . Thus,

$$b_{k+1} = (1 + \max_{0 \leq j \leq k+1} \|B_j\|)^{-1} \sigma_{k+1} < (1 + \max_{0 \leq j \leq k} \|B_j\|)^{-1} \sigma_k = b_k \leq b_{\text{max}}. \quad \square$$

Additionally, instead of [Model Assumption 5.2](#), we assume that model Hessians grow at most linearly with  $|\mathcal{S}_k|$ , which covers multiple quasi-Newton approximations—see [Section 1](#).

**MODEL ASSUMPTION 6.2.** *There are  $\mu > 0$  and  $0 \leq p \leq 1$  such that, for all  $k \in \mathbb{N}$ ,*

$$(6.3) \quad \max_{0 \leq j \leq k} \|B_j\| \leq \mu(1 + |\mathcal{S}_k|^p).$$

Because  $|\mathcal{S}_k|$  is non-decreasing with  $k$ , (6.3) is equivalent to  $\|B_k\| \leq \mu(1 + |\mathcal{S}_k|^p)$  for all  $k \in \mathbb{N}$ . The following theorem considers the case with a finite number of successful iterations. The proof follows [[4](#), Theorem 4.2] and is recalled here for completeness.

**THEOREM 6.3.** *Let [Problem Assumption 4.1](#) and [Model Assumptions 3.1, 3.2, 6.1](#) and [6.2](#) be satisfied. If [Algorithm 4.1](#) generates finitely many successful iterations, then  $x_k = x^*$  for all sufficiently large  $k$  where  $x^*$  is a stationary point.*

*Proof.* Assume by contradiction that  $x^*$  is not a stationary point. Because the number of successful iterations is finite, according to [Model Assumption 6.2](#), there is  $k_f \in \mathbb{N}$  such that  $\|B_k\| \leq \mu(1 + |\mathcal{S}_{k_f}|^p)$  for all  $k \geq k_f$ , where  $k_f$  is the index of the last successful iteration. The mechanism of [Algorithm 4.1](#) ensures that  $\sigma_k$  increases on unsuccessful iterations. Hence, there must exist an unsuccessful iteration  $k > k_f$  such that  $\sigma_k \geq b_{\text{succ}}(1 + \mu(1 + |\mathcal{S}_{k_f}|^p)) \geq b_{\text{succ}}(1 + \|B_k\|)$ , with  $b_{\text{succ}}$  defined in [Lemma 6.1](#). Because  $x^*$  is not stationary, we can apply [Lemma 6.1](#), which shows that  $k$  is very successful, and contradicts our assumption.  $\square$

We know from [Theorem 5.9](#) that  $\liminf_{k \rightarrow +\infty} \nu_k^{-1/2} \xi_{k,\text{cp}}^{1/2} = 0$  when [Algorithm 4.1](#) generates infinitely many successful iterations. Let  $\epsilon > 0$  and  $k_\epsilon$  be the first iteration of [Algorithm 4.1](#) such that  $\nu_k^{-1/2} \xi_{k,\text{cp}}^{1/2} \leq \epsilon$ . Define

$$(6.4a) \quad \mathcal{S}(\epsilon) := \mathcal{S}_{k_\epsilon-1} = \{k \in \mathcal{S} \mid k < k_\epsilon\},$$

$$(6.4b) \quad \mathcal{U}(\epsilon) := \mathcal{U}_{k_\epsilon-1} = \{k \in \mathbb{N} \mid k \notin \mathcal{S} \text{ and } k < k_\epsilon\}.$$

The next theorems bound  $k_\epsilon$ . The proofs are similar to [[18](#), Theorem 2].

**THEOREM 6.4.** *Let [Problem Assumption 4.1](#) and [Model Assumptions 3.1](#), [3.2](#), [6.1](#) and [6.2](#) be satisfied. Assume that [Algorithm 4.1](#) generates infinitely many successful iterations and that there is  $(f+h)_{\text{low}} \in \mathbb{R}$  such that  $(f+h)(x_k) \geq (f+h)_{\text{low}}$  for all  $k \in \mathbb{N}$ . If  $0 \leq p < 1$ ,*

$$(6.5) \quad |\mathcal{S}(\epsilon)| \leq ((1-p)\kappa_1\epsilon^{-2} + 1)^{1/(1-p)} - 1 = O(\epsilon^{-2/(1-p)}),$$

where

$$\kappa_1 = \frac{((f+h)(x_0) - (f+h)_{\text{low}})(b_{\text{max}} + 2\mu(1 + b_{\text{max}}))}{\eta_1\theta_1(1 - \theta_1)},$$

and  $b_{\text{max}}$  is as in [Theorem 6.2](#). If  $p = 1$ ,

$$(6.6) \quad |\mathcal{S}(\epsilon)| \leq \exp(\kappa_1\epsilon^{-2}) - 1.$$

*Proof.* Let  $k \in \mathcal{S}(\epsilon)$ , then  $\nu_k^{-1/2} \xi_{k,\text{cp}}^{1/2} \geq \epsilon$  and

$$(6.7) \quad (f+h)(x_k) - (f+h)(x_k + s_k) \geq \eta_1(1 - \theta_1)\xi_{\text{cp}}(x_k; \nu_k^{-1}) \geq \eta_1(1 - \theta_1)\nu_k\epsilon^2.$$

[Theorem 6.2](#) implies

$$\begin{aligned} \nu_k &= \frac{\theta_1}{\|B_k\| + \sigma_k} \geq \frac{\theta_1}{\max_{0 \leq j \leq k} \|B_j\| + b_{\text{max}}(1 + \max_{0 \leq j \leq k} \|B_j\|)} \\ &= \frac{\theta_1}{b_{\text{max}} + (1 + b_{\text{max}}) \max_{0 \leq j \leq k} \|B_j\|}. \end{aligned}$$

[Model Assumption 6.2](#) then implies

$$(6.8) \quad \nu_k \geq \frac{\theta_1}{b_{\text{max}} + \mu(1 + b_{\text{max}})(1 + |\mathcal{S}_k|^p)} = \frac{\theta_1}{|\mathcal{S}_k|^p} \zeta(|\mathcal{S}_k|^p),$$

where  $\zeta : \mathbb{R}_+ \rightarrow \mathbb{R}$ ,  $\zeta(x) := x/(b_{\text{max}} + \mu(1 + b_{\text{max}})(x + 1))$ .

Because  $\zeta$  is non-decreasing and  $|\mathcal{S}_k| \geq 1$  (as we have infinitely many successful iterations),  $\zeta(|\mathcal{S}_k|^p) \geq \zeta(1) = (1 + 2\mu(1 + b_{\max}))^{-1}$ . Thus, (6.8) becomes

$$\nu_k \geq \frac{\theta_1}{b_{\max} + 2\mu(1 + b_{\max})} \frac{1}{|\mathcal{S}_k|^p},$$

which combines with (6.7) to yield

$$(6.9) \quad (f + h)(x_k) - (f + h)(x_k + s_k) \geq \frac{\eta_1 \theta_1 (1 - \theta_1) \epsilon^2}{b_{\max} + 2\mu(1 + b_{\max})} \frac{1}{|\mathcal{S}_k|^p} := C \frac{1}{|\mathcal{S}_k|^p}.$$

We sum over all  $k \in \mathcal{S}(\epsilon)$ , and obtain

$$(f + h)(x_0) - (f + h)_{\text{low}} \geq C \sum_{k \in \mathcal{S}(\epsilon)} \frac{1}{|\mathcal{S}_k|^p} = C \sum_{k=0}^{|\mathcal{S}(\epsilon)|-1} \frac{1}{|\mathcal{S}_{\phi(k)}|^p},$$

where  $\phi$  is an increasing map from  $\{0, \dots, |\mathcal{S}(\epsilon)| - 1\}$  to  $\mathcal{S}(\epsilon)$ . Thus, by definition of  $\phi$  and  $\mathcal{S}_{\phi(k)}$ ,  $|\mathcal{S}_{\phi(k+1)}| = |\mathcal{S}_{\phi(k)}| + 1$  and  $|\mathcal{S}_{\phi(0)}| = 1$ . In other words,  $|\mathcal{S}_{\phi(k)}| = k + 1$ , and

$$(f + h)(x_0) - (f + h)_{\text{low}} \geq C \sum_{k=0}^{|\mathcal{S}(\epsilon)|-1} \frac{1}{(k+1)^p} = C \sum_{k=1}^{|\mathcal{S}(\epsilon)|} \frac{1}{k^p}.$$

Because  $\int_k^{k+1} \frac{1}{t^p} dt \leq \int_k^{k+1} \frac{1}{k^p} dt = \frac{1}{k^p}$ ,

$$(6.10) \quad (f + h)(x_0) - (f + h)_{\text{low}} \geq C \sum_{k=1}^{|\mathcal{S}(\epsilon)|} \int_k^{k+1} \frac{1}{t^p} dt = C \int_1^{|\mathcal{S}(\epsilon)|+1} \frac{1}{t^p} dt.$$

There are two cases to consider:

- if  $0 \leq p < 1$ ,  $(f + h)(x_0) - (f + h)_{\text{low}} \geq C \frac{(|\mathcal{S}(\epsilon)|+1)^{1-p}-1}{1-p}$ , which is (6.5);
- if  $p = 1$ ,  $(f + h)(x_0) - (f + h)_{\text{low}} \geq C \log(|\mathcal{S}(\epsilon)| + 1)$ , which is (6.6).  $\square$

Finally, we derive a bound on the cardinality of  $\mathcal{U}(\epsilon)$ .

**THEOREM 6.5.** *Let [Problem Assumption 4.1](#) and [Model Assumptions 3.1, 3.2](#) and [6.1](#) hold. Assume that [Algorithm 4.1](#) generates infinitely many successful iterations. Then*

$$(6.11) \quad |\mathcal{U}(\epsilon)| \leq |\log_{\gamma_1}(\gamma_3)| |\mathcal{S}(\epsilon)| + \log_{\gamma_1}(1 + \mu(1 + |\mathcal{S}(\epsilon)|^p)) + \frac{\log(b_{\max}/\sigma_0)}{\log(\gamma_1)},$$

where  $\mu$  and  $p$  are defined in [Model Assumption 5.2](#),  $b_{\max}$  as in [Theorem 6.2](#), and  $|\mathcal{S}(\epsilon)|$  is as in [Theorem 6.4](#).

*Proof.* The mechanism of [Algorithm 4.1](#) guarantees that for all  $k \in \mathbb{N}$ ,  $|\mathcal{U}_k| \leq |\log_{\gamma_1}(\gamma_3)| |\mathcal{S}_k| + \log_{\gamma_1}(\sigma_k/\sigma_0)$ . Hence, [Theorem 6.2](#) yields

$$\begin{aligned} |\mathcal{U}(\epsilon)| &\leq |\log_{\gamma_1}(\gamma_3)| |\mathcal{S}(\epsilon)| + \log_{\gamma_1} \left( \frac{b_{\max}(1 + \max_{0 \leq j \leq k_\epsilon - 1} \|B_j\|)}{\sigma_0} \right) \\ &\leq |\log_{\gamma_1}(\gamma_3)| |\mathcal{S}(\epsilon)| + \log_{\gamma_1}(1 + \mu(1 + |\mathcal{S}(\epsilon)|^p)) + \log_{\gamma_1} \left( \frac{b_{\max}}{\sigma_0} \right). \quad \square \end{aligned}$$

The complexity bound in [Theorem 6.4](#) is of the same order as that of [\[4, Lemma 4.3\]](#) for trust-region methods when  $p = 0$  in [Model Assumption 6.2](#), which corresponds to bounded model Hessians. Unlike [\[3, Lemma 3.6\]](#), the constant  $\theta_2$ , as defined in the switch on [Algorithm 4.1](#), does not appear in our complexity bound. Thus, large values of  $\theta_2$  in [Algorithm 4.1](#) will not worsen the complexity bound. In the general case where  $p > 0$ , our bound is better than that in [\[24, Theorem 4.2\]](#), as their step computation rule makes the bound dependent on  $\theta_2$ . As  $p$  approaches 1, the bound in [\[24, Theorem 4.2\]](#) goes to infinity, whereas ours, though exponential, remains finite, as in [\[18\]](#). Finally, the same example as in [\[18, §3.1\]](#) shows that our complexity bounds are also tight.

**7. Algorithmic refinements.** We describe a special case of [Algorithm 4.1](#) and an extension for which the convergence theory continues to hold, that we exploit in the numerical experiments of [Section 8](#), and that prove to be efficient in practice. As both refinements have already been studied by [Leconte and Orban \[25\]](#) in the context of their trust-region method, we keep our description to a minimum.

**7.1. Special case: diagonal model Hessians.** If we select  $B_k$  to be diagonal in [Algorithm 4.1](#), a specialized implementation emerges whenever  $h$  is separable and  $\psi(\cdot; x_k)$  is chosen to be separable at each iteration. For a number of choices of separable  $h$  that are of interest in applications, the step  $s_k$  may be computed analytically without requiring an iterative subproblem solver. We refer to this implementation as R2DH, where “DH” stands to *diagonal Hessians*. This section is modeled after [\[25, Section 4\]](#), to which we refer the reader for further information.

Diagonal quasi-Newton methods originate from [\[16, 21, 31\]](#). In order for a variational problem to possess a solution that defines a diagonal update, the classic secant equation is replaced with the *weak* secant equation  $s_k^T B_{k+1} s_k = s_k^T y_k$ , where  $y_k = \nabla f(x_{k+1}) - \nabla f(x_k)$ . A handful of diagonal updates have been proposed in the literature. The most efficient is probably the *spectral* update  $B_{k+1} = \tau_{k+1} I$ , where  $\tau_{k+1} := s_k^T y_k / s_k^T s_k$  is defined as in the spectral gradient method [\[8\]](#). Because  $B_k$  is a multiple of the identity,  $h$  and  $\psi(\cdot; x_k)$  need not be separable as the computation of  $s_k$  boils down to the evaluation of a proximal operator with step length  $1/\sqrt{\tau_k + \sigma_k}$ —see [Definition 2.2](#). [Zhu et al. \[38\]](#) derive an update akin to the well-known PSB formula that may be indefinite. We refer to it below as *PSB*. [Andrei \[1\]](#) derives an update based on a different variational problem that may also be indefinite. We refer to it below as *Andrei*. Additionally, we include a new diagonal variant inspired from the BFGS formula using a diagonal update. The main idea comes from applying [\[12, Lemma 5.1\]](#) to the last term of the BFGS update, i.e.,  $y_k y_k^T / s_k^T y_k$ , to obtain the diagonal update

$$D_{k+1} = \frac{\sum_{i=1}^n |(y_k)_i|}{s_k^T y_k} \text{diag}(|y_k|).$$

This update remains positive as long as  $s_k^T y_k > 0$ . We refer to this variant below as *DBFGS*. Although DBFGS does not always satisfy the secant equation, our numerical results demonstrate its competitiveness against other state-of-the-art diagonal-based methods. Note that the three updates (i.e., PSB, Andrei and DBFGS) generate  $B_k$  that is not a multiple of the identity, and hence  $h$  and  $\psi(\cdot; x_k)$  should be separable.

R2DH may act as standalone solver for [\(1.1\)](#) or as subproblem solver to compute  $s_k$  in [Algorithm 4.1](#). Our results in [Section 8](#) illustrate that, in both use cases, R2DH typically outperforms R2 [\[2, 3, Algorithm 6.1\]](#).

**7.2. Non-monotone variants.** Inspired by the success of the non-monotone spectral gradient method [8], [Leconte and Orban](#) [25, Section 6] explain how to modify an algorithm similar to [Algorithm 4.1](#) to incorporate a non-monotone strategy.

Let  $q \in \mathbb{N}$  be a given *memory parameter*. Define  $q_k = 1$  if  $q = 0$  and  $q_k := \min(k, q)$  if  $q > 0$ . Define also  $\mathcal{S}_{q_k}^+$  the set of the  $q_k$  *most recent* successful iterations. By convention, we set  $\mathcal{S}_0^+ = \{0\}$ . An iteration  $k$  now considers the objective value at each iteration in  $\mathcal{S}_{q_k}^+$ . Define

$$(7.1) \quad (f + h)_{\max, k} := \max\{(f + h)(x_j) \mid j \in \mathcal{S}_{q_k}^+\}.$$

[Algorithm 4.1](#) corresponds to  $q = 0$ . The non-monotone strategy consists in enforcing decrease with respect to  $(f + h)_{\max, k}$  instead of  $(f + h)(x_k)$ . In other words, we redefine

$$\rho_k := \frac{(f + h)_{\max, k} - (f + h)(x_k + s_k)}{(f + h)_{\max, k} - (\varphi + \psi)(s_k; x_k)}.$$

As in [25, Section 6], the new expression of  $\rho_k$  does not interfere with convergence properties or complexity bounds, except that it changes the constants in the latter. This is a positive result, especially in comparison to [23], where additional assumptions, including uniform continuity of  $f + h$ , are required to establish convergence of their non-monotone proximal-gradient method.

**8. Numerical experiments.** Our implementation of [Algorithm 4.1](#) and all solvers used in the experiments are available in the [RegularizedOptimization](#) Julia module [5, 19]. We compare the performance of R2N and its variants against PANOC [35]. By default, R2N uses an L-BFGS approximation with memory 5, as implemented in the [LinearOperators](#) Julia module [26], and uses parameters  $\theta_1 = (1 + \varepsilon_M^{1/5})^{-1} \approx 0.999$ ,  $\theta_2 = 1/\varepsilon_M \approx 10^{15}$ ,  $\eta_1 = \varepsilon_M^{1/4} \approx 10^{-4}$ ,  $\eta_2 = 0.9$ , and  $\sigma_0 = \varepsilon_M^{1/3} \approx 10^{-6}$ , where  $\varepsilon_M$  is the machine epsilon. The reason for defining values based on  $\varepsilon_M$  is that our code may be run in various floating-point arithmetics. Here, however, all tests are run in double precision. If iteration  $k$  of [Algorithm 4.1](#) is very successful,  $\sigma_{k+1} = \sigma_k/3$ ; if iteration  $k$  is unsuccessful,  $\sigma_{k+1} = 3\sigma_k$ . Otherwise,  $\sigma_{k+1} = \sigma_k$ .

All solvers use the same stopping criterion and terminate when

$$(8.1) \quad \nu_k^{-1} \|s_{k, \text{cp}}\| < \epsilon_a,$$

where  $\epsilon_a = \varepsilon_M^{3/10} \approx 2 \cdot 10^{-5}$  is an absolute tolerance, or exceed the budget of 5,000 iterations and 3,600 seconds of CPU time. This absolute stopping criterion (8.1) aligns with the one used by default in PANOC [35]. As discussed in [Section 3](#),  $\nu_k^{-1} \|s_{k, \text{cp}}\|$  can also be considered as a first-order stationarity measure for (1.1). [Theorem 5.9](#) and (3.6) imply that  $\liminf_{k \rightarrow +\infty} \nu_k^{-1} \|s_{k, \text{cp}}\| = 0$ . To solve the subproblem in Line 7 of [Algorithm 4.1](#), we use either R2 [3, Algorithm 6.1], or one of several R2DH variants (Spec, PSB, Andrei, or DBFGS) as described in [Section 7](#), as well as the non-monotone spectral R2DH (R2DH-Spec-NM) with memory 5. R2 initializes  $\nu_0 = 1.0$ . R2N and R2DH initialize  $\nu_0$  according to Line 5 of [Algorithm 4.1](#). The subproblem solvers terminate as soon as

$$\hat{\nu}_k^{-1/2} \hat{\xi}_{\text{cp}}(x_k + s, \hat{\nu}_k)^{1/2} \leq \begin{cases} 10^{-3} & \text{if } k = 0, \\ \min\left(\left(\nu_k^{-1} \xi_{\text{cp}}\right)^{3/4}, 10^{-3} \left(\nu_k^{-1} \xi_{\text{cp}}\right)^{1/2}\right) & \text{if } k > 0, \end{cases}$$

where  $\xi_{\text{cp}} = \xi_{\text{cp}}(x_k, \nu_k)$ ,  $\hat{\nu}_k$  and  $\hat{\xi}_{\text{cp}}$  are the step size and first-order stationarity measure related to the subproblem solver. Note that R2 and all the R2DH variants can also be used to solve (1.1) directly. PANOC [35] is run with all default parameters. All quasi-Newton approximations are initialized to the identity. In all experiments, we use  $\psi(s; x) := h(x + s)$ .

Our objective is to minimize the number of objective and gradient evaluations, as they are generally expensive to compute, while assuming that the proximal operators of common regularizers such that  $\ell_0$  and  $\ell_1$  norms are comparatively cheap to evaluate. We include also other test problems with the nuclear norm and the rank regularizers.

In our figures, we set  $(f + h)^*$  to the best value found by all the solvers. We plot  $\Delta(f + h)(x_k) = (f + h)(x_k) - (f + h)^*$  against the iterations to illustrate progress towards that best value. We also report the following solver statistics in tables: the final value of  $f$  at convergence; the final  $h/\lambda$ , where  $\lambda$  is a weight on the regularizer  $h$ ; the final stationarity measure  $\nu^{-1}\|s\|$ ; the number of evaluations of the smooth objective ( $\#f$ ); the number of evaluations of the gradient ( $\#\nabla f$ ); the number of proximal operator evaluations ( $\#\text{prox}$ ); and the elapsed time  $t$  in seconds.

**8.1. Basis pursuit denoise (BPDN).** The first set of experiments focuses on the basis pursuit denoise problem as described in [3], which is common in statistical and compressed sensing applications. The goal is to recover a sparse signal  $x_{\text{true}} \in \mathbb{R}^n$  from noisy observed data  $b \in \mathbb{R}^m$ . This problem can be formulated as

$$(8.2) \quad \underset{x}{\text{minimize}} \quad \frac{1}{2}\|Ax - b\|_2^2 + \lambda\|x\|_0,$$

where  $A$  is  $m \times n$  and randomly generated with orthonormal rows. We set  $m = 2,000$ ,  $n = 5,120$ , and  $b := Ax_{\text{true}} + \varepsilon$ , where  $\varepsilon \sim \mathcal{N}(0, 0.01)$ . The true signal  $x_{\text{true}}$  is a vector of zeros, except for 100 of its components. We set  $\lambda = 0.1\|A^T b\|_\infty$ . All algorithms start from the same randomly generated, hence non-sparse,  $x_0$ . For this problem, we compare R2 with R2DH variants (Spec, PSB, Andrei, and DBFGS). The objective is to find a good subsolver for R2N, as the subproblem in Line 7 of Algorithm 4.1 has a structure similar to (8.2).

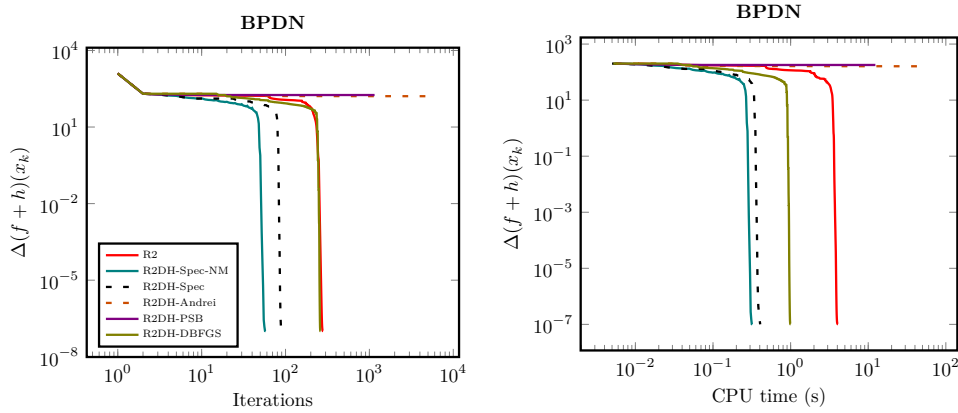


FIG. 8.1. BPDN objective vs. iterations (left) and CPU time (right).

Figure 8.1 shows that all solvers reach similar accuracy, except for the PSB and Andrei variants. R2DH-Spec-NM displays the best performance, followed by the R2DH-Spec and closely by R2DH-DBFGS variants, although it requires more evaluations to

TABLE 8.1  
Comparison of different solvers on the BPDN problem.

Solver	$f$	$h/\lambda$	$\Delta(f+h)$	$\nu^{-1}\ s\ $	$\#f$	$\#\nabla f$	$\#\text{prox}$	$t(s)$
R2	9.86e-02	100	4.26e-10	1.6e-05	281	273	280	4.08
R2DH-Spec-NM	9.86e-02	100	0.00e+00	3.8e-06	58	58	57	0.32
R2DH-Spec	9.86e-02	100	4.36e-10	1.8e-05	89	59	88	0.41
R2DH-Andrei	3.00e+00	2922	1.61e+02	2.2e+00	5001	4987	9991	55.49
R2DH-PSB	1.34e-09	3300	1.79e+02	2.0e-05	1153	1153	2311	12.36
R2DH-DBFGS	9.86e-02	100	1.24e-10	1.3e-05	262	153	261	1.00

achieve stationarity. Table 8.1 shows that all R2DH variants surpass R2 in all measures, except for R2DH-PSB and R2DH-Andrei, which either require more evaluations to attain the same level of accuracy or appear to converge to a different stationary point. R2 and all other R2DH variants identify a similarly-sparse solution. R2DH-Andrei requires significantly more evaluations and time than other R2DH variants and hits the iteration limit before (8.1) is triggered. Note that R2DH-DBFGS requires fewer function and gradient evaluations, as well as less time, than R2, but struggles to compete with R2DH-Spec-NM and R2DH-Spec. R2DH-Spec-NM is more efficient than R2DH-Spec, it avoids the unsuccessful iterations that R2DH-Spec falls into. Given its strong performance, in the following experiments we adopt R2DH-Spec-NM both as our default implementation of R2DH and as the R2DH subsolver.

**8.2. Matrix completion.** We address the matrix completion problem from [37] with rank nuclear norm regularizers to recover a low-rank matrix from noisy observations. The problem is formulated as

$$(8.3) \quad \underset{X}{\text{minimize}} \quad \frac{1}{2} \|P_{\Omega}(X - M)\|_F^2 + \lambda h(X),$$

where  $X \in \mathbb{R}^{n \times n}$  and  $n = 120$ . Here,  $\lambda = 10^{-1}$  is a weight, and  $h(X)$  is either  $\text{rank}(X)$  or  $\|X\|_*$ ,  $M$  is formed by applying a standard two-component Gaussian mixture model (GMM) to a low-rank matrix  $X_r$ . Specifically,  $M$  is computed as:

$$M = (1 - c)(X_r + \mathcal{N}(0, \sigma_A^2)) + c(X_r + \mathcal{N}(0, \sigma_B^2)),$$

where  $\mathcal{N}(0, \sigma_A^2)$  represents the noise component with variance  $\sigma_A^2$ , and  $\mathcal{N}(0, \sigma_B^2)$  represents the influence of outliers with a larger variance  $\sigma_B^2$ . The parameter  $c$  controls the relative proportion of noise and outliers in the observed matrix  $M$ . Finally,  $P_{\Omega}$  is a linear operator that extracts entries  $(i, j) \in \Omega$  and sets unobserved entries to zero, where

$$\Omega = \{(i, j) \in \{1, \dots, m\} \times \{1, \dots, n\} \mid R_{ij} < s_r\},$$

$R_{ij} \sim \mathcal{U}(0, 1)$ , and  $s_r$  is a threshold that determines the sparsity of the observed matrix.

For all solvers, we select a random initial matrix and set the rank of  $X_r$  to 40. The parameters are set to  $c = 0.2$ ,  $\sigma_A^2 = 0.0001$ ,  $\sigma_B^2 = 0.01$  and  $s_r = 0.8$ . Given that the smooth part of (8.3) is a linear least-squares residual, we apply the Levenberg-Marquardt (LM) algorithm from Aravkin et al. [4, Algorithm 4.1], which is a specific instance of R2N with  $B_k = J_k^T J_k$ , where  $J_k$  is the Jacobian of the least-squares residual at iteration  $k$ . Notably, R2DH can serve as a subproblem solver within LM—this combination is referred to as LM-R2DH, in contrast to the default LM-R2. We compare the performance of R2, R2DH, LM-R2, LM-R2DH and PANOC in Figure 8.2 and

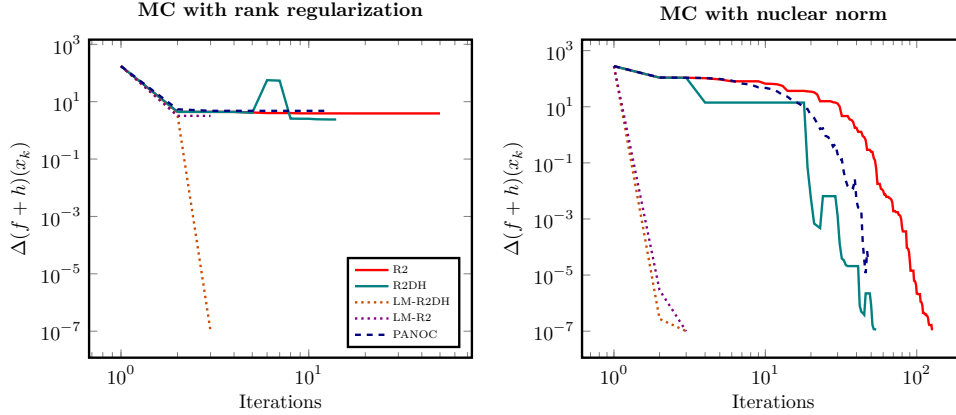


FIG. 8.2. Objectives vs. iterations for MC with rank (left) and nuclear norm (right) regularizers.

TABLE 8.2

Comparison of different solvers for matrix completion problem with rank regularizer.

Solver	$f$	$h/\lambda$	$\Delta(f+h)$	$\nu^{-1}\ s\ $	$\#f$	$\#J$	$\#\text{prox}$	$t(s)$
R2	1.61e-09	111	3.90e+00	1.7e-05	50	43	49	0.24
R2DH	1.09e-13	96	2.40e+00	4.7e-07	14	14	13	0.07
LM-R2DH	1.48e-11	72	0.00e+00	5.4e-06	3	283	103	0.59
LM-R2	2.67e-11	104	3.20e+00	1.6e-06	3	585	201	1.18
PANOC	1.07e-09	120	4.80e+00	5.1e-06	32	32	19	0.10

Tables 8.2 and 8.3. In Tables 8.2 and 8.3, the column  $\#\nabla f$  is replaced by the number of Jacobian or adjoint products  $\#J$ .

Figure 8.2 shows that LM-R2DH stands out in terms of final objective value for the rank regularizer followed by R2DH, while for the nuclear norm regularizer, all solvers achieve similar final objective value. Variants of LM require the fewest objective evaluations, although they demand many Jacobian-vector products, as seen in Tables 8.2 and 8.3. For the rank regularizer, Table 8.2 shows that, while R2DH requires more objective evaluations than either LM variant, it performs significantly fewer Jacobian-vector products. It is followed by PANOC, which requires more objective evaluations and Jacobian-vector products than R2DH. In terms of the objective value, LM-R2DH outperforms all the other solvers, but at the cost of additional Jacobian-vector products and proximal evaluations, and provides the solution with the best objective value, and, in particular, the lowest-rank solution. Finally, for the nuclear norm regularizer, LM variants behave almost identically according to Figure 8.2 and Table 8.3. R2DH outperforms R2 and PANOC in all measures, as shown in Table 8.3.

**8.3. General regularized problems.** In this section, we illustrate the performance of R2N on two test problems. The first problem addresses an image recognition task using a support vector machine (SVM) similar to those in [3]. The objective is to use this nonlinear SVM to classify digits from the MNIST dataset, specifically distinguishing between “1” and “7”, while excluding all other digits. A sparse support is imposed using an  $\ell_0$  regularizer. The optimization problem is given by

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{1} - \tanh(b \odot \langle A, x \rangle)\|^2 + \lambda \|x\|_0,$$

TABLE 8.3

Comparison of different solvers for matrix completion problem with nuclear norm regularizer.

Solver	$f$	$h/\lambda$	$\Delta(f+h)$	$\nu^{-1}\ s\ $	$\#f$	$\#J$	$\#\text{prox}$	$t(s)$
R2	1.00e-02	7.5e+00	1.54e-08	1.9e-05	128	79	127	0.65
R2DH	1.00e-02	7.5e+00	1.50e-08	7.0e-06	54	26	53	0.23
LM-R2DH	1.00e-02	7.5e+00	0.00e+00	3.2e-08	3	589	201	1.08
LM-R2	1.00e-02	7.5e+00	1.07e-13	5.9e-08	3	555	201	1.05
PANOC	1.00e-02	7.5e+00	2.83e-05	8.0e-06	103	103	55	0.41

where  $\lambda = 10^{-1}$  and  $A \in \mathbb{R}^{m \times n}$ , with  $n = 784$  representing the vectorized size of each image. In our tests, we use the training dataset, which includes  $m = 13,007$  images. Here,  $\odot$  denotes the elementwise product between vectors, and  $\mathbf{1} = (1, \dots, 1)$ .

The second problem is from [13, 34] and arises in image denoising and deblurring applications. The related optimization problem is given by

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \sum_{i=1}^n \log \left( (Ax - b)_i^2 + 1 \right) + \lambda \|x\|_1,$$

where  $\lambda = 10^{-4}$  and  $A \in \mathbb{R}^{n \times n}$  with  $n = 256^2$  is a Gaussian blur operator. The term  $b$  denotes the blurred image with added Gaussian noise. In our test,  $b$  is the blurred version of the cameraman image  $x^*$  with added Gaussian noise, i.e.,  $b = Ax^* + \text{noise}$ . The smooth part related to the two optimization problems is neither quadratic nor linear least squares, but a general non-convex problem.

We compare the performance of five methods: R2, R2DH, R2N-R2 (R2N with R2 as a subsolver), R2N-R2DH (R2N with R2DH as a subsolver) and PANOC [35].

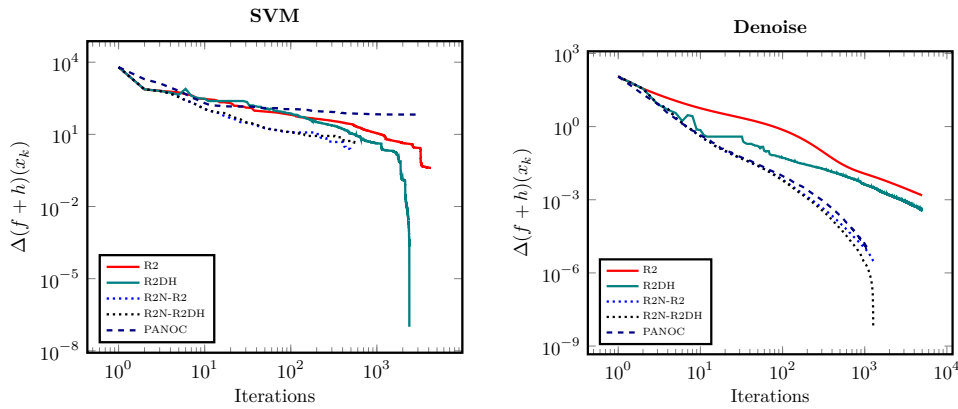


FIG. 8.3. Plots of the objective vs. iterations related to SVM (left) and denoise (right).

As shown in Tables 8.4 and 8.5, for both problems, the R2N variants outperform R2, R2DH and PANOC in terms of objective and gradient evaluations, though they require more proximal operator evaluations. Both R2N-R2DH and R2N-R2 have comparable performance and reach very good solutions compared to the other methods.

Note that for the non-linear SVM problem, as indicated in Figure 8.3 and Table 8.4, although R2DH reduces the objective function the most, it requires a higher number of evaluations of  $f$  and  $\nabla f$  than both R2N-R2DH and R2N-R2. PANOC is the least efficient method for this problem, as it requires the most evaluations and time while achieving the worst objective value.

TABLE 8.4  
Comparison of different solvers on the nonlinear SVM problem.

Solver	$f$	$h/\lambda$	$\Delta(f+h)$	$\nu^{-1}\ s\ $	$\#f$	$\#\nabla f$	$\#\text{prox}$	$t(s)$
R2	1.40e+01	122	4.02e−01	2.0e−05	4267	3303	4266	31.72
R2DH	1.40e+01	118	0.00e+00	2.0e−05	2384	1369	2383	15.12
R2N-R2	1.20e+01	161	2.29e+00	1.9e−05	513	313	51201	4.29
R2N-R2DH	1.60e+01	144	4.59e+00	2.0e−05	561	297	54975	4.49
PANOC	1.40e+01	784	6.66e+01	1.7e−05	7338	7338	4211	73.84

TABLE 8.5  
Comparison of different solvers on the denoising problem.

Solver	$f$	$h/\lambda$	$\Delta(f+h)$	$\nu^{-1}\ s\ $	$\#f$	$\#\nabla f$	$\#\text{prox}$	$t(s)$
R2	5.89e−02	3.6e+03	1.51e−03	9.8e−04	5001	4998	5000	66.86
R2DH	5.88e−02	3.6e+03	3.74e−04	5.5e−04	5001	2924	5000	51.25
R2N-R2	5.88e−02	3.6e+03	2.53e−06	1.9e−05	1327	1327	132601	325.60
R2N-R2DH	5.88e−02	3.6e+03	0.00e+00	2.0e−05	1269	1269	126139	301.78
PANOC	5.87e−02	3.6e+03	8.95e−06	1.9e−05	2192	2192	1104	36.39

For the denoising problem, the R2N variants require the fewest objective and gradient evaluations, but incur significantly more proximal-operator evaluations and longer runtime. These costs arise primarily from solving the subproblem in Line 7 of Algorithm 4.1. Moreover, since the dimension of the denoising problem is 65, 536, the cost of evaluating the proximal operator is not negligible. PANOC requires roughly twice as many function and gradient evaluations as the R2N variants, while using fewer proximal-operator evaluations and less time.

**8.4. FitzHugh-Nagumo inverse problem.** We consider an inverse problem for recovering the parameters of a nonlinear ordinary differential equation (ODE) model. Let  $x \in \mathbb{R}^p$  with  $p = 5$  denote the model parameters, and let  $F : \mathbb{R}^p \rightarrow \mathbb{R}^{2(N+1)}$  map  $x$  to the time series of state variables obtained by solving the FitzHugh–Nagumo (FH) neuron activation model [3, 4]

$$(8.4) \quad \frac{dV}{dt} = (V - \frac{1}{3}V^3 - W + x_1)x_2^{-1}, \quad \frac{dW}{dt} = x_2(x_3V - x_4W + x_5),$$

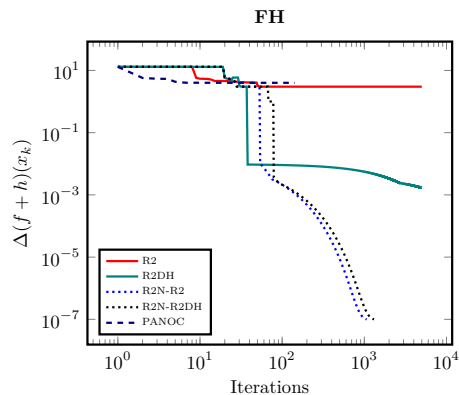
initialized at  $(V, W) = (2, 0)$  and integrated over  $t \in [0, 20]$  over a uniform grid of  $N$  steps. We take  $N = 1000$ , i.e.,  $\Delta t = 0.02s$ . We denote  $V(t_i; x)$  and  $W(t_i; x)$  the numerical solution at grid points  $\{t_i\}_{i=0}^N$ , and set  $F(x) = (v(x), w(x))$  with  $v(x) = (V(t_0; x), \dots, V(t_N; x))$  and  $w(x) = (W(t_0; x), \dots, W(t_N; x))$ . Synthetic observations are generated as  $b = F(x_{\text{true}}) + \mathcal{N}(0, 0.1^2 I)$ , using the sparse ground truth  $x_{\text{true}} = (0, 1, 0, 0, 0)$ . We estimate  $x$  using the sparsity-regularized problem

$$(8.5) \quad \underset{x \in \mathbb{R}^p}{\text{minimize}} \quad \frac{1}{2} \|F(x) - b\|_2^2 + \|x\|_0,$$

where  $F$  and its Jacobian are computed by solving and performing automatic differentiation on (8.4).

The FH problem is challenging due in particular to the fact that the gradient of the smooth term in (8.5) is not Lipschitz continuous.

For this problem, the R2N variants and R2DH recover a sparse solution, as shown in Tables 8.6 and 8.7. PANOC requires fewer objective, gradient, and proximal evaluations than R2N and R2DH, but struggles to converge to a sparse solution, is very slow, and requires the most time. R2N-R2 (closely followed by R2N-R2DH)

FIG. 8.4. *FH objective vs. iterations.*TABLE 8.6  
*Comparison of different solvers on the FH problem.*

Solver	$f$	$h/\lambda$	$\Delta(f+h)$	$\nu^{-1}\ s\ $	$\#f$	$\#\nabla f$	$\#\text{prox}$	$t(s)$
R2	1.17e+00	4	3.00e+00	3.5e-02	5001	3758	5000	5.67
R2DH	1.17e+00	1	1.69e-03	1.2e-02	5001	2494	5000	7.20
R2N-R2	1.17e+00	1	0.00e+00	2.0e-05	1067	1036	106601	1.38
R2N-R2DH	1.17e+00	1	7.26e-11	2.0e-05	1313	1270	125856	1.60
PANOC	1.16e+00	5	4.00e+00	1.9e-05	380	380	237	46.07

shows the best performance and successfully recovers a sparse solution. Although the R2N variants require more proximal operator evaluations than the other methods, the cost of evaluating the proximal operator is negligible, since the FH problem is low-dimensional with only  $p = 5$  parameters.

TABLE 8.7  
*Solution recovered by the different solvers on the FH problem.*

True	R2	R2DH	R2N-R2	R2N-R2DH	PANOC
0.00	0.00	0.00	0.00	0.00	2.51e-10
1.00	1.30	1.47	1.78	1.78	1.71
0.00	-0.25	0.00	0.00	0.00	-263.31
0.00	0.81	0.00	0.00	0.00	2181.53
0.00	0.43	0.00	0.00	0.00	463.48

**9. Discussion.** We proposed method R2N, a modified quasi-Newton method for nonsmooth regularized problems. R2N generalizes both R2 [3] and LM [4] and enjoys convergence properties without assuming Lipschitz continuity of  $\nabla f$  or boundedness of the model Hessians. Inspired by Diouane et al. [18], who work on trust-region methods for smooth optimization, we propose a complexity analysis of R2N to handle potentially unbounded model Hessians. Unlike traditional complexity analyses that assume uniformly bounded model Hessians, our study covers practical cases, including quasi-Newton updates such as PSB, BFGS, and SR1 by bounding the model Hessian growth with a power of the number of successful iterations—a reasonable bound as, in practice, it is uncommon to update quasi-Newton approximations on unsuccessful iterations. Nevertheless, Diouane et al. [18] show that similar complexity bounds continue to hold when the model Hessians are bounded by a power of the number of

iterations, and not just the number of successful iterations. Because their analysis uses similar arguments, their complexity bounds continue to hold for R2N.

Numerical illustrations show the strong potential of our implementation of R2N and some of its variants (both as a main solver and as a subproblem solver) compared to PANOC [35]. In particular, diagonal variants are competitive with, and often outperform, R2 when used as a subsolver inside R2N. One of the main advantages of R2N in practice is that proximal operators are easier to compute than in TR [3]. We illustrated that advantage by solving rank and nuclear norm-regularized problems. One way to further enhance the performance of R2N is to use more efficient subproblem solvers, for example, by generalizing those proposed in [6, 7, 22] for convex  $h$  and reducing the required number of proximal operator evaluations. Alternatively, in certain cases, the subproblems can be solved exactly, as in [17].

R2N convergence analysis arguments can be used to update and strengthen the existing convergence analysis of methods R2, TR, TRDH [25], LM and LMTR [4]. In follow-up research, we aim to identify the nature of limit points under our assumptions.

**Acknowledgement.** The authors would like to thank Maxence Gollier for his help with an improved implementation of R2N, as well as two anonymous reviewers for their constructive comments and suggestions that helped improve the paper.

#### REFERENCES

- [1] N. Andrei. [A diagonal quasi-Newton updating method for unconstrained optimization](#). *Numer. Algor.*, 81:575–590, 2019.
- [2] A. Aravkin, R. Baraldi, G. Leconte, and D. Orban. [Corrigendum: A proximal quasi-Newton trust-region method for nonsmooth regularized optimization](#). Cahier G-2021-12-SM, GERAD, Montréal, QC, Canada, 2024.
- [3] A. Y. Aravkin, R. Baraldi, and D. Orban. [A proximal quasi-Newton trust-region method for nonsmooth regularized optimization](#). *SIAM J. Optim.*, 32(2):900–929, 2022.
- [4] A. Y. Aravkin, R. Baraldi, and D. Orban. [A Levenberg–Marquardt method for nonsmooth regularized least squares](#). *SIAM J. Sci. Comput.*, 46(4):A2557–A2581, 2024.
- [5] R. Baraldi, G. Leconte, and D. Orban. [RegularizedOptimization.jl: Algorithms for regularized optimization](#), 2024.
- [6] S. Becker and J. Fadili. [A quasi-Newton proximal splitting method](#). In F. Pereira, C. Burges, L. Bottou, and K. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 25. Curran Associates, Inc., 2012.
- [7] S. Becker, J. Fadili, and P. Ochs. [On quasi-Newton forward-backward splitting: Proximal calculus and convergence](#). *SIAM J. Optim.*, 29(4):2445–2481, 2019.
- [8] E. G. Birgin, J. M. Martínez, and M. Raydan. [Spectral projected gradient methods: Review and perspectives](#). *J. Stat. Softw.*, 60(3):1–21, 2014.
- [9] J. Bolte, S. Sabach, and M. Teboulle. [Proximal alternating linearized minimization for nonconvex and nonsmooth problems](#). *Math. Program.*, 146:459–494, 2014.
- [10] R. I. Boț, E. R. Csetnek, and S. László. [An inertial forward–backward algorithm for the minimization of the sum of two nonconvex functions](#). *EURO J. Comput. Optim.*, 4(1):3–25, 2016.
- [11] C. Cartis, N. I. M. Gould, and Ph. L. Toint. [On the evaluation complexity of composite function minimization with applications to nonconvex nonlinear programming](#). *SIAM J. Optim.*, 21(4):1721–1739, 2011.
- [12] E. Chouzenoux, J.-C. Pesquet, and A. Repetti. [Variable metric forward–backward algorithm for minimizing the sum of a differentiable function and a convex function](#). *J. Optimiz. Theory App.*, 162(1):107–132, 2014.
- [13] E. Chouzenoux, S. Martin, and J.-C. Pesquet. [A local MM subspace method for solving constrained variational problems in image recovery](#). *J. Math. Imaging Vis.*, 65(2):253–276, 2023.
- [14] A. R. Conn, N. I. M. Gould, and Ph. L. Toint. [Trust-region methods](#). Number 1 in MOS-SIAM Series on Optimization. SIAM, Philadelphia, USA, 2000.
- [15] A. De Marchi. [Proximal gradient methods beyond monotony](#). *J. Nonsmooth Anal. Optim.*, 4 (Original research articles), 2023.

- [16] J. E. Dennis, Jr. and H. Wolkowicz. **Sizing and least-change secant methods**. *SIAM J. Numer. Anal.*, 30(5):1291–1314, 1993.
- [17] Y. Diouane, M. Gollier, and D. Orban. **A nonsmooth exact penalty method for equality-constrained optimization: complexity and implementation**. Cahier G-2024-65, GERAD, Montréal, Canada, 2024.
- [18] Y. Diouane, M. L. Habiboullah, and D. Orban. **Complexity of trust-region methods in the presence of unbounded Hessian approximations**. Cahier G-2024-43, GERAD, Montréal, Canada, 2024. To appear in *Math. Program.*
- [19] Y. Diouane, M. Gollier, M. L. Habiboullah, and D. Orban. **RegularizedOptimization.jl: A julia framework for regularized and nonsmooth optimization**. Cahier G-2025-75, GERAD, Montréal, QC, Canada, 2025.
- [20] M. Fukushima and H. Mine. **A generalized proximal point algorithm for certain non-convex minimization problems**. *Int. J. Syst. Sci.*, 12(8):989–1000, 1981.
- [21] J.-C. Gilbert and C. Lemaréchal. **Some numerical experiments with variable-storage quasi-Newton algorithms**. *Math. Program.*, 45:407–435, 1989.
- [22] C. Kanzow and T. Lechner. **Efficient regularized proximal quasi-Newton methods for large-scale nonconvex composite optimization problems**. *Pac. J. Optim.*, 20(3):537–568, 2024.
- [23] C. Kanzow and P. Mehlitz. **Convergence properties of monotone and nonmonotone proximal gradient methods revisited**. *J. Optimiz. Theory App.*, 195(2):624–646, 2022.
- [24] G. Leconte and D. Orban. **Complexity of trust-region methods with unbounded Hessian approximations for smooth and nonsmooth optimization**. Cahier G-2023-65, GERAD, Montréal, QC, Canada, 2023.
- [25] G. Leconte and D. Orban. **The indefinite proximal gradient method**. *Comput. Optim. Appl.*, 91(2):861–903, 2025.
- [26] G. Leconte, D. Orban, A. Soares Siqueira, and contributors. **LinearOperators.jl: Linear Operators for Julia**, 2023.
- [27] J. D. Lee, Y. Sun, and M. A. Saunders. **Proximal Newton-type methods for minimizing composite functions**. *SIAM J. Optim.*, 24(3):1420–1443, 2014.
- [28] H. Li and Z. Lin. **Accelerated proximal gradient methods for nonconvex programming**. In *Proceedings of the 28th International Conference on Neural Information Processing Systems - Volume 1, NIPS’15*, pages 379–387, Cambridge, MA, USA, 2015. MIT Press.
- [29] P.-L. Lions and B. Mercier. **Splitting algorithms for the sum of two nonlinear operators**. *SIAM J. Numer. Anal.*, 16(6):964–979, 1979.
- [30] R. Liu, S. Pan, Y. Wu, and X. Yang. **An inexact regularized proximal Newton method for nonconvex and nonsmooth optimization**. *Comput. Optim. Appl.*, 88(2):603–641, 2024.
- [31] J. L. Nazareth. **If quasi-Newton then why not quasi-Cauchy?** *SIAG/OPT Views-and-News*, 6: 11–14, 1995.
- [32] M. J. D. Powell. **On the convergence of a wide range of trust region methods for unconstrained optimization**. *IMA J. Numer. Anal.*, 30(1):289–301, 2010.
- [33] R. Rockafellar and R. J.-B. Wets. *Variational Analysis*. Springer Verlag, 1998.
- [34] L. Stella, A. Themelis, and P. Patrinos. **Forward-backward quasi-Newton methods for nonsmooth optimization problems**. *Comput. Optim. Appl.*, 67(3):443–487, 2017.
- [35] L. Stella, A. Themelis, P. Sotasakis, and P. Patrinos. **A simple and efficient algorithm for nonlinear model predictive control**. In *2017 IEEE 56th Annual Conference on Decision and Control (CDC)*, pages 1939–1944, 2017.
- [36] A. Themelis, L. Stella, and P. Patrinos. **Forward-backward envelope for the sum of two nonconvex functions: Further properties and nonmonotone line search algorithms**. *SIAM J. Optim.*, 28(3): 2274–2303, 2018.
- [37] Q. Yu and X. Zhang. **A smoothing proximal gradient algorithm for matrix rank minimization problem**. *Comput. Optim. Appl.*, pages 1–20, 2022.
- [38] M. Zhu, J. L. Nazareth, and H. Wolkowicz. **The quasi-Cauchy relation and diagonal updating**. *SIAM J. Optim.*, 9(4):1192–1204, 1999.