

A LEVENBERG-MARQUARDT METHOD FOR NONSMOOTH REGULARIZED LEAST SQUARES

ALEKSANDR Y. ARAVKIN*, ROBERT BARALDI†, AND DOMINIQUE ORBAN‡

Abstract. We develop a Levenberg-Marquardt method for minimizing the sum of a smooth nonlinear least-squares term $f(x) = \frac{1}{2}\|F(x)\|_2^2$ and a nonsmooth term h . Both f and h may be nonconvex. Steps are computed by minimizing the sum of a regularized linear least-squares model and a model of h using a first-order method such as the proximal gradient method. We establish global convergence to a first-order stationary point of both a trust-region and a regularization variant of the Levenberg-Marquardt method under the assumptions that F and its Jacobian are Lipschitz continuous and h is proper and lower semi-continuous. In the worst case, both methods perform $O(\epsilon^{-2})$ iterations to bring a measure of stationarity below $\epsilon \in (0, 1)$. We report numerical results on three examples: a group-lasso basis-pursuit denoise example, a nonlinear support vector machine, and parameter estimation in neuron firing. For those examples to be implementable, we describe in detail how to evaluate proximal operators for separable h and for the group lasso with trust-region constraint. In all cases, the Levenberg-Marquardt methods perform fewer outer iterations than a proximal-gradient method with adaptive step length and a quasi-Newton trust-region method, neither of which exploit the least-squares structure of the problem. Our results also highlight the need for more sophisticated subproblem solvers than simple first-order methods.

Key words. Regularized optimization, nonsmooth optimization, nonconvex optimization, nonlinear least squares, Levenberg-Marquardt method, proximal gradient method.

AMS subject classifications. 49J52, 65K10, 90C53, 90C56,

1. Introduction.

We consider the problem

$$(1.1) \quad \underset{x}{\text{minimize}} \quad f(x) + h(x), \quad f(x) = \frac{1}{2}\|F(x)\|_2^2,$$

where $F : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is continuously differentiable and $h : \mathbb{R}^n \rightarrow \mathbb{R}$ is proper and lower semi-continuous; we allow h to be nonsmooth and nonconvex. In practice, f is often a data-misfit term while h is a regularizer designed to promote desirable properties in the solution, such as sparsity. Numerous applications investigated in the nonsmooth regularized optimization literature actually have the structure (1.1), including basis pursuit denoising [14, 28], sparse factorization and dictionary learning [2], and sparse total least squares [30]. Yet nonsmooth numerical methods do not exploit the least-squares structure, nor accommodate general nonsmooth regularizers.

We describe two methods for (1.1): a quadratic regularization variant and trust-region variant inspired by the method of Levenberg [19] and Marquardt [21], denoted

*Department of Applied Mathematics, University of Washington, Seattle WA., USA. E-mail: saravkin@uw.edu.

†Optimization and Uncertainty Quantification, Sandia National Laboratories, P.O. Box 5800, Albuquerque, NM, 87125, USA. E-mail: rjbaral@sandia.gov. This research was sponsored by the Department of Energy Office of Science, Office of Advanced Scientific Computing Research's John von Neumann Fellowship. Sandia National Laboratories is a multimission laboratory managed and operated by National Technology and Engineering Solutions of Sandia, LLC., a wholly owned subsidiary of Honeywell International, Inc., for the U.S. Department of Energy's National Nuclear Security Administration under contract DE-NA0003525. This paper describes objective technical results and analysis. Any subjective views or opinions that might be expressed in the paper do not necessarily represent the views of the U.S. Department of Energy or the United States Government..

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

LM and LMTR respectively. Steps are computed by approximately minimizing simpler nonsmooth iteration-dependent Gauss-Newton-type models. Our algorithmic realizations utilize first-order methods, such as the proximal gradient method or the quadratic regularization method of Aravkin et al. [1], to solve the subproblems. The trust-region approach allows for any arbitrary trust-region norm, which, in practice, is influenced by nonconvex subproblem tractability. For both algorithms, we establish global convergence in terms of an optimality measure describing achievable decrease by a single proximal gradient step. Additionally, we derive a worst-case complexity bound of $\mathcal{O}(1/\epsilon^2)$ iterations to bring the stationarity measure below a tolerance of $\epsilon \in (0, 1)$ for LM and LMTR, i.e., the presence of a nonsmooth term in the objective yields a complexity bound of the same order as in the smooth case.

We provide implementation details and illustrate the performance of our methods on several numerical examples, including basis pursuit denoise with group-lasso regularization, nonlinear support vector machine with $\ell_{1/2}^{1/2}$ -norm regularization, and a sparse parameter estimation example taken from the Fitzhugh-Nagumo model of neuron firing. Our methods exhibit favorable performance under certain conditions with respect to previous work Aravkin et al. [1]. We additionally provide efficient, open-source software implementations of LM and LMTR as a package in the Julia language [3]. We find that exploiting the least-squares structure yields few LM and LMTR outer iterations, a well-known benefit in smooth optimization. The cost incurred is a large number of inner iterations, i.e, spent solving the subproblem. Thus, the results highlight the need for more sophisticated methods to minimize the sum of a linear least-squares term and a nonsmooth regularizer.

Related research. The present research is based on the framework laid out by Aravkin, Baraldi, and Orban [1]. The convergence and complexity of our trust-region Levenberg-Marquardt implementation follow directly from the general results of [1]. To the best of our knowledge, the trust-region literature does not explicitly cover the case of a nonlinear least-squares smooth objective with a nonsmooth regularizer other than a penalty term even though numerous applications exhibit that structure. See [13] for background and an extensive treatment.

A large portion of the literature focuses on h convex and/or globally Lipschitz continuous, e.g., Cartis et al. [11], Grapiglia et al. [17] and references therein. We do not attempt to give a comprehensive account of that literature here as we focus on significantly weaker assumptions. While many methods exist in the first-order literature, e.g., [12], few can effectively utilize any significant curvature information. Proximal Newton methods [18] require solutions to nontrivial proximal operators and positive semi-definiteness of the Hessian. The small number of references that allow both f and h to be nonconvex that we are aware of include: Li and Lin [20], who design accelerations of the proximal gradient method under the assumption that $f + h$ is coercive; Bolte et al. [8] who design an alternating method for cases where $h(x) = h_1(x_1) + h_2(x_2)$ and (x_1, x_2) is a partition of x ; Stella et al. [26] who propose a linesearch limited-memory BFGS method named PANOC; Themelis et al. [27] who propose a nonmonotone linesearch proximal quasi-Newton method named ZeroFPR based on the forward-backward envelope; and Boç et al. [9], who study a proximal method with momentum. The last three converge if $f + h$ satisfies the Kurdyka-Łojasiewicz (KL) assumption. Moreover, while all include (1.1) as a special case, few exploit any curvature information and none are specific to the least-squares structure. The algorithms presented here, like those of [1], require no such coercivity or KL assumptions.

Notation. We use $\|\cdot\|$ to represent a generic, but fixed, norm on \mathbb{R}^n or \mathbb{R}^m . The unit ball defined by that norm is \mathbb{B} , and $x + \Delta\mathbb{B}$ is the ball centered at x of radius $\Delta > 0$. For an integer $q \geq 1$, $\|\cdot\|_q$ is the ℓ_q -norm and \mathbb{B}_q is the unit ball in the ℓ_q -norm. If $A \subseteq \mathbb{R}^n$, $\chi(\cdot | A)$ is the indicator of A , i.e., the function whose value is 0 if $x \in A$ and $+\infty$ otherwise. Unless otherwise noted, if A is a matrix, $\|A\|$ denotes the spectral norm of A , i.e., its largest singular value. We use $J(x) : \mathbb{R}^n \rightarrow \mathbb{R}^{n \times m}$ to denote the Jacobian of F at x .

2. Background.

DEFINITION 2.1 (Limiting subdifferential). *Consider $\phi : \mathbb{R}^n \rightarrow \overline{\mathbb{R}}$ and $\bar{x} \in \mathbb{R}^n$ with $\phi(\bar{x}) < \infty$. We say that $v \in \mathbb{R}^n$ is a regular subgradient of ϕ at \bar{x} , and we write $v \in \hat{\partial}\phi(\bar{x})$ if*

$$\liminf_{x \rightarrow \bar{x}} \frac{\phi(x) - \phi(\bar{x}) - v^T(x - \bar{x})}{\|x - \bar{x}\|_2} \geq 0.$$

The set of regular subgradients is also called the Fréchet subdifferential. We say that v is a general subgradient of ϕ at \bar{x} , and we write $v \in \partial\phi(\bar{x})$, if there are sequences $\{x_k\}$ and $\{v_k\}$ such that

$$x_k \rightarrow \bar{x}, \quad \phi(x_k) \rightarrow \phi(\bar{x}), \quad v_k \in \hat{\partial}\phi(x_k) \text{ and } v_k \rightarrow v.$$

The set of general subgradients is called the limiting subdifferential.

PROPOSITION 2.2 (25, Theorem 10.1). *If $\phi : \mathbb{R}^n \rightarrow \overline{\mathbb{R}}$ is proper and has a local minimum at \bar{x} , then $0 \in \hat{\partial}\phi(\bar{x}) \subseteq \partial\phi(\bar{x})$. If ϕ is convex, the latter condition is also sufficient for \bar{x} to be a global minimum. If $\phi = f + h$ where f is continuously differentiable on a neighborhood of \bar{x} and h is finite at \bar{x} , then $\partial\phi(\bar{x}) = \nabla f(\bar{x}) + \partial h(\bar{x})$.*

If $0 \in \hat{\partial}\phi(\bar{x})$, we say that \bar{x} is *first-order stationary* for ϕ . Under our assumptions,

$$(2.1) \quad x \text{ is first-order stationary for (1.1)} \iff 0 \in J(x)^T F(x) + \partial h(x).$$

The proximal gradient method [16] applied to a regularized objective $f(x) + h(x)$ where f is differentiable is defined by the iteration

$$(2.2) \quad x_{k+1} \in \underset{\nu h}{\text{prox}}(x_k - \nu \nabla f(x_k)) \quad (k \geq 0),$$

where $\nu > 0$ is a steplength and the *proximal operator* is defined as

$$(2.3) \quad \underset{\nu h}{\text{prox}}(y) := \underset{u}{\text{argmin}} \frac{1}{2} \|u - y\|_2^2 + \nu h(u).$$

Without further assumptions on h , (2.3) is a set that may be empty, or contain one or more elements. The iteration (2.2) has the following descent property

LEMMA 2.3 (8, Lemma 2). *Let ∇f be Lipschitz continuous with Lipschitz constant $L \geq 0$, h be proper lower semi-continuous and $\inf h > -\infty$. Let $x_k \in \text{dom } h$, $0 < \nu < 1/L$, and x_{k+1} be defined according to (2.2). Then,*

$$(2.4) \quad (f + h)(x_{k+1}) \leq (f + h)(x_k) - \frac{1}{2}(\nu^{-1} - L) \|x_{k+1} - x_k\|_2^2.$$

117 **3. Linear Least Squares.** For fixed $\sigma \geq 0$ and $x \in \mathbb{R}^n$, define

118 (3.1a) $\varphi(s; x) := \frac{1}{2} \|J(x)s + F(x)\|_2^2,$

119 (3.1b) $\psi(s; x) \approx h(x + s) \quad \text{with} \quad \psi(0; x) = h(x),$

120 (3.1c) $m(s; x, \sigma) := \varphi(s; x) + \frac{1}{2} \sigma \|s\|_2^2 + \psi(s; x).$

122 Consider the parametric problem and its optimal set

123 (3.2a) $p(x, \sigma) := \min_s m(s; x, \sigma) \leq \varphi(0; x) + \psi(0; x) = f(x) + h(x)$

124 (3.2b) $P(x, \sigma) := \operatorname{argmin}_s m(s; x, \sigma).$

126 The form of (3.2) is representative of a Levenberg-Marquardt subproblem for (1.1) in
127 which f and h are modeled separately.

128 In particular, $\varphi(0; x) = f(x)$ and $\nabla_s \varphi(0; x) = \nabla f(x)$. We make the following
129 additional assumption.

130 **MODEL ASSUMPTION 3.1.** For any $x \in \mathbb{R}^n$, $\psi(\cdot; x)$ is proper, lsc and prox-bounded,
131 i.e., there exists $\lambda_x \in \mathbb{R}_+ \cup \{+\infty\}$ such that $\psi(\cdot; x) + \frac{1}{2} \lambda_x^{-1} \|\cdot\|_2^2$ is bounded below. In
132 addition, $\psi(0; x) = h(x)$, and $\partial \psi(0; x) = \partial h(x)$.

133 In **Model Assumption 3.1**, we assume that our choice of λ_x is the supremum of all
134 possible choices, and we refer to it as the *threshold of prox-boundedness* of $\psi(\cdot; x)$. In
135 particular, $\psi(\cdot; x)$ is bounded below if and only if $\lambda_x = +\infty$.

136 By **Proposition 2.2**, if $\sigma \geq \lambda_x^{-1}$,

137
$$s \in P(x, \sigma) \implies 0 \in \nabla \varphi(s; x) + \sigma s + \partial \psi(s; x).$$

138 We define

139 (3.3) $\xi(x, \sigma) := (f + h)(x) - p(x, \sigma).$

140 The following stationarity criterion follows directly from the definitions above.

141 **LEMMA 3.1.** Let **Model Assumption 3.1** be satisfied and $\sigma \geq \lambda_x^{-1}$. Then $\xi(x, \sigma) =$
142 $0 \iff 0 \in P(x, \sigma) \implies x$ is first-order stationary for (1.1). In addition, x is first-order
143 stationary for (1.1) if and only if $s = 0$ is first-order stationary for (3.1c).

144 *Proof.* Note first that $\xi(x, \sigma) = 0 \iff p(x, \sigma) = (f + h)(x) = \varphi(0; x) + \psi(0; x),$
145 which occurs if and only if $0 \in P(x, \sigma)$. **Proposition 2.2** then implies $0 \in \partial m(0; x, \sigma) =$
146 $\nabla \varphi(0; x) + \partial \psi(0; x)$ and is equivalent to (2.1). \square

147 The next result states some properties of (3.2).

148 **PROPOSITION 3.2.** Let **Model Assumption 3.1** be satisfied. $\operatorname{dom} p = \operatorname{dom} P =$
149 $\operatorname{dom} \psi \times \{\sigma \mid \sigma \geq \lambda_x^{-1}\}$. In addition, for any $x \in \mathbb{R}^n$,

- 150 1. $p(x, \cdot)$ is proper lsc and for each $\sigma > \lambda_x^{-1}$, $P(x, \sigma)$ is nonempty and compact;
- 151 2. if $\{\sigma_k\} \rightarrow \bar{\sigma} > \lambda_x^{-1}$ in such a way that $\{p(x, \sigma_k)\} \rightarrow p(x, \bar{\sigma})$, and for each $k,$
152 $s_k \in P(x, \sigma_k)$, then $\{s_k\}$ is bounded and all its limit points are in $P(x, \bar{\sigma})$;
- 153 3. $p(x, \cdot)$ is continuous at any $\bar{\sigma} > \lambda_x^{-1}$ and $\{p(x, \sigma_k)\} \rightarrow p(x, \bar{\sigma})$ holds in part 2
154 if $\bar{\sigma} > 0$.

155 *Proof.* Parts 1–2 follow from applying [25, Theorem 1.17] by noting that (3.1c) is
156 level-bounded in s locally uniformly in (x, σ) because $\psi(\cdot; x) + \frac{1}{2} \lambda_x^{-1} \|s\|_2^2$ is bounded
157 and $\varphi(s; x) + \frac{1}{2} (\sigma - \lambda_x^{-1}) \|s\|_2^2$ is level bounded in s locally uniform in (x, σ) . Part 3
158 also follows from [25, Theorem 1.17] by noting that (3.1c) is continuous in σ at any
159 $\bar{\sigma} > \lambda_x^{-1}$. \square

By [Proposition 3.2](#) part 3, $\xi(x, \cdot)$ is continuous at any $\bar{\sigma} > \lambda_x^{-1}$.

Although [\(3.1a\)](#) is a natural model of f about x , convergence properties may be stated in terms of the simpler first-order model

$$(3.4a) \quad \varphi_1(s; x) := f(x) + \nabla f(x)^T s = \frac{1}{2} \|F(x)\|_2^2 + (J(x)^T F(x))^T s,$$

$$(3.4b) \quad m_1(s; x, \sigma) := \varphi_1(s; x) + \frac{1}{2} \sigma \|s\|^2 + \psi(s; x).$$

The first step of the proximal gradient method [\(2.2\)](#) applied to the minimization of both $\varphi(s; x) + \psi(s; x)$ and $\varphi_1(s; x) + \psi(s; x)$ with steplength $\nu > 0$ is

$$(3.5) \quad \begin{aligned} s_1 &\in \underset{\nu\psi(\cdot; x)}{\text{prox}}(-\nu J(x)^T F(x)) \\ &= \underset{s}{\text{argmin}} \frac{1}{2} \|s + \nu J(x)^T F(x)\|_2^2 + \nu\psi(s; x) \\ &= \underset{s}{\text{argmin}} (J(x)^T F(x))^T s + \frac{1}{2} \nu^{-1} \|s\|_2^2 + \psi(s; x) \\ &= \underset{s}{\text{argmin}} m_1(s; x, \nu^{-1}). \end{aligned}$$

If $\nu^{-1} \geq \sigma$, then $m_1(s; x, \sigma) \leq m_1(s; x, \nu^{-1})$. Therefore, if s_1 results from [\(3.5\)](#), it also induces decrease in [\(3.4b\)](#).

In parallel to [Lemma 3.1](#) and [Proposition 3.2](#), we may define

$$(3.6a) \quad p_1(x, \sigma) := \min_s m_1(s; x, \sigma) \leq \varphi_1(0; x) + \psi(0; x) = f(x) + h(x)$$

$$(3.6b) \quad P_1(x, \sigma) := \underset{s}{\text{argmin}} m_1(s; x, \sigma),$$

$$(3.6c) \quad \xi_1(x, \sigma) := (f + h)(x) - p_1(x, \sigma) \geq 0,$$

and we have the following results, stating corresponding properties of p_1 and ξ_1 . The proofs replicate those in [Proposition 3.2](#) and [Lemma 3.3](#).

LEMMA 3.3. *Let [Model Assumption 3.1](#) be satisfied and $\sigma \geq \lambda_x^{-1}$. Then $\xi_1(x, \sigma) = 0 \iff 0 \in P_1(x, \sigma) \implies x$ is first-order stationary for [\(1.1\)](#). In addition, x is first-order stationary for [\(1.1\)](#) if and only if $s = 0$ is first-order stationary for [\(3.4b\)](#).*

PROPOSITION 3.4. *Let [Model Assumption 3.1](#) be satisfied. $\text{dom } p_1 = \text{dom } P_1 = \text{dom } \psi \times \{\sigma \mid \sigma \geq \lambda_x^{-1}\}$. In addition, for any $x \in \mathbb{R}^n$,*

1. $p_1(x, \cdot)$ is proper lsc and for each $\sigma > \lambda_x^{-1}$, $P_1(x, \sigma)$ is nonempty and compact;
2. if $\{\sigma_k\} \rightarrow \bar{\sigma} > \lambda_x^{-1}$ in such a way that $\{p_1(x, \sigma_k)\} \rightarrow p_1(x, \bar{\sigma})$, and for each k , $s_k \in P_1(x, \sigma_k)$, then $\{s_k\}$ is bounded and all its limit points are in $P_1(x, \bar{\sigma})$;
3. $p_1(x, \cdot)$ is continuous at any $\bar{\sigma} > \lambda_x^{-1}$ and $\{p_1(x, \sigma_k)\} \rightarrow p_1(x, \bar{\sigma})$ holds in part 2 if $\bar{\sigma} > 0$.

Because $L = 0$ for φ_1 , [Lemma 2.3](#) implies that the decrease achieved by s_1 is $(\varphi_1 + \psi)(s_1; x) \leq (\varphi_1 + \psi)(0; x) - \frac{1}{2} \nu^{-1} \|s_1\|^2$, which can be rearranged as

$$(3.7) \quad (f + h)(x) - (\varphi_1 + \psi)(s_1; x) \geq \frac{1}{2} \nu^{-1} \|s_1\|^2 \geq \frac{1}{2} \sigma \|s_1\|^2.$$

In the special case where $\psi = 0$, $s_1 = -\nu^{-1} \nabla f(x)$, so that [\(3.7\)](#) reduces to

$$\xi_1(x, \sigma) \geq \xi_1(x, \nu^{-1}) \geq f(x) - \varphi_1(s_1; x) \geq \frac{1}{2} \sigma \nu^{-1} \|\nabla f(x)\|^2 \geq \frac{1}{2} \sigma^2 \|\nabla f(x)\|^2,$$

which suggests that $\sigma^{-1} (\xi_1(x, \nu^{-1}))^{1/2}$ may be used as stationarity measure.

198

4. Nonlinear Least Squares.

199

200

201

202

203

204

205

4.1. A regularization approach. We first examine the formulation of the method of [Levenberg](#) and [Marquardt](#) in which the model (3.1c) is employed to compute a step. Specifically, consider [Algorithm 4.1](#). The step s_k is computed by approximately minimizing (3.1c) in stage 7 but the quality of the step is measured without taking the regularization term $\frac{1}{2}\sigma_k\|s_k\|^2$ into account in stage 8. The subproblem step s_k may be computed by continuing the iterations of the proximal gradient method initialized at $s_{k,1}$. This gives rise to one possible implementation of [Algorithm 4.1](#).

Algorithm 4.1 Nonsmooth regularized Levenberg-Marquardt method.

- 1: Choose constants $0 < \eta_1 \leq \eta_2 < 1$ and $0 < \gamma_3 \leq 1 < \gamma_1 \leq \gamma_2$.
- 2: Choose $x_0 \in \mathbb{R}^n$ where h is finite, $\sigma_0 > 0$, compute $F(x_0)$ and $h(x_0)$.
- 3: **for** $k = 0, 1, \dots$ **do**
- 4: Choose a steplength $\nu_k < 1/(\|J(x_k)\|^2 + \sigma_k)$.
- 5: Compute $s_{k,1}$ as defined in (3.5) and $\xi_1(x_k, \nu_k^{-1})$ as defined in (3.6c).
- 6: Define $m(s; x_k, \sigma_k)$ as in (3.1c).
- 7: Compute an approximate solution s_k of (3.2b).
- 8: Compute the ratio

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

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

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

11: **end for**

206

207

208

209

210

211

It may occur that $\sigma_k \leq \lambda_{x_k}^{-1}$. In such a case, $\psi(s_k; x_k) = -\infty$ so that the rules of extended arithmetic imply $\rho_k = 0$, whether $h(x_k + s_k) = +\infty$ or is finite. Thus s_k will be rejected at stage 9 and σ_{k+1} will be chosen larger than σ_k at stage 10. After a finite number of such increases, σ_k will exceed $\lambda_{x_k}^{-1}$ and a step with finite $\psi(s_k; x_k)$ will result.

Our main working assumption is the following.

212

213

214

PROBLEM ASSUMPTION 4.1. *The residual F and its Jacobian J are bounded and Lipschitz continuous on $\Omega := \{x \in \mathbb{R}^n \mid (f+h)(x) \leq (f+h)(x_0)\}$ and h is proper and lower semi-continuous.*

215

216

217

218

219

220

While [Problem Assumption 4.1](#) is a strong demand on all of \mathbb{R}^n and, in particular, rules out the case of linear least squares, it is a common assumption in the convergence analysis of the Levenberg-Marquardt method. If Ω is a compact set, then F is Lipschitz continuous on Ω if it is \mathcal{C}^1 on Ω , and J is Lipschitz continuous on Ω if F is \mathcal{C}^2 on Ω .

Under [Problem Assumption 4.1](#), ∇f is Lipschitz continuous on Ω , i.e., there exists $L > 0$ such that

221

$$(4.1) \quad |f(x+s) - (f(x) + \nabla f(x)^T s)| \leq \frac{1}{2}L\|s\|_2^2 \quad \text{for all } x, x+s \in \Omega.$$

We emphasize that in what follows, knowledge of L , or an estimate thereof, is not required. Our next assumption on the model is the following.

MODEL ASSUMPTION 4.1. *There exists a constant $\kappa_m > 0$ such that for all x and $s \in \mathbb{R}^n$, $|(f + h)(x + s) - (\varphi + \psi)(s; x)| \leq \kappa_m \|s\|^2$.*

Model Assumption 4.1 is essentially an assumption on the nonsmooth part ψ of the model. Indeed, (3.1a) and (4.1) combine to yield

$$\begin{aligned} |f(x + s) - \varphi(s; x)| &\leq |f(x + s) - (f(x) + \nabla f(x)^T s)| + \frac{1}{2} \|J(x)s\|^2 \\ &\leq \frac{1}{2} (L + \|J(x)\|^2) \|s\|^2. \end{aligned}$$

where we used the definition of $f(x)$, the identity $\nabla f(x) = J(x)^T F(x)$, and (4.1). Thus if J is bounded on Ω , we obtain

$$|f(x + s) - \varphi(s; x)| \leq \frac{1}{2} (L + \sup_{x \in \Omega} \|J(x)\|^2) \|s\|^2.$$

In particular, Model Assumption 4.1 is satisfied with $\kappa_m = \frac{1}{2} (L + \sup_{x \in \Omega} \|J(x)\|^2)$ if we select $\psi(s; x) := h(x + s)$.

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

MODEL ASSUMPTION 4.2. *There exists $\lambda > 0$ such that $\lambda_{x_k} \geq \lambda$ for all $k \in \mathbb{N}$.*

Model Assumption 4.2 is satisfied if h itself is prox-bounded and we select $\psi(s; x_k) := h(x_k + s)$ at each iteration.

Our first result ensures that σ_k is bounded above in Algorithm 4.1.

THEOREM 4.1. *Let Problem Assumption 4.1 and Model Assumptions 3.1, 4.1 and 4.2 be satisfied, and let*

$$(4.2) \quad \sigma_{\text{succ}} := \max(2\kappa_m / (1 - \eta_2), \lambda^{-1}) > 0.$$

If x_k is not first-order stationary and $\sigma_k \geq \sigma_{\text{succ}}$, then iteration k is very successful and $\sigma_{k+1} \leq \sigma_k$.

Proof. Let s_k be the step computed at iteration k of Algorithm 4.1. If $\sigma_k < \lambda_{x_k}^{-1}$, $\rho_k = 0$ as explained above, s_k is rejected and σ_k is increased. Hence, we assume that $\sigma_k \geq \lambda^{-1} \geq \lambda_{x_k}^{-1}$. Because x_k is not first-order stationary, $s_k \neq 0$. Because s_k is an approximate solution of (3.2b), we must have

$$\varphi(0; x_k) + \psi(0; x_k) \geq \varphi(s_k; x_k) + \frac{1}{2} \sigma_k \|s_k\|^2 + \psi(s_k; x_k)$$

and therefore,

$$(4.3) \quad \varphi(0; x_k) + \psi(0; x_k) - (\varphi(s_k; x_k) + \psi(s_k; x_k)) \geq \frac{1}{2} \sigma_k \|s_k\|^2.$$

Model Assumption 4.1 and (4.3) combine to yield

$$|\rho_k - 1| = \frac{|f(x_k + s_k) + h(x_k + s_k) - (\varphi(s_k; x_k) + \psi(s_k; x_k))|}{\varphi(0; x_k) + \psi(0; x_k) - (\varphi(s_k; x_k) + \psi(s_k; x_k))} \leq \frac{2\kappa_m \|s_k\|^2}{\sigma_k \|s_k\|^2}.$$

After simplifying by $\|s_k\|^2$, we obtain $\sigma_k \geq \sigma_{\text{succ}} \implies \rho_k \geq \eta_2$. \square

Note that [Theorem 4.1](#) does not explicitly include [Problem Assumption 4.1](#) in its assumptions, though it is likely to be required for [Model Assumption 4.1](#) to hold.

Interestingly, [Theorem 4.1](#) holds without assuming that the step s_k satisfies a sufficient decrease condition. Upon examination of the proof, the reason turns out to be that any step that results in simple decrease in $m(s; \sigma, x)$ results in sufficient decrease in $\varphi(\cdot; x) + \psi(\cdot; x)$, independently of the method used to compute s_k .

[Theorem 4.1](#) ensures existence of a constant $\sigma_{\max} > 0$ such that

$$(4.4) \quad \sigma_k \leq \sigma_{\max} := \min(\sigma_0, \gamma_2 \sigma_{\text{succ}}) > 0 \quad \text{for all } k \in \mathbb{N}.$$

Our next result concerns the situation where a finite number of successful iterations occur. The proof is almost identical to that of [[13](#), Theorem 6.4.4] and [[1](#), Theorem 3.5] and is omitted.

THEOREM 4.2. *Let [Problem Assumption 4.1](#) and [Model Assumptions 3.1](#) and [4.1](#) be satisfied. If [Algorithm 4.1](#) only generates finitely many successful iterations, then $x_k = x^*$ for all sufficiently large k and x^* is first-order critical.*

By [Rockafellar and Wets](#) [[25](#), Theorem 1.25], $p_1(x, \sigma)$ increases when σ increases, and thus, $\xi_1(x, \sigma)$ decreases when σ increases. Thus, it follows from (4.4) that

$$(4.5) \quad \xi_1(x_k, \sigma_k) \geq \xi_1(x_k, \sigma_{\max}) \quad \text{for all } k \in \mathbb{N}.$$

[Lemma 3.1](#), (4.5) and the remarks at the end of [section 3](#) suggest using $\xi_1(x_k, \sigma_{\max})^{\frac{1}{2}}$ as stationarity measure. Indeed, for given $\epsilon > 0$, $\xi_1(x_k, \sigma_{\max}) \leq \epsilon / \sigma_{\max} \implies \sigma_k \xi_1(x_k, \sigma_{\max}) \leq \epsilon$.

Because we must choose the steplength ν_k as in [Step 4](#) of [Algorithm 4.1](#), we compute $\xi_1(x_k, \nu_k^{-1})$ rather than $\xi_1(x_k, \sigma_k)$. Concretely, for given $0 < \theta < 1$, we set

$$(4.6) \quad \nu_k := \theta / (\|J_k\|^2 + \sigma_k).$$

Under [Problem Assumption 4.1](#), there exists $\kappa_J > 0$ such that $\|J(x)\| \leq \kappa_J$ for all $x \in \Omega$. Because [Algorithm 4.1](#) only generates $x_k \in \Omega$, the above and (4.4) yield

$$(4.7) \quad \nu_k \geq \theta / (\kappa_J^2 + \sigma_{\max}) := \nu_{\min} > 0 \quad \text{for all } k \in \mathbb{N}.$$

Therefore, $\nu_k^{-1} \leq \nu_{\min}^{-1}$ for all $k \geq 0$, and

$$(4.8) \quad \xi_1(x_k, \nu_k^{-1}) \geq \xi_1(x_k, \nu_{\min}^{-1}) \quad \text{for all } k \in \mathbb{N}.$$

For a stopping tolerance $\epsilon \in (0, 1)$, we seek to determine $k(\epsilon) \in \mathbb{N}$ such that

$$(4.9) \quad \xi_1(x_k, \nu_{\min}^{-1})^{\frac{1}{2}} > \epsilon \quad \text{for all } k < k(\epsilon) \quad \text{and} \quad \xi_1(x_{k(\epsilon)}, \nu_{\min}^{-1})^{\frac{1}{2}} \leq \epsilon.$$

Define the sets

$$(4.10a) \quad \mathcal{S} := \{k \in \mathbb{N} \mid \rho_k \geq \eta_1\},$$

$$(4.10b) \quad \mathcal{S}(\epsilon) := \{k \in \mathcal{S} \mid k < k(\epsilon)\},$$

$$(4.10c) \quad \mathcal{U}(\epsilon) := \{k \in \mathbb{N} \mid k \notin \mathcal{S} \text{ and } k < k(\epsilon)\}.$$

In order to conduct the complexity analysis, it is necessary to assume that the step computation at [stage 7](#) of [Algorithm 4.1](#) is related to $\xi_1(x_k, \sigma_k)$. We make the following assumption.

295 STEP ASSUMPTION 4.1. *There exists $\kappa_{\text{mdc}} \in (0, 1)$ such that s_k computed at*
 296 *stage 7 of Algorithm 4.1 satisfies*

$$297 \quad (4.11) \quad \varphi(0; x_k) + \psi(0; x_k) - (\varphi(s_k; x_k) + \psi(s_k; x_k)) \geq \kappa_{\text{mdc}} \xi_1(x_k, \nu_k^{-1}).$$

298 Step Assumption 4.1 is similar to sufficient decrease conditions used in trust-region
 299 methods—see [13]. Aravkin et al. [1] provide a concrete use of such condition in a
 300 trust-region method for nonsmooth regularized optimization. Clearly, the sufficient
 301 decrease assumption is satisfied after a single step of the proximal gradient method
 302 applied to (3.1c). Hence, it is also satisfied at a minimizer of (3.1c). Thus, in step 7
 303 of Algorithm 4.1, one strategy is to continue the proximal-gradient iterations until a
 304 stopping condition is attained.

305 The following results parallel those of Aravkin et al. [1], which are in turn inspired
 306 from those of Cartis et al. [11] and references therein.

307 LEMMA 4.3. *Let Problem Assumption 4.1 and Model Assumptions 3.1 and 4.1 be*
 308 *satisfied and s_k be computed according to Step Assumption 4.1, where ν_k is chosen*
 309 *according to (4.6). Assume there are infinitely many successful iterations and that*
 310 *$f(x) + h(x) \geq (f + h)_{\text{low}}$ for all $x \in \mathbb{R}^n$. Then, for all $\epsilon \in (0, 1)$,*

$$311 \quad (4.12) \quad |\mathcal{S}(\epsilon)| \leq \frac{(f + h)(x_0) - (f + h)_{\text{low}}}{\eta_1 \kappa_{\text{mdc}} \epsilon^2} = O(\epsilon^{-2}).$$

312 *Proof.* For $k \in \mathcal{S}(\epsilon)$, Step Assumption 4.1 and (4.8) imply

$$\begin{aligned} 313 \quad (f + h)(x_k) - (f + h)(x_k + s_k) &\geq \eta_1 (\varphi(0; x_k) + \psi(0; x_k) - (\varphi(s_k; x_k) + \psi(s_k; x_k))) \\ 314 &\geq \eta_1 \kappa_{\text{mdc}} \xi_1(x_k, \nu_k^{-1}) \\ 315 &\geq \eta_1 \kappa_{\text{mdc}} \xi_1(x_k, \nu_{\text{min}}^{-1}) \\ 316 &\geq \eta_1 \kappa_{\text{mdc}} \epsilon^2. \end{aligned}$$

318 The rest of the proof mirrors that of [1, Lemma 3.6]. \square

319 LEMMA 4.4. *Under the assumptions of Lemma 4.3,*

$$320 \quad (4.13) \quad |\mathcal{U}(\epsilon)| \leq \frac{\log(\sigma_{\text{max}}/\sigma_0)}{\log(\gamma_1)} + |\mathcal{S}(\epsilon)| \frac{|\log(\gamma_3)|}{\log(\gamma_1)} = O(\epsilon^{-2}).$$

321 *Proof.* For each $k \in \mathcal{U}(\epsilon)$, $\sigma_{k+1} \geq \gamma_1 \sigma_k$, while for each $k \in \mathcal{S}(\epsilon)$, $\sigma_{k+1} \geq \gamma_3 \sigma_k$.
 322 Thus if $k(\epsilon)$ is the iteration for which (4.9) occurs for the first time,

$$323 \quad \sigma_0 \gamma_1^{|\mathcal{U}(\epsilon)|} \gamma_3^{|\mathcal{S}(\epsilon)|} \leq \sigma_{k(\epsilon)-1} \leq \sigma_{\text{max}}.$$

324 Taking logarithms, we have

$$325 \quad |\mathcal{U}(\epsilon)| \log(\gamma_1) + |\mathcal{S}(\epsilon)| \log(\gamma_3) \leq \log(\sigma_{\text{max}}/\sigma_0).$$

326 Rearranging and recalling that $0 < \gamma_3 < 1$ yields (4.13). \square

327 Combining Lemmas 4.3 and 4.4 yields the overall iteration complexity bound.

328 THEOREM 4.5. *Under the assumptions of Lemma 4.3,*

$$329 \quad (4.14) \quad |\mathcal{S}(\epsilon)| + |\mathcal{U}(\epsilon)| = O(\epsilon^{-2}).$$

330 Stated differently, Theorem 4.5 ensures that either $(f + h)(x_k) \rightarrow -\infty$ or that
 331 $\liminf_{k \rightarrow \infty} \xi_1(x_k, \nu_{\text{min}}^{-1}) = 0$.

4.2. A trust-region approach. We now apply Algorithm 3.1 of Aravkin et al. [1] to (1.1). We assume that each $f_i : \mathbb{R}^n \rightarrow \mathbb{R}$ is \mathcal{C}^1 , so that their Problem Assumption 3.1 is satisfied. A natural model for f about x is the Gauss-Newton model (3.1a), which satisfies $\varphi(0; x) = f(x)$ and $\nabla_s \varphi(0; x) = \nabla f(x) = J(x)^T F(x)$. The model $\psi(s; x)$ of $h(x + s)$ is required to satisfy the same Model Assumption 4.1, which holds provided ∇f is Lipschitz continuous or each f_i is \mathcal{C}^2 with bounded Hessian. In Aravkin et al. [1, Algorithm 3.1], the first proximal gradient step s_1 is computed by solving

$$(4.15) \quad \begin{aligned} & \underset{s}{\text{minimize}} \quad \frac{1}{2} \|F(x)\|_2^2 + (J(x)^T F(x))^T s + \frac{1}{2} \nu^{-1} \|s\|^2 + \psi(s; x) \\ & \text{subject to} \quad \|s\| \leq \Delta, \end{aligned}$$

i.e.,

$$s_1 \in \underset{\nu\psi(\cdot; x) + \chi(\cdot | \Delta\mathbb{B})}{\text{prox}} \quad (-\nu J(x)^T F(x)),$$

where $0 < \nu < 1/(\|J(x)\|^2 + \alpha^{-1} \Delta^{-1})$ for a preset constant $\alpha > 0$. Subsequent steps continue the proximal gradient iterations to compute an approximate solution of

$$(4.16) \quad \underset{s}{\text{minimize}} \quad \frac{1}{2} \|J(x)s + F(x)\|_2^2 + \psi(s; x) \quad \text{subject to} \quad \|s\| \leq \min(\beta \|s_1\|, \Delta),$$

where $\beta \geq 1$. The above describes a trust-region variant of the method of Levenberg [19] and Marquardt [21] for regularized nonlinear least-squares problems. The assumption that $\psi(\cdot; x)$ is prox-bounded can be removed because $\psi(\cdot; x) + \chi(\cdot | \Delta\mathbb{B})$ is always bounded below, hence prox-bounded with $\lambda_x = \infty$. An approximate solution of (4.16) must satisfy Step Assumption 4.1 with $\xi_1(x, \sigma)$ replaced with

$$\hat{\xi}_1(\Delta; x, \nu) := f(x) + h(x) - \hat{p}_1(\Delta; x, \nu),$$

where $\hat{p}_1(\Delta; x, \nu)$ is the optimal value of (4.15).

Under the above assumptions, Aravkin et al. establish that the trust-region radius Δ never drops below the threshold

$$\Delta_{\min} := \min \left(\Delta_0, \hat{\gamma}_1 \frac{\kappa_{\text{mdc}}(1 - \eta_2)}{2\kappa_m \alpha \beta^2} \right),$$

where $\Delta_0 > 0$ is the initial trust-region radius, $\hat{\gamma}_1 \in (0, 1)$ is the fraction by which Δ is reduced on rejected steps, $\eta_2 \in (0, 1)$ is the threshold above which Δ is increased on accepted steps, and κ_{mdc} and κ_m play similar roles as the constants of the same name in Model Assumption 4.1 and Step Assumption 4.1.

Aravkin et al. use $\hat{\xi}_1(\Delta_{\min}; x, \nu)$ as stationarity measure. They show that for any $\epsilon \in (0, 1)$, the number of iterations necessary to achieve

$$\hat{\xi}_1(\Delta_{\min}; x, \nu)^{\frac{1}{2}} \leq \epsilon$$

is $O(\epsilon^{-2})$ provided that $f + h$ is bounded below. We refer the reader to [1] for complete details.

5. Proximal operators. In Algorithm 4.1 or the algorithm of Section 4.2, a typical model of the nonsmooth term h is $\psi(s; x) := h(x + s)$. If those algorithms are to use Aravkin et al.'s quadratic regularization method [1, Algorithm 6.1] to compute a step, the latter will in turn form a model of $\psi(\cdot; x)$ at each iteration. In order to simplify notation, let $\psi_k(s) := \psi(s; x_k) = h(x_k + s)$ be the model used at iteration k of Algorithm 4.1 or the algorithm of Section 4.2.

370 **5.1. General proximal operators.** In Algorithm 4.1, the nonsmooth term
 371 in the objective of the subproblem is $\psi_k(s)$. The typical model about s_j reduces to
 372 $\omega_j(t) = \psi_k(s_j + t) = h(x_k + s_j + t)$ and, instead of (5.2), the step computed is

$$373 \quad (5.1) \quad t_j \in \operatorname{argmin}_t \frac{1}{2}\nu^{-1}\|t - q\|^2 + h(x_k + s_j + t).$$

374 The same change of variable as above yields

$$375 \quad v_j \in \operatorname{argmin}_v \frac{1}{2}\nu^{-1}\|v - \bar{q}\|^2 + h(v) = \operatorname{prox}_{\nu h}(\bar{q}),$$

376 whether h is separable or not. Thus we obtain

$$377 \quad t_j \in \operatorname{prox}_{\nu h}(\bar{q}) - (x_k + s_j).$$

378 The nonsmooth term in the objective of the subproblem of the algorithm of
 379 Section 4.2 is $\psi_k(s) + \chi(s; \Delta_k)$. About iterate s_j of [1, Algorithm 6.1], the user supplies
 380 a model $\omega_j(t) := \omega(t; s_j) \approx \psi_k(s_j + t) + \chi(s_j + t \mid \Delta_k \mathbb{B})$, and the typical choice is
 381 $\omega_j(t) = \psi_k(s_j + t) + \chi(s_j + t \mid \Delta_k \mathbb{B}) = h(x_k + s_j + t) + \chi(s_j + t \mid \Delta_k \mathbb{B})$. The step
 382 computed is $t_j \in \operatorname{prox}_{\nu \omega_j}(q)$ for certain fixed $\nu > 0$ and $q \in \mathbb{R}^n$, i.e.,

$$383 \quad (5.2) \quad t_j \in \operatorname{argmin}_t \frac{1}{2}\nu^{-1}\|t - q\|^2 + h(x_k + s_j + t) + \chi(s_j + t \mid \Delta_k \mathbb{B}).$$

384 The change of variables $v := x_k + s_j + t$ allows us to rewrite (5.2) as

$$385 \quad (5.3) \quad v_j \in \operatorname{argmin}_v \frac{1}{2}\nu^{-1}\|v - \bar{q}\|^2 + h(v) + \chi(v - x_k \mid \Delta_k \mathbb{B}),$$

386 where $\bar{q} := x_k + s_j + q$, from which we recover $t_j = v_j - (x_k + s_j)$.

387 **5.2. Separable shifted proximal operators.** If h is separable and the trust
 388 region is defined by the ℓ_∞ -norm, the problem decomposes and the i -th component of
 389 v_j is

$$390 \quad (5.4) \quad \begin{aligned} v_{j,i} &\in \operatorname{argmin}_{v_i} \frac{1}{2}\nu^{-1}(v_i - \bar{q}_i)^2 + h_i(v_i) + \chi(v_i - x_{k,i} \mid [-\Delta_k, \Delta_k]) \\ &= \operatorname{argmin}_{v_i} \frac{1}{2}\nu^{-1}(v_i - \bar{q}_i)^2 + h_i(v_i) + \chi(v_i \mid [x_{k,i} - \Delta_k, x_{k,i} + \Delta_k]). \end{aligned}$$

391 Two situations may occur. In the first situation, $x_{k,i} - \Delta_k < v_{j,i} < x_{k,i} + \Delta_k$, so that
 392 $v_{j,i} \in \operatorname{prox}_{\nu h_i}(\bar{q}_i)$, i.e.,

$$393 \quad t_{j,i} \in \operatorname{prox}_{\nu h_i}(\bar{q}_i) - (x_{k,i} + s_{j,i}).$$

394 In the second situation, at least one unconstrained solution lies outside of $[x_{k,i} -$
 395 $\Delta_k, x_{k,i} + \Delta_k]$, so that constrained global minima of (5.4) are either one or both
 396 bounds, and/or unconstrained local minima that lie between the bounds.

397 When h is convex, the constrained solution is the feasible point nearest the unique
 398 unconstrained global solution, i.e.,

$$399 \quad v_{j,i} \in \operatorname{proj}_{[x_{k,i} - \Delta_k, x_{k,i} + \Delta_k]} (\operatorname{prox}_{\nu h_i}(\bar{q}_i)),$$

400 i.e.,

$$401 \quad t_{j,i} \in \operatorname{proj}_{[x_{k,i} - \Delta_k, x_{k,i} + \Delta_k]} (\operatorname{prox}_{\nu h_i}(\bar{q}_i)) - (x_{k,i} + s_{j,i}).$$

402 EXAMPLE 5.1 ($\ell_{1/2}^{1/2}$ pseudonorm). Consider $\psi(s) = \|s\|_{1/2}^{1/2} = \sum_j |s_j|^{1/2}$. When
 403 the trust-region bounds are inactive, *Cao et al. [10]* express the solution of (5.4) as

$$404 \quad v_{j,i} = \begin{cases} \frac{2}{3}|\bar{q}_i| \left(1 + \cos\left(\frac{2}{3}\pi - \frac{2}{3}\mu_\lambda(\bar{q}_i)\right)\right) & \bar{q}_i > p(\lambda) \\ 0 & |\bar{q}_i| \leq p(\lambda) \\ -\frac{2}{3}|\bar{q}_i| \left(1 + \cos\left(\frac{2}{3}\pi - \frac{2}{3}\mu_\lambda(\bar{q}_i)\right)\right) & \bar{q}_i < -p(\lambda) \end{cases}$$

405 where

$$406 \quad \mu_\lambda(\bar{q}_i) := \arccos\left(\frac{\lambda}{4} \left(\frac{|\bar{q}_i|}{3}\right)^{-3/2}\right), \quad p(\lambda) := \frac{54^{1/3}}{4}(2\lambda)^{2/3}.$$

407 When the trust-region constraint is active, *Cao et al. [10]* state that the above yields
 408 the inflection points of (5.4). We simply check the inflection points as well as the
 409 bounds. If the inflection points are within the bounds, we choose the minimum; if not,
 410 we select the minimum value of the cost function at the bounds.

411 **5.3. Nonseparable shifted proximal operators for convex h .** In this section
 412 we consider examples of nonseparable shifted proximal operators. The starting point
 413 is (5.3) where we assume that h is closed, proper, and convex. We rewrite

$$414 \quad \chi(v - x \mid \Delta\mathbb{B}) = \sup_z \langle v - x, z \rangle - \sigma_{\Delta\mathbb{B}}(z),$$

415 where we write x and Δ instead of x_k and Δ_k for simplicity, and where the support
 416 function

$$417 \quad \sigma_{\Delta\mathbb{B}}(z) := \sup_d \langle d, z \rangle + \chi(d \mid \Delta\mathbb{B}).$$

418 We substitute into (5.3) and obtain the saddle point problem

$$419 \quad (5.5) \quad \min_v \sup_z \frac{1}{2}\nu^{-1}\|v - \bar{q}\|^2 + h(v) + \langle v - x, z \rangle - \sigma_{\Delta\mathbb{B}}(z).$$

420 The objective of (5.5) is convex in v and concave in z . The saddle-point conditions
 421 can be written

$$422 \quad \begin{aligned} 0 &\in \nu^{-1}(v - \bar{q}) + \partial h(v) + z = \nu^{-1}(v - (\bar{q} - \nu z)) + \partial h(v) \\ 0 &\in v - x - \partial \sigma_{\Delta\mathbb{B}}(z). \end{aligned}$$

423 The first condition implies that $v \in \text{prox}_{\nu h}(\bar{q} - \nu z)$. By convexity of h , v is unique so
 424 that we are left with

$$425 \quad (5.6) \quad 0 \in v - x - \partial \sigma_{\Delta\mathbb{B}}(z), \quad \text{where} \quad \text{prox}_{\nu h}(\bar{q} - \nu z) = \{v\}.$$

426 **5.3.1. Special case: ℓ_2 -norm.** For $h(\cdot) := \lambda \|\cdot\|_2$,

$$427 \quad (5.7) \quad \text{prox}_{\nu\lambda\|\cdot\|_2}(y) = \begin{cases} 0 & \text{if } \|y\| \leq \nu\lambda \\ \left(1 - \frac{\nu\lambda}{\|y\|_2}\right)y & \text{if } \|y\| > \nu\lambda \end{cases}.$$

428 We now show how to solve (5.3) by converting (5.6) to a scalar root finding
 429 problem. For given z , let

$$430 \quad \zeta = \zeta(z) := \|\bar{q} - \nu z\|_2.$$

435 There are two possibilities.

436 **Case A:** If $\zeta \leq \nu\lambda$, (5.7) yields

$$437 \quad \text{prox}_{\nu\lambda\|\cdot\|_2}(\bar{q} - \nu z) = \{v\} = \{0\}.$$

438 The optimal value of (5.3) in this case is $\frac{1}{2}\nu^{-1}\|\bar{q}\|^2$.

439 **Case B:** If $\zeta > \nu\lambda$, (5.7) yields

$$440 \quad (5.8) \quad \text{prox}_{\nu\lambda\|\cdot\|_2}(\bar{q} - \nu z) = \{v\} = \left\{ \left(1 - \frac{\nu\lambda}{\zeta}\right) (\bar{q} - \nu z) \right\},$$

441 and (5.6) becomes

$$442 \quad 0 \in x - \left(1 - \frac{\nu\lambda}{\zeta}\right) (\bar{q} - \nu z) + \partial\sigma_{\Delta\mathbb{B}}(z)$$

$$443 \quad = (\zeta - \nu\lambda) \frac{\nu}{\zeta} \left(z - \left(\frac{1}{\nu} \bar{q} - \frac{\zeta}{\nu(\zeta - \nu\lambda)} x \right) \right) + \partial\sigma_{\Delta\mathbb{B}}(z),$$

$$444$$

445 which we interpret as

$$446 \quad (5.9) \quad z = z(\zeta) := \text{prox}_{\frac{\zeta}{\nu(\zeta - \nu\lambda)}\sigma_{\Delta\mathbb{B}}} \left(\frac{1}{\nu} \bar{q} - \frac{\zeta}{\nu(\zeta - \nu\lambda)} x \right).$$

447 Recall that [6, Theorem 6.46]

$$448 \quad (5.10) \quad \text{prox}_{\alpha\sigma_{\Delta\mathbb{B}}}(y) = y - \alpha \text{proj}_{\Delta\mathbb{B}}(\alpha^{-1}y), \quad (\alpha > 0).$$

449 Therefore, the projection into $\Delta\mathbb{B}$ must be computable. In our implementation, we
450 use $\mathbb{B} = \mathbb{B}_\infty$.

451 We may now search for ζ such that

$$452 \quad (5.11) \quad g(\zeta) := \zeta - \|\bar{q} - \nu z(\zeta)\|_2 = 0.$$

453 Because projections into convex sets are Lipschitz continuous, so is g over $(\nu\lambda, +\infty)$.

454 Since (5.3) is strongly convex, there is a unique solution, and so g has at most
455 one root such that $\zeta > \nu\lambda$. Any such root of g yields v given by (5.8) and $z(\zeta)$ given
456 by (5.9) that jointly satisfy (5.6). If g has no such root, the Case A must occur.

457 The combination of (5.9) and (5.10) yields

$$458 \quad (5.12) \quad \bar{q} - \nu z(\zeta) = \frac{\zeta}{\zeta - \nu\lambda} \left[x + \text{proj}_{\Delta\mathbb{B}} \left(\frac{\zeta - \nu\lambda}{\zeta} \bar{q} - x \right) \right].$$

459 As $\zeta \uparrow \infty$, $(\zeta - \nu\lambda)/\zeta \uparrow 1$, and by continuity, the term between square brackets
460 in (5.12) converges to $x + \text{proj}_{\Delta\mathbb{B}}(\bar{q} - x)$. Therefore, $\|\bar{q} - \nu z(\zeta)\|_2 \rightarrow \|x + \text{proj}_{\Delta\mathbb{B}}(\bar{q} - x)\|_2$
461 and for sufficiently large ζ , we must have $g(\zeta) > 0$.

462 To study $g(\zeta)$ as $\zeta \downarrow \nu\lambda$, we consider several mutually-exclusive cases.

463 1. If $x \notin \Delta\mathbb{B}$, then, $\text{proj}_{\Delta\mathbb{B}}(-x) \neq -x$. As $\zeta \downarrow \nu\lambda$, $(\zeta - \nu\lambda)/\zeta \downarrow 0$, and by
464 continuity, the term between square brackets converges to $x + \text{proj}_{\Delta\mathbb{B}}(-x) \neq 0$.
465 Therefore, $\|\bar{q} - \nu z(\zeta)\|_2 \rightarrow \infty$ and for sufficiently small ζ , we must have
466 $g(\zeta) < 0$.

467 2. Consider next the case where $x \in \text{int } \Delta\mathbb{B}$. For ζ sufficiently close to $\nu\lambda$,

468 (5.13)
$$\text{proj}_{\Delta\mathbb{B}}\left(\frac{\zeta - \nu\lambda}{\zeta}\bar{q} - x\right) = \frac{\zeta - \nu\lambda}{\zeta}\bar{q} - x,$$

469 and $\bar{q} - \nu z(\zeta) = \bar{q}$, i.e., $z(\zeta) = 0$. In this case,

470 (a) if $\|\bar{q}\|_2 > \nu\lambda$, then $g(\zeta) < 0$ for ζ close enough to $\nu\lambda$,

471 (b) if $\|\bar{q}\|_2 \leq \nu\lambda$, then $g(\zeta) > 0$ for all $\zeta > \nu\lambda$;

472 3. If $\|x\|_\infty = \Delta$ and $\text{proj}_{\Delta\mathbb{B}}(\bar{q} - x) = -x$, then $\text{proj}_{\Delta\mathbb{B}}(\alpha\bar{q} - x) = -x$ for any
473 $\alpha > 0$. In this case, the term between square brackets in (5.12) is always zero,
474 and $\bar{q} - \nu z(\zeta) = 0$. Thus for all $\zeta > \nu\lambda$, $g(\zeta) = \zeta > 0$.

475 4. If $\|x\|_\infty = \Delta$ but $\text{proj}_{\Delta\mathbb{B}}(\bar{q} - x) \neq -x$, there are two possible situations. Either
476 the ray $\alpha\bar{q} - x$ intersects $\text{int } \Delta\mathbb{B}$, or it does not. If it does, (5.13) occurs for
477 all ζ sufficiently close to $\nu\lambda$, $\bar{q} - \nu z(\zeta) = \bar{q}$, and cases 2a-2b apply. If it does
478 not, we have from Lipschitz continuity that

479
$$\left\|x + \text{proj}_{\Delta\mathbb{B}}\left(\frac{\zeta - \nu\lambda}{\zeta}\bar{q} - x\right)\right\|_2 = \left\|\text{proj}_{\Delta\mathbb{B}}\left(\frac{\zeta - \nu\lambda}{\zeta}\bar{q} - x\right) - \text{proj}_{\Delta\mathbb{B}}(-x)\right\|_2 \leq \frac{\zeta - \nu\lambda}{\zeta}\|\bar{q}\|_2.$$

480 Thus, $\|\bar{q} - \nu z(\zeta)\|_2 \leq \|\bar{q}\|_2$, and

481 (a) if $\|\bar{q}\|_2 > \nu\lambda$, then $g(\zeta) \geq \zeta - \|\bar{q}\|_2 > 0$ for $\zeta > \|\bar{q}\|_2$, and so there may
482 exist a root in $(\nu\lambda, \|\bar{q}\|_2]$. By (5.12), and the fact that $\|y\|_2 \leq \sqrt{n}\|y\|_\infty$
483 for all y , we also have

484
$$\|\bar{q} - \nu z(\zeta)\|_2 \leq \frac{\zeta}{\zeta - \nu\lambda} \left(\|x\|_2 + \left\|\text{proj}_{\Delta\mathbb{B}}\left(\frac{\zeta - \nu\lambda}{\zeta}\bar{q} - x\right)\right\|_2 \right) \leq \frac{(\|x\|_2 + \Delta\sqrt{n})\zeta}{\zeta - \nu\lambda},$$

485 so that $g(\zeta) > 0$ for $\zeta > \nu\lambda + 2\Delta\sqrt{n}$. Thus, the search interval may
486 potentially be reduced to $(\nu\lambda, \min(\nu\lambda + \|x\|_2 + \Delta\sqrt{n}, \|\bar{q}\|_2)]$.

487 (b) if $\|\bar{q}\| \leq \nu\lambda$, then $g(\zeta) > 0$ for all $\zeta > \nu\lambda$.

488 Thus, in cases 1 and 2a, a root is guaranteed to exist in $(\nu\lambda, +\infty)$ and can be
489 found by a bisection method. The upper bound may be found by observing that (5.12)
490 implies

491
$$\|\bar{q} - \nu z(\zeta)\| \leq \frac{\zeta}{\zeta - \nu}(\|x\| + \Delta),$$

492 so that

493
$$g(\zeta) = \zeta - \|\bar{q} - \nu z(\zeta)\| \geq \zeta - \frac{\zeta}{\zeta - \nu\lambda}(\|x\| + \Delta),$$

494 and $g(\zeta) > 0$ as soon as $\zeta > \|x\| + \Delta + \nu\lambda$.

495 In case 1, a lower bound follows by applying the reverse triangle inequality to (5.12):

496
$$\|\bar{q} - \nu z(\zeta)\| \geq \frac{\zeta}{\zeta - \nu\lambda}(\|x\| - \Delta),$$

497 so that $g(\zeta) < 0$ as soon as $\zeta < \nu\lambda + \|x\| - \Delta$.

498 In case 2a, the lower bound is simply $\|\bar{q}\|$.

499 In cases 2b, 3 and 4b, there can be no root in $(\nu\lambda, +\infty)$ and Case A must occur.

500 Only case 4a requires a root search, with or without sign change. If no root exists
501 in the search interval, Case A must occur.

502 **5.3.2. Special case: Group lasso.** The group lasso penalty is a sum of ℓ_2 -norms
 503 of subvectors:

$$504 \quad R_g(x) = \sum_i \|x_{[i]}\|_2,$$

505 where the $x_{[i]}$ partition x into non-overlapping groups. The proximal operator of R_g
 506 consists in applying (5.7) to each subvector:

$$507 \quad (5.14) \quad \text{prox}_{\lambda R_g}(z)_{[i]} = \left(1 - \frac{\lambda}{\|z_{[i]}\|_2}\right)_+ z_{[i]}.$$

508 Thus, the strategy of the previous section may be applied to each group.

509 **6. Implementation and numerical experiments.** Our implementation of
 510 Algorithm 3.1 of [1] and Algorithm 4.1 for (1.1) employs Aravkin, Baraldi, and
 511 Orban’s quadratic regularization method, named R2, to compute a step. R2 may be
 512 viewed as an implementation of the proximal gradient method with adaptive step
 513 size. The trust-region variant uses $\Delta_0 = 1$, terminates the outer iterations as soon as
 514 $\xi(\Delta_k; x_k, \nu_k)^{1/2} < \epsilon_a + \epsilon_r \xi_{1,0}^{1/2}$, where $\epsilon_a > 0$ and $\epsilon_r > 0$ are an absolute and a relative
 515 tolerance, and $\xi_{1,0}$ is the value of ξ_1 observed at the first iteration. A round of inner
 516 iterations terminates as soon as

$$517 \quad (6.1) \quad \hat{\xi}_1(x_k + s, \hat{\sigma}_k) \leq \begin{cases} 10^{-1} & \text{if } k = 0, \\ \max(\epsilon, \min(10^{-1}, \xi_1(x_k, \sigma_k)/10)) & \text{if } k > 0, \end{cases}$$

518 where $\hat{\sigma}_k$ and $\hat{\xi}_1$ are the regularization parameter and first-order stationarity measure
 519 used inside R2. In Algorithm 4.1, we use $\sigma_0 = 0.01$, and we terminate the outer
 520 iterations as soon as $\xi_1(x_k, \sigma_k)^{1/2} < \epsilon$ for a tolerance $\epsilon > 0$ because σ_{\max} is unknown.
 521 The inner iterations stop in the same manner as (6.1). All algorithms are implemented
 522 in the Julia language [7] version 1.8 as part of the RegularizedOptimization.jl package
 523 [3]. The shifted proximal operators are implemented in the ShiftedProximalOperators.jl
 524 package [5], while test problems are in the RegularizedProblems.jl package [4]. By
 525 contrast with the numerical results of Aravkin et al. [1], test cases are explicitly
 526 implemented as nonlinear least-squares problems, with access to the residual $F(x)$ and
 527 its Jacobian, and not simply the gradient of $f(x) := \frac{1}{2}\|F(x)\|_2^2$. Jacobian-vector and
 528 transposed-Jacobian-vector products are either implemented manually or computed
 529 via forward [24] and reverse [23] automatic differentiation, respectively.

530 We perform comparisons with R2 and with the quasi-Newton trust-region method
 531 of Aravkin et al. [1], named TR, and which does not exploit the structure of (1.1). The
 532 trust region is defined in ℓ_∞ -norm and the quadratic model uses a limited-memory SR1
 533 Hessian approximation with memory 5. In all experiments, we use $\psi(s; x) := h(x + s)$.

534 A direct comparison between the four methods is difficult because LM and LMTR
 535 do not utilize the same gradient; they instead take Jacobian-vector and transposed-
 536 Jacobian-vector products. To provide a meaningful comparison, in the tables below, we
 537 state: 1) the number of objective (or residual) evaluations; 2) the number of gradient
 538 evaluations (for R2 and TR) ; 3) the number of transposed-Jacobian-vector products
 539 (for LM and LMTR), listed under gradient evaluations; 4) the solve time in seconds. Our
 540 rationale is as follows. LM and LMTR pass a model to R2 whose objective evaluation
 541 requires one Jv , and whose gradient uses a Jv and a $J^T v$. Note however that the
 542 latter Jv can be cached and reused. Thus, R2 requires one Jv at each iteration, and
 543 additionally one $J^T v$ at each successful iteration.

In the figures, we plot descent as a function of residual/objective evaluations.

The summary of the numerical results below is that exploiting the least-squares structure results in a large reduction in outer iterations. However, solving the subproblem with a first-order method such as R2 consumes many $J^T v$. Our experiments thus highlight the need for more sophisticated subproblem solvers dedicated to (3.1c) and (4.16).

6.1. Group LASSO. In the group-LASSO problem, we observe noisy data from a linear system $b = Ax_T + \varepsilon$, where $A \in \mathbb{R}^{m \times n}$ has orthonormal rows, and x_T is segmented into g groups with every element in that group set to one of $\{-1, 0, 1\}$. The group-LASSO problem is given by

$$(6.2) \quad \min_x \frac{1}{2} \|Ax - b\|_2^2 + \lambda \|x\|_{1,2},$$

where $h(x) = \|x\|_{1,2} = \sum_{i=1}^g \|x_{[i]}\|_2$, i.e., the sum of the ℓ_2 -norm of the groups. The groups consisting of all zeros are labeled as “inactive”, whereas the groups set to ± 1 are “active”. We let $m = 512$, $n = 200$ and $\lambda = 10^{-2}$. We designate $g = 5$ such groups of possible 16 (each with 32 elements) to be “active”. The noise $\varepsilon \sim \mathcal{N}(0, 0.01)$. Thus (6.2) has the form (1.1), where $F(x) = Ax - b$. We set the absolute and relative exit tolerances to be 10^{-4} each. The number of subproblem iterations is capped at 100 for each outer iteration.

Figure 6.1 shows the solutions of each algorithm, and Table 6.1 reports the statistics. All algorithms arrive at approximately the same solution. R2 requires the most function evaluations whereas the others require about the same. Table 6.1 suggests that a tradeoff exists between the number of proximal operator evaluations and the number of gradient/Jacobian-vector evaluations. TR takes many proximal iterations, whereas LMTR and LM take far fewer. This tradeoff is further exemplified in the next test cases.

TABLE 6.1

Group-LASSO (6.2) statistics for R2, TR, LM, and LMTR, and $h(x) = \|x\|_{1,2}$. The $\#\nabla f$ is the number of $J^T v$ for LM and LMTR.

Alg	$f(x)$	$h(x)$	$(f + h)(x)$	$\ x - x_T\ _2$	$\# f$	$\# \nabla f$	$\# \text{prox}$	t (s)
R2	0.00	0.26	0.27	0.45	113	67	113	0.02
TR	0.00	0.26	0.27	0.47	17	17	339	2.56
LM	0.00	0.26	0.27	0.46	10	647	265	0.05
LMTR	0.00	0.26	0.27	0.46	5	327	130	0.98

We additionally plot descent history in Figure 6.4a. The plots are roughly similar, with the trust region methods TR and LMTR performing the best.

6.2. Nonlinear support vector machine. We now solve an image recognition problem of the form (1.1), where

$$(6.3) \quad F(x) = \mathbf{1} - \tanh(b \odot \langle A, x \rangle), \quad \mathbf{1} = [1, \dots, 1]^T,$$

$A \in \mathbb{R}^{m \times n}$, $n = 784$ is the vectorized image size, the number of images is $m = 13007$ in the training set and $m = 2163$ in the test set, and \odot denotes the elementwise product

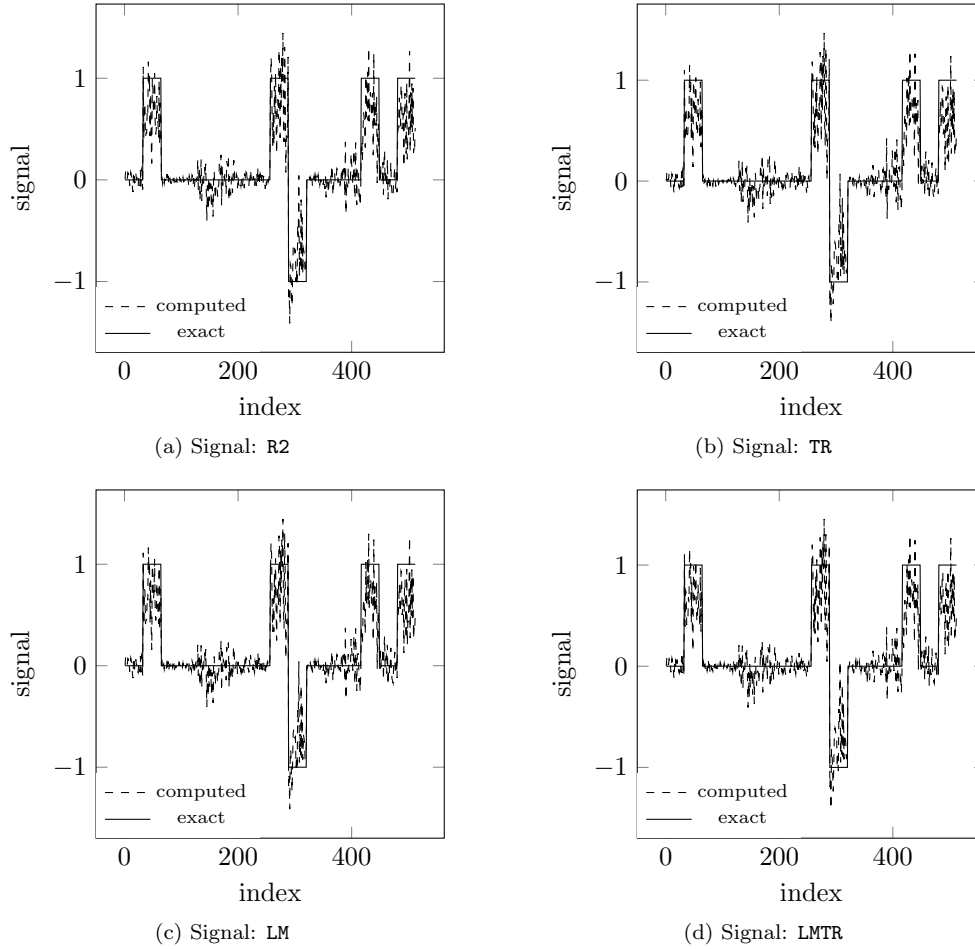


FIG. 6.1. *Group-LASSO* (6.2) solutions with R2, TR, LM, and LMTR with $h = \lambda \|\cdot\|_{1,2}$.

576 between vectors. We wish to use this nonlinear SVM to classify digits of the MNIST
 577 dataset as either 1 or 7, with all other digits removed. We additionally impose the
 578 condition that the support is sparse, and therefore use $h(x) = \|x\|_{1/2}^{1/2}$ as a regularizer.
 579 Hence, our overall problem is

$$(6.4) \quad \min_x \frac{1}{2} \|\mathbf{1} - \tanh(b \odot \langle A, x \rangle)\|^2 + \lambda \|x\|_{1/2}^{1/2}$$

581 with $\lambda = 10^{-1}$. We initialize the problem at $x = \mathbf{1}^n$ so that approximately 50% of the
 582 data is misclassified. We set the stopping tolerances again to 10^{-4} and the maximum
 583 number of inner iterations to 100.

584 **Figure 6.2** shows the solution map of each algorithm, which can be interpreted
 585 as the pixels most important in determining whether the image is indeed a 1 or 7.
 586 All algorithms produce a sparse solution; only about 8% of pixels in the support
 587 vector are nonzero. The problem is large and nonconvex; hence, the final solutions
 588 share pixels but altogether, they are different. This can be seen in **Table 6.3**, which
 589 reports the statistics. R2 again requires the most function evaluations. TR requires

590 about 10 times more than LM and LMTR. We again observe that a tradoff exists between
 591 number of proximal operator evaluations and the number of gradient/Jacobian-vector
 592 evaluations. Here, proximal operator evaluations are cheaper than gradient or Jv
 593 evaluations, so wallclock time is higher for LM and LMTR.

594 We plot descent history against number of function/residual iterations in [Figure 6.4b](#).
 595 Here we can see LM and LMTR performing the best in terms of descent.

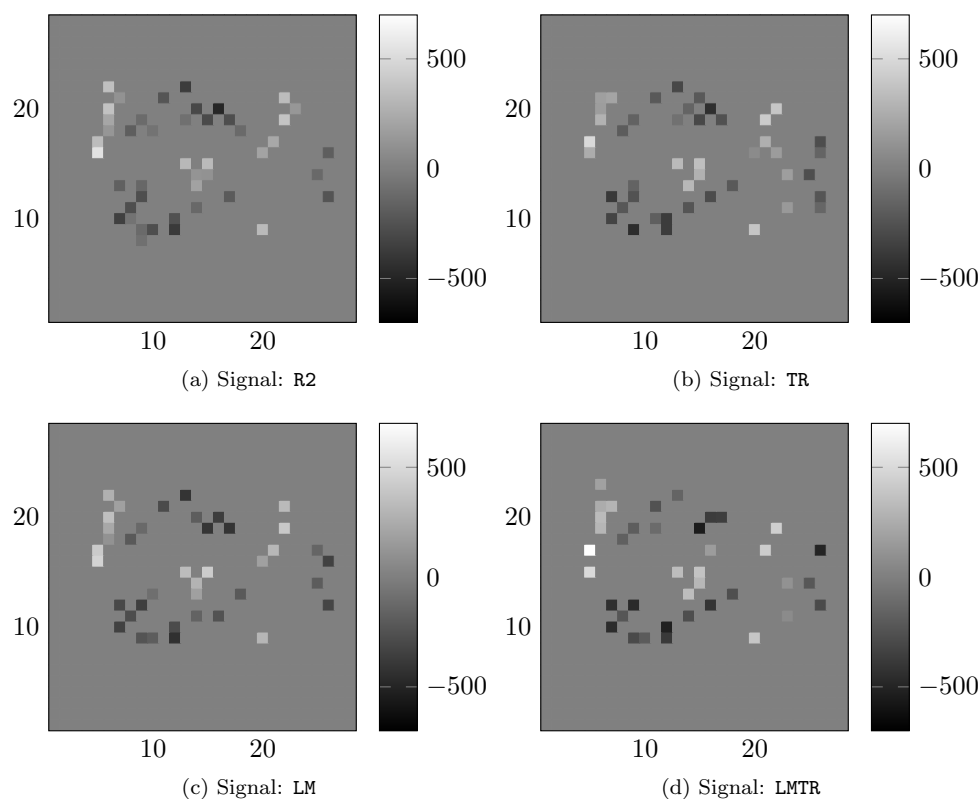


FIG. 6.2. *Nonlinear SVM (6.4) solutions with R2, TR, LM, LMTR.*

TABLE 6.3

Nonlinear SVM (6.4) statistics for R2, TR, LM, and LMTR. Training/test error is with respect to the ℓ_2 -norm.

Alg	f	h	$f + h$	(Train, Test)	# f	# ∇f	# prox	t (s)
R2	57.11	66.28	123.39	(99.8001, 99.3528)	1359	1085	1359	19.07
TR	50.28	70.70	120.99	(99.8309, 99.1216)	361	204	12749	9.21
LM	54.36	65.86	120.21	(99.8385, 99.3528)	23	3569	1276	34.60
LMTR	49.43	68.26	117.69	(99.8155, 99.1216)	24	3927	1420	42.56

596 **6.3. FitzHugh-Nagumo inverse problem.** The problem has the form (1.1),
 597 with $F : \mathbb{R}^5 \rightarrow \mathbb{R}^{2n+2}$ defined as $F(x) = (v(x) - \bar{v}(\bar{x}), w(x) - \bar{w}(\bar{x}))$, where $v(x) =$

598 $(v_1(x), \dots, v_{n+1}(x))$ and $w(x) = (w_1(x), \dots, w_{n+1}(x))$ are sampled values of discretized
 599 functions $V(t; x)$ and $W(t; x)$ satisfying the FitzHugh [15] and Nagumo et al. [22]
 600 model for neuron activation

$$601 \quad (6.5) \quad \frac{dV}{dt} = (V - V^3/3 - W + x_1)x_2^{-1}, \quad \frac{dW}{dt} = x_2(x_3V - x_4W + x_5),$$

602 parametrized by x . The sampling is defined by a discretization of the time interval
 603 $t \in [0, 20]$ and initial conditions $(V(0), W(0)) = (2, 0)$. The data $(\bar{v}(x), \bar{w}(x))$ is
 604 generated by solving (6.5) with $\bar{x} = (0, 0.2, 1, 0, 0)$, which corresponds to a simulation
 605 of the Van der Pol [29] oscillator. In our experiments, we use $n = 100$ and solve

$$606 \quad (6.6) \quad \min_x \frac{1}{2} \|F(x)\|_2^2 + \lambda \|x\|_1,$$

607 where $h(x) = \lambda \|x\|_1$ with $\lambda = 10$ to enforce sparsity in the parameters. Our absolute
 608 stopping criteria is 10^{-2} , whereas our the relative stopping criteria is set to 10^{-4} .

609 The solution found by each solver is given in Table 6.5 TR has the correct nonzero
 610 parameters, but the values are farther off. The corresponding simulations are shown
 611 in Figure 6.3; each method is able to fit the data.

TABLE 6.5
 Final parameters for the FH problem (6.6) found by R2, TR, LM, and LMTR.

True	R2	TR	LM	LMTR
0.00	0.00	0.00	0.00	0.00
0.20	0.26	0.33	0.25	0.25
1.00	0.84	0.70	0.86	0.85
0.00	0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.00	0.00

612 Table 6.7 reports the statistics for each algorithm, which exhibit the same pattern
 613 of results as before. The final objective values are fairly similar. LMTR uses the smallest
 614 amount of objective evaluations, whereas LM has a harder time solving (6.6). Because
 615 the gradient of the smooth term in (6.6) is not Lipschitz continuous, we had to set a
 616 σ_{\min} for both R2 and LM, which increased iteration count. Similar to the SVM example,
 617 we can see that LM and LMTR take more time than TR, which again stems from proximal
 618 operators being much cheaper to compute than Jv products for this example. Notably,
 619 TR seems to fit the data worse but attain a lower value of the regularizer.

620 Finally, Figure 6.4c shows descent of our objective function value against objective
 621 function iteration. LMTR again performs the best, whereas LM and TR were similar
 622 in this metric. This again enunciates the tradeoff between objective, gradient, and
 623 proximal operator expense. Expensive proximal evaluations would be the limiting
 624 factor in TR and R2; one can think of Total Variation regularization as a test case,
 625 since the proximal operator is itself a minimization problem.

626 **7. Discussion.** Similarly to smooth optimization, exploiting the least-squares
 627 structure of f can decrease significantly the number of outer iterations. The challenge
 628 highlighted by our numerical results, which is the subject of ongoing research, is
 629 to either identify a closed-form minimizer of (3.1c) for relevant choices of ψ , or to

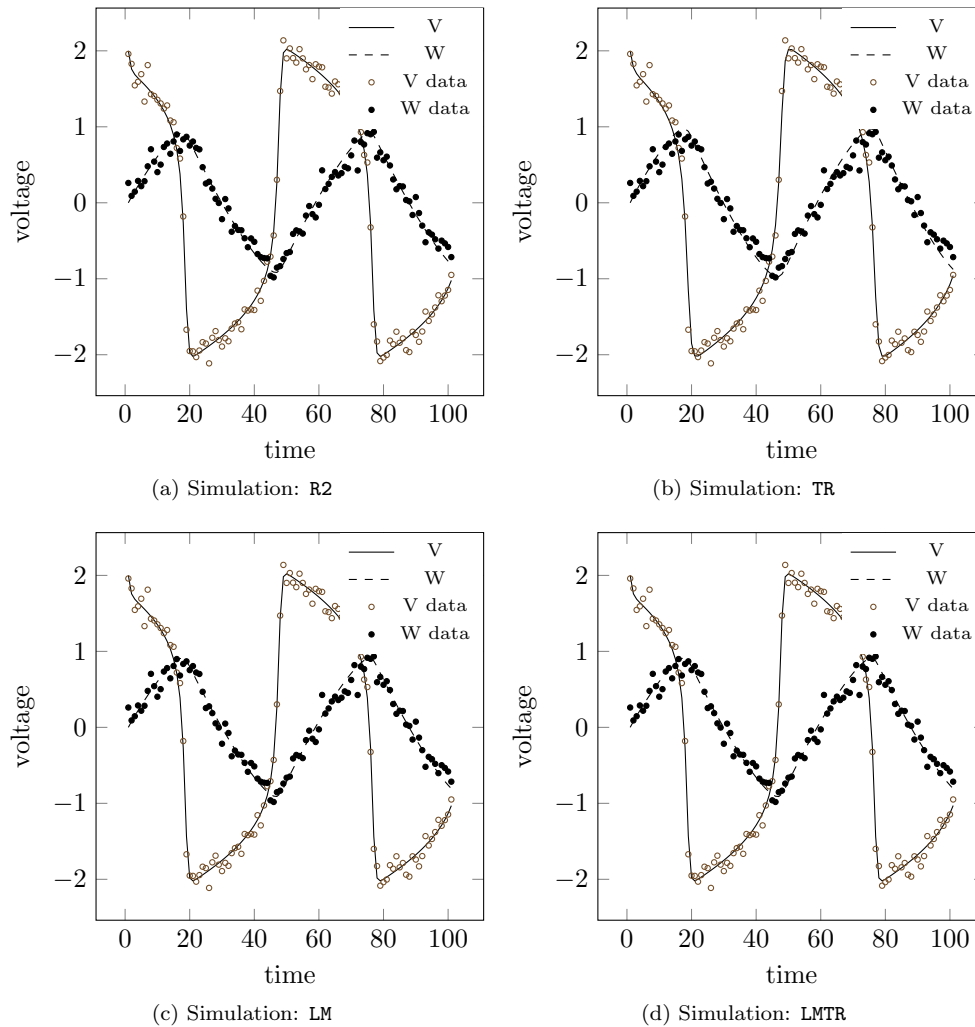


FIG. 6.3. Simulation of the FH problem (6.6) solutions found by R2, TR, LM, LMTR.

630 devise methods that can produce a higher-quality step than R2 with fewer transposed-
 631 Jacobian-vector products. As long as the subproblem solver yields a step satisfying
 632 [Step Assumption 4.1](#), our convergence properties and worst-case complexity bounds
 633 are guaranteed to hold. Thus, any improvement in the step computation mechanism
 634 will immediately translate into a more efficient solver overall. In ongoing research,
 635 we are exploring other improvements, including inexact evaluations of f and ∇f ,
 636 nonmonotone methods, and inexact evaluation of proximal operators.

637

REFERENCES

- 638 [1] A. Aravkin, R. Baraldi, and D. Orban. [A proximal quasi-Newton trust-region method for](#)
 639 [nonsmooth regularized optimization](#). *SIAM J. Optim.*, (2):900–929, 2022.
 640 [2] F. Bach, R. Jenatton, J. Mairal, and G. Obozinski. *Optimization with Sparsity-Inducing*
 641 *Penalties*, volume 4 of *Foundations and Trends in Machine Learning*. now publishers, 2012.

TABLE 6.7
 Statistics for the FH problem (6.6) for R2, TR, LM, and LMTR.

Alg	f	h	$f + h$	$\ x - x_T\ _2$	# f	# ∇f	# prox	t (s)
R2	1.24	10.91	12.15	1.58	4230	3428	4230	40.40
TR	1.87	10.31	12.17	1.93	134	77	2452	0.67
LM	1.20	11.03	12.23	1.55	101	4236	1402	20.17
LMTR	1.20	11.02	12.22	1.55	32	2006	741	10.50

- 642 [3] R. Baraldi and D. Orban. RegularizedOptimization.jl: Algorithms for regularized optimization.
 643 <https://github.com/JuliaSmoothOptimizers/RegularizedOptimization.jl>, February 2022.
- 644 [4] R. Baraldi and D. Orban. RegularizedProblems.jl: Test cases for regularized optimization.
 645 <https://github.com/JuliaSmoothOptimizers/RegularizedProblems.jl>, February 2022.
- 646 [5] R. Baraldi and D. Orban. ShiftedProximalOperators.jl: Proximal operators for regularized opti-
 647 mization. <https://github.com/JuliaSmoothOptimizers/ShiftedProximalOperators.jl>, February
 648 2022.
- 649 [6] A. Beck. *First Order Methods in Optimization*. SIAM, Philadelphia, USA, 2017.
- 650 [7] J. Bezanson, A. Edelman, S. Karpinski, and V. B. Shah. Julia: A fresh approach to numerical
 651 computing. *SIAM Rev.*, 59(1):65–98, 2017.
- 652 [8] J. Bolte, S. Sabach, and M. Teboulle. Proximal alternating linearized minimization for nonconvex
 653 and nonsmooth problems. *Math. Program.*, (146):459–494, 2014.
- 654 [9] R. I. Boł, E. R. Csetnek, and S. László. An inertial forward–backward algorithm for the
 655 minimization of the sum of two nonconvex functions. *EURO J. Comput. Optim.*, (4):3–25, 2016.
- 656 [10] W. Cao, J. Sun, and Z. Xu. Fast image deconvolution using closed-form thresholding formulas
 657 of l_q ($q = 12, 23$) regularization. *Journal on visual communication and image representation*,
 658 24(1), 2013.
- 659 [11] C. Cartis, N. I. M. Gould, and Ph. L. Toint. On the evaluation complexity of composite function
 660 minimization with applications to nonconvex nonlinear programming. *SIAM J. Optim.*, 21(4):
 661 1721–1739, 2011.
- 662 [12] P. L. Combettes and J.-C. Pesquet. Proximal splitting methods in signal processing. In *Fixed-*
 663 *point algorithms for inverse problems in science and engineering*, pages 185–212. Springer,
 664 2011.
- 665 [13] A. R. Conn, N. I. M. Gould, and Ph. L. Toint. *Trust-Region Methods*. Number 1 in MOS-SIAM
 666 Series on Optimization. SIAM, Philadelphia, USA, 2000.
- 667 [14] D. L. Donoho. Compressed sensing. *IEEE T. Inform. Theory*, 52(4):1289–1306, 2006.
- 668 [15] R. FitzHugh. Mathematical models of threshold phenomena in the nerve membrane. *B. Math.*
 669 *Biophys.*, 17(4):257–278, 1955.
- 670 [16] M. Fukushima and H. Mine. A generalized proximal point algorithm for certain non-convex
 671 minimization problems. *Int. J. Syst. Sci.*, 12(8):989–1000, 1981.
- 672 [17] G. Grapiglia, J. Yuan, and Y. Yuan. Nonlinear stepsize control algorithms: Complexity bounds
 673 for first- and second-order optimality. *J. Optim. Theory and Applics.*, (171):980–997, 2016.
- 674 [18] J. D. Lee, Y. Sun, and M. A. Saunders. Proximal Newton-type methods for minimizing composite
 675 functions. *SIAM J. Optim.*, 24(3):1420–1443, 2014.
- 676 [19] K. Levenberg. A method for the solution of certain problems in least squares. *Q. Appl. Math.*,
 677 (2):164–168, 1944.
- 678 [20] H. Li and Z. Lin. Accelerated proximal gradient methods for nonconvex programming. In
 679 *Proceedings of the 28th International Conference on Neural Information Processing Systems -*
 680 *Volume 1, NIPS’15*, pages 379–387, Cambridge, MA, USA, 2015. MIT Press.
- 681 [21] D. W. Marquardt. An algorithm for least-squares estimation of nonlinear parameters. *Journal*
 682 *of the Society for Industrial and Applied Mathematics*, 11(2):431–441, 1963.
- 683 [22] J. Nagumo, S. Arimoto, and S. Yoshizawa. An active pulse transmission line simulating nerve
 684 axon. *Proceedings of the IRE*, 50(10):2061–2070, 1962.
- 685 [23] J. Revels. Reverse mode automatic differentiation for Julia. [https://github.com/JuliaDiff/](https://github.com/JuliaDiff/ReverseDiff.jl)
 686 [ReverseDiff.jl](https://github.com/JuliaDiff/ReverseDiff.jl), 2022.
- 687 [24] J. Revels, M. Lubin, and T. Papamarkou. Forward-mode automatic differentiation in Julia,
 688 2016. <https://arxiv.org/abs/1607.07892>.

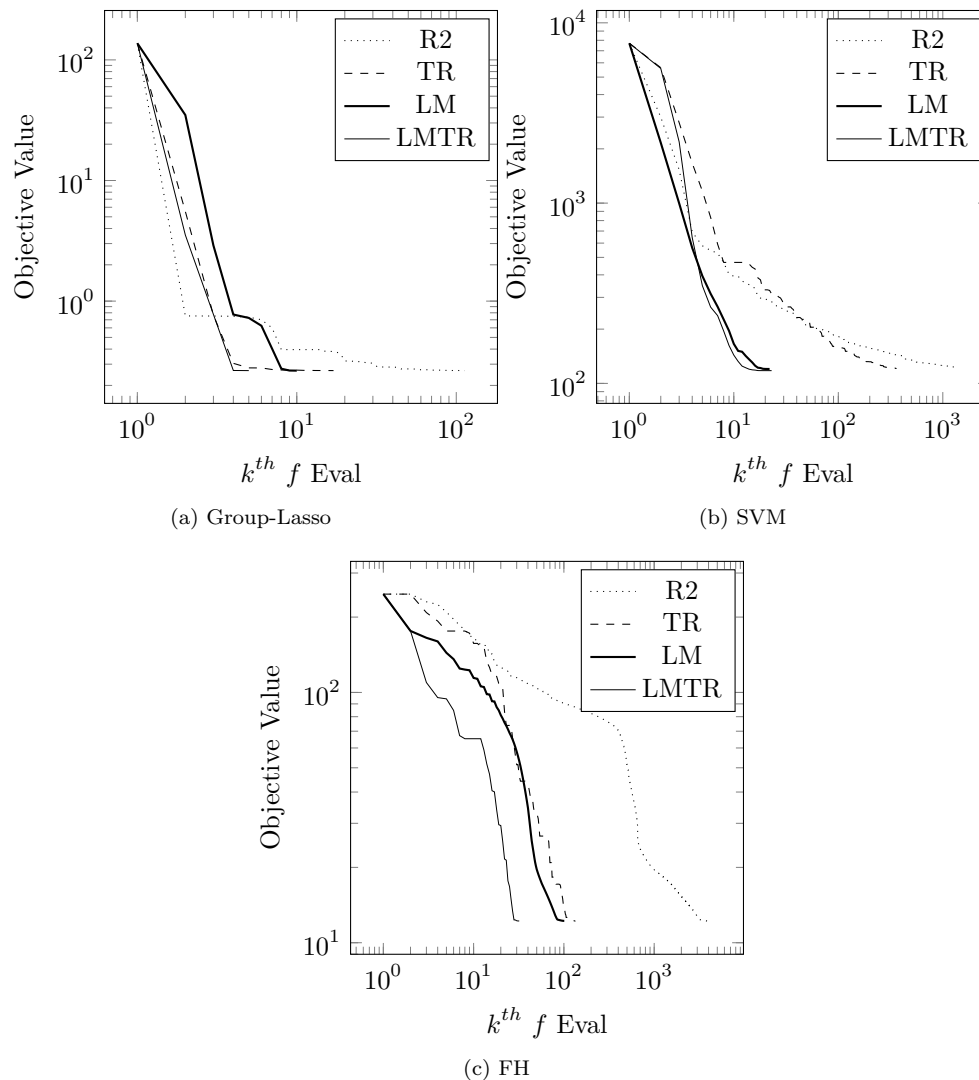


FIG. 6.4. Objective decrease per objective or residual evaluation.

- 689 [25] R. Rockafellar and R. Wets. *Variational Analysis*, volume 317. Springer Verlag, 1998.
- 690 [26] L. Stella, A. Themelis, P. Sotasakis, and P. Patrinos. **A simple and efficient algorithm for**
- 691 **nonlinear model predictive control**. In *2017 IEEE 56th Annual Conference on Decision and*
- 692 *Control (CDC)*, pages 1939–1944, 2017.
- 693 [27] A. Themelis, L. Stella, and P. Patrinos. **Forward-backward envelope for the sum of two nonconvex**
- 694 **functions: Further properties and nonmonotone linesearch algorithms**. *SIAM J. Optim.*, 28(3):
- 695 2274–2303, 2018.
- 696 [28] R. Tibshirani. **Regression shrinkage and selection via the lasso**. *J. Roy. Statist. Soc. Ser. B*, 58
- 697 (1):267–288, 1996.
- 698 [29] B. Van der Pol. **Lxxxviii. On “relaxation-oscillations”**. *The London, Edinburgh, and Dublin*
- 699 *Philosophical Magazine and Journal of Science*, 2(11):978–992, 1926.
- 700 [30] H. Zhu, G. Leus, and G. B. Giannakis. **Sparsity-cognizant total least-squares for perturbed**
- 701 **compressive sampling**. *IEEE T. Signal Proces.*, 59(5):2002–2016, 2011.